# FlexiAct: Towards Flexible Action Control in Heterogeneous Scenarios

SHIYI ZHANG<sup>\*</sup>, Tsinghua Shenzhen International Graduate School, Tsinghua University, China JUNHAO ZHUANG<sup>\*</sup>, Tsinghua Shenzhen International Graduate School, Tsinghua University, China ZHAOYANG ZHANG<sup>†</sup>, Tencent ARC Lab, China YING SHAN, Tencent ARC Lab, China

YANSONG TANG<sup>‡</sup>, Tsinghua Shenzhen International Graduate School, Tsinghua University, China



Fig. 1. Visualization for our FlexiAct results. Given a target image, FlexiAct transfers actions from a reference video to the target subject, achieving accurate motion adaptation and appearance consistency even in heterogeneous scenarios with varying spatial structures or cross-domain subjects.

Action customization involves generating videos where the subject performs actions dictated by input control signals. Current methods use pose-guided or global motion customization but are limited by strict constraints on spatial structure such as layout, skeleton, and viewpoint consistency, reducing adaptability across diverse subjects and scenarios. To overcome these limitations, we propose FlexiAct, which transfers actions from a reference video to an arbitrary target image. Unlike existing methods, FlexiAct allows for variations in layout, viewpoint, and skeletal structure between the subject of the reference video and the target image, while maintaining identity consistency. Achieving this requires precise action control, spatial structure adaptation, and consistency preservation. To this end, we introduce RefAdapter, a lightweight image-conditioned adapter that excels in spatial adaptation and consistency preservation, surpassing existing methods in balancing appearance consistency and structural flexibility. Additionally, based on our observations, the denoising process exhibits varying levels of attention to motion (low frequency) and appearance details (high frequency) at different timesteps. So we propose FAE (Frequency-aware Action Extraction), which, unlike existing methods that rely on separate spatial-temporal architectures, directly achieves action extraction during the denoising process. Experiments demonstrate that our method effectively transfers actions to subjects with diverse layouts, skeletons, and viewpoints. We release our code and model weights to support further research at FlexiAct.

 $\texttt{CCS Concepts:} \bullet \textbf{Computing methodologies} \to \textbf{Computer vision}.$ 

<sup>\*</sup>Equal contribution

<sup>&</sup>lt;sup>†</sup>Project lead: Zhaoyang Zhang (zhaoyangzhang@link.cuhk.edu.hk)

<sup>&</sup>lt;sup>‡</sup>Corresponding author: Yansong Tang (tang.yansong@sz.tsinghua.edu.cn)

Additional Key Words and Phrases: Artificial Intelligence Generated Content, Computer Vision, Video Customization

#### **ACM Reference Format:**

Shiyi Zhang\*, Junhao Zhuang, Zhaoyang Zhang, Ying Shan, and Yansong Tang. 2025. FlexiAct: Towards Flexible Action Control in Heterogeneous Scenarios. In Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers (SIGGRAPH Conference Papers '25), August 10–14, 2025, Vancouver, BC, Canada. ACM, New York, NY, USA, 11 pages. https://doi.org/10.1145/3721238.3730683

# 1 INTRODUCTION

Action transfer involves applying specific actions to a target subject and is widely used in films, games, and animation. However, it often requires substantial financial and human resources. For example, professional motion capture systems can cost tens of thousands of dollars and require skilled technicians. Similarly, creating a 30-second animation at 12 frames per second can take about 20 work days from six professional animators. These high costs pose significant challenges, limiting access for many potential creators.

In response to these limitations, significant efforts have been devoted to achieving motion control in video generation, which can be broadly categorized into two main approaches: (1) Predefined signals methods using signals like pose and depth maps, such as AnymateAnyone [Hu 2024] and StableAnimator [Tu et al. 2024b], and (2) Global motion methods like Motion Director [Zhao et al. 2023] and Motion Inversion [Wang et al. 2024b]. Despite advancements, these methods exhibit notable limitations. Predefined signal methods require strict alignment of spatial structures (e.g., shape, skeleton, viewpoint) between the target image and reference video, which is often not feasible in real-world scenarios. They also struggle with obtaining pose information for non-human subjects. Conversely, global motion methods typically generate motions with fixed layouts and cannot transfer motion across diverse subjects. Some approaches [Zhao et al. 2023] employ identity-specific Low-Rank Adaptations (LoRAs) [Hu et al. 2021] for animation, yet they encounter difficulties with appearance consistency and flexibility.

To this end, we introduce FlexiAct, an Image-to-Video (I2V) framework for flexible action customization in heterogeneous scenarios. Given a reference video and an arbitrary target image, our method transfers actions from the reference video to the target image without alignment in layout, shape, or viewpoint, preserving both action dynamics and appearance details. FlexiAct builds upon CogVideoX-I2V [Yang et al. 2024] with a two-stage training on two novel components: RefAdapter and Frequency-aware Action Extraction (FAE), addressing the following challenges: (1) **Spatial structure adaptation:** Adapting actions to target images with different poses, layouts, or viewpoints. (2) **Precise action extraction and control:** Accurately decoupling and replicating action from the reference video.

RefAdapter addresses the first challenge, which is an imageconditioned architecture generating videos given the input images. It combines the accuracy of I2V frameworks with the flexibility of conditional injection architectures like IP-Adapter [Ye et al. 2023a], ensuring appearance consistency between the video and the conditioning image while avoiding strict constraints on the first video frame. This enables FlexiAct to adapt reference motion to various spatial structures using arbitrary frames as image conditions. RefAdapter requires only low training costs, finetuning a small set of LoRA, avoiding the large parameter replication in ReferenceNet [Hu 2024] and ControlNet [Zhang et al. 2023].

For precise action control, we propose Frequency-aware Action Extraction (FAE). This method incorporates a set of learnable embeddings to capture entangled video information from the reference video during training. As illustrated in Figure 2, we observe that these embeddings dynamically adjust their attention to different frequency components across denoising timesteps. Specifically, they prioritize motion information (low-frequency features) in early timesteps and shift focus to appearance details (high-frequency features) in later timesteps. Leveraging this property, FAE performs action extraction directly during the denoising process by modulating attention weights at different timesteps, eliminating the need for separate spatial-temporal architectures.

To validate the effectiveness of FlexiAct, we establish a benchmark for heterogeneous scenarios. Experiments demonstrate FlexiAct's flexible and general action transfer capabilities. As shown in Figure 1, FlexiAct accurately transfers action from a reference video to subjects with varying layouts, viewpoints, shapes, and even domains, while maintaining appearance consistency.

In summary, our paper makes the following key contributions:

- We propose FlexiAct, a flexible action transfer method that **first** adapts reference actions to arbitrary subjects with diverse spatial structures while ensuring action and appearance consistency.
- We introduce RefAdapter, which achieves spatial structure adaptation and appearance consistency with a few trainable parameters.
- We propose Frequency-aware Action Extraction, which precisely extracts action and controls the video synthesis during sampling.
- Our extensive experiments demonstrate FlexiAct's capabilities across diverse scenarios, including various subjects and domains.

#### 2 RELATED WORK

#### 2.1 Global Motion Customization

Global Motion Customization focuses on transferring the overall motion dynamics from a reference video, such as camera movements, object trajectories, and actions, to generate videos with consistent global motion patterns [Jeong et al. 2024, 2023; Ling et al. 2024; Yatim et al. 2023; Zhao et al. 2023]. The challenge of this task lies in effectively extracting motion from the reference video. Recent work like Motion Director [Zhao et al. 2023] addresses this by adopting spatial-temporal LoRA[Hu et al. 2021] to decouple the appearance and the motion. Meanwhile, Diffusion Motion Transfer [Yatim et al. 2023] extracts motion via a handcrafted loss during inference. On the other hand, Video Motion Customization [Jeong et al. 2023] encodes motion directly into the text-to-video model. Motion Inversion [Wang et al. 2024b] introduces two types of embeddings to decouple the appearance and motion. However, most of these methods fail to adapt motion to specific subjects, as they primarily focus on generating videos that approximate the layout of the reference video. In contrast, we propose a framework that can handle the action transfer in heterogeneous scenarios. Furthermore, inspired by the insights into noise schedules from [Lin et al. 2024; Lu et al. 2024; Qiu et al. 2023; Si et al. 2024; Wu et al. 2023b] and our observations regarding different denoising timesteps, we propose the first denoising process-based action extraction framework.



Fig. 2. Attention maps between our frequency-aware embeddings and video tokens in the MMDIT at different denoising timesteps. Our embeddings focus on low-frequency motion information (e.g., motion regions) in early denoising stages and shift to high-frequency details in later stages.

#### 2.2 Predefined signal-based Action Customization

Action customization methods based on predefined signals, such as pose, depth, and edges, transfer motion from these signals to animate target images [Esser et al. 2023; He et al. 2022; Wang et al. 2023d]. They focus on how to precisely control subjects with identical spatial structures using predefined signals. Early approaches [Huang et al. 2021; Siarohin et al. 2019, 2021] primarily utilized GANs [Goodfellow et al. 2020] for reference animation. Recent advancements have shifted to diffusion models, with Disco [Wang et al. 2024a] pioneering this transition for image animation. Subsequent works, such as MagicAnimate [Xu et al. 2024] and AnimateAnyone [Hu 2024], employ ReferenceNet and pose nets to decouple pose and appearance modeling. Further innovations include Champ [Zhu et al. 2024], which integrates 3D SMPL signals for enhanced controllability, and Unianimate [Wang et al. 2024c], which incorporates Mamba [Dao and Gu 2024] into diffusion models for improved efficiency. Additionally, MimicMotion [Zhang et al. 2024a] introduces a regional loss to mitigate distortion, while ControlNeXt [Peng et al. 2024] replaces the computationally intensive ControlNet [Zhang et al. 2023] with a lightweight convolution-based pose net. However, these methods remain heavily reliant on predefined signals, limiting their effectiveness when the target image and reference video exhibit significant spatial discrepancies, such as variations in shape or pose. Moreover, they face challenges in non-human scenarios, where predefined motion signals are often unavailable or difficult to obtain. In contrast, our method does not rely on predefined signals with numerous constraints, enabling it to handle more general scenarios, such as transferring actions between subjects with different shapes, skeletons, viewpoints, and even across domains.

# 2.3 Customized Video Generation via Condition Injection

With the advancement of text-to-video models [Bar-Tal et al. 2024; Blattmann et al. 2023; Brooks et al. 2024; Chen et al. 2023b, 2024;

Esser et al. 2023; Guo et al. 2023; Ma et al. 2024; Wang et al. 2023c,b,a; Yang et al. 2024; Yuan et al. 2024; Zhang et al. 2024b; Zhou et al. 2024], customized video generation has emerged as a critical and highly active research topic. Among these methods, some focus on injecting control signals into the video generation process through condition injection, which can generally be categorized into two types: one based on cross-attention injection, such as IP-Adapter[Ye et al. 2023b], and the other on module duplication for layer-wise injection, such as ReferenceNet[Hu 2024]. Cross-attention approaches, though lightweight, often fail to ensure appearance consistency due to coarse-grained representations (e.g., CLIP image features [Radford et al. 2021]). Module duplication enables finer control but incurs high training costs from parameter replication. In contrast, based on the I2V model, our RefAdapter strikes a balance, achieving ReferenceNet-level fine-grained appearance control while requiring only a small number of training parameters. Additionally, RefAdapter can reduce the strict first-frame dependency of I2V models.

# 3 METHOD

# 3.1 Overview

As illustrated in Figure 3, FlexiAct builds upon CogVideoX-I2V with a two-stage training on two components: RefAdapter and Frequency-aware Action Extraction (FAE). RefAdapter facilitates action adaptation to subjects with varying spatial structures while maintaining appearance consistency. FAE dynamically adjusts attention weights to frequency-aware embeddings at different denoising timesteps, enabling effective action extraction. These two components are trained separately to ensure that action extraction does not interfere with RefAdapter's consistency preservation. Section 3.1 introduces FlexiAct's base model. Section 3.3 introduces RefAdapter and its training methodology. Section 3.4 details the training and inference pipeline of FlexiAct.



Fig. 3. **Overview of FlexiAct.** (1) The upper part illustrates RefAdapter's training, which conditions arbitrary frames to enable transitions across diverse spatial structures. (2) The lower part outlines FAE's training and inference, where attention weights of video tokens to the frequency-aware embedding are dynamically adjusted based on timesteps, facilitating action extraction.

# 3.2 Basis Image-to-Video Diffusion Model

We use CogVideoX-I2V [Yang et al. 2024] as our basis image-tovideo (I2V) model. CogVideoX-I2V is an MMDiT-based [Esser et al. 2024] video diffusion model that operates in a latent space. Given an image  $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$  and a textual prompt, CogVideoX-I2V generates a video  $\mathbf{V} \in \mathbb{R}^{T \times H \times W \times 3}$ . CogVideoX-I2V utilizes a 3D VAE to map condition images and videos into the latent space. For video inputs, the 3D VAE encoder ( $\epsilon$  in Figure 3) compresses both temporal and spatial dimension, producing a latent  $L_{video} \in \mathbb{R}^{\frac{T}{4} \times \frac{H}{8} \times \frac{W}{8} \times C}$ , where *C* denotes the channel number. For image inputs (*T* = 1), the encoder preserves the temporal dimension, yielding  $L_{image} \in \mathbb{R}^{1 \times \frac{H}{8} \times \frac{W}{8} \times C}$ , which is then zero-padded along the temporal dimension to match  $L_{video}$ 's shape  $(1 \rightarrow \frac{T}{4})$ . During inference, this padded  $L_{image}$  is concatenated with a random noise  $\mathcal{N} \in \mathbb{R}^{\frac{T}{4} \times \frac{H}{8} \times \frac{W}{8} \times C}$  along the channel dimension for subsequent denoising process. During training, the first frame of the ground truth video serves as the input image; its padded Limage is concatenated with the noisy Lvideo along the channel dimension to predict the added noise. This process is optimized with MSE loss between the added and predicted noise, consistent with classic diffusion models [Ho et al. 2020].

## 3.3 RefAdapter

The upper part of Figure 3 illustrates the RefAdapter training process. We note that directly using I2V for spatial structure adaptation is challenging because: (1) the action extraction process compromises the I2V model's consistency preservation, and (2) I2V is a strongly constrained image-conditioned framework. During training, I2V uses the first video frame as the condition image, ensuring video consistency but potentially hindering smooth action transfer if the initial spatial structure differs from the reference video.

To address this, we introduce a gap between the condition image and the initial spatial structure by using a randomly sampled frame from the unedited video as the condition image during training. Specifically, we propose **RefAdapter**, which includes LoRA injected into CogVideoX-I2V's MMDiT layers. RefAdapter is trained on 42,000 videos from [Ju et al. 2024] in a 40,000-step one-time training. The training process of RefAdapter mostly follows CogVideoX-I2V, with key distinctions: (1) The condition image is randomly selected from the entire untrimmed video instead of the first frame, maximizing spatial structure discrepancy. (2) We replace the first embedding along the temporal dimension of  $L_{video}$  with  $L_{image}$ , enabling the model to use the first embedding as a reference for guiding video generation, rather than constraining it as the video's starting point. Without this replacement, the generated video's initial state would remain constrained to match the condition image.

#### 3.4 Frequency-aware Action Extraction

To extract action information from a reference video, a straightforward approach involves training motion embeddings to fit motion information, akin to Genie [Bruce et al. 2024] or textual inversion [Gal et al. 2022]. However, our initial attempts with this method yield suboptimal results. Inspired by [Lin et al. 2024; Lu et al. 2024; Qiu et al. 2023; Si et al. 2024; Wu et al. 2023b], we delve into the relationship between the embeddings and action information during the denoising process. By examining the attention maps between motion embeddings and video tokens across different denoising timesteps, as visualized in Figure 2, we observe that our embeddings predominantly focus on low-frequency action information in the early stages, gradually shifting their attention to high-frequency details in the later stages. Leveraging this insight, we propose the Frequency-aware Action Extraction. FlexiAct: Towards Flexible Action Control in Heterogeneous Scenarios

Specifically, the lower part of Figure 3 illustrates the training and inference process of FAE. We train **Frequency-aware Embedding** for individual reference videos, which includes learnable parameters concatenated to MMDiT layers' inputs. This training differs from CogVideoX-I2V by applying random cropping on the input video to prevent the Frequency-aware Embedding from focusing on the reference video's layout. RefAdapter is not loaded during this training to protect its conditioning ability. After training, the Frequency-aware Embedding captures both motion and appearance from the reference video.

During inference, FAE extracts action information and adapts it to the target image. As shown in Figure 2, attention maps between frequency-aware embeddings and video tokens reveal that at larger timesteps (e.g., step=800), the embeddings focus on motion information (low frequency), with high attention on the moving parts of the subject. At intermediate timesteps (e.g., step=500), the focus shifts to fine-grained details of the subject. At later timesteps (e.g., step=200), attention is distributed across the entire image, indicating a focus on global details like the background.

Based on this observation, during inference, we increase the attention weight of video tokens on the frequency-aware embeddings at larger timesteps while maintaining the original weights at other timesteps. This enhances the generated video's ability to perceive and replicate the motion of the reference video. The reweighting strategy of FAE can be formulated as:

$$W_{bias} = \begin{cases} \alpha, & t_l \le t \le T \\ \frac{\alpha}{2} \left[ \cos\left(\frac{\pi}{t_h - t_l}(x - t_l)\right) + 1 \right], & t_h \le t < t_l \\ 0, & 0 \le t < t_h \end{cases}$$
(1)

where  $W_{\text{bias}}$  denotes the bias value applied to the original attention weight. *t* denotes denoising timestep. The parameter  $\alpha$  controls the strength of the bias, while  $t_l$  and  $t_h$  represent the low-frequency and high-frequency timesteps, respectively. A transition function is employed between these timesteps to smoothly vary  $W_{\text{bias}}$  from low-frequency to high-frequency timesteps. During inference, the attention weight of video tokens to the frequency-aware embedding is dynamically adjusted as  $W_{\text{attn}} = W_{ori} + W_{\text{bias}}$  in all DiT layers.

Experimental results demonstrate that the attention reweighting strategy improves the generated video's ability to reproduce the reference action. In practice, we set  $\alpha = 1$ ,  $t_h = 700$ , and  $t_l = 800$ , and demonstrate the impact of the transition function on the generation results in Figure 11. Without the transition function, changing the bias at t = 700 would cause the model to learn appearance information from the reference video (e.g., the clothing in Figure 11) while altering the bias at t = 800 would result in inaccurate motion. This indicates that a transition process between t = 700 and t = 800 is necessary to achieve a balance between appearance and motion.

#### 3.5 Training and Inference Pipeline of FlexiAct

As shown in Figure 3, we first train RefAdapter on a broader dataset (upper part). Subsequently, we train the frequency-aware embedding based on the individual reference video. RefAdapter does not participate in this training process. During the Inference phase, the RefAdapter is loaded, and an arbitrary target image is provided. FAE SIGGRAPH Conference Papers '25, August 10-14, 2025, Vancouver, BC, Canada



Target ImagePose-based MethodGlobal Method (I2V)

Fig. 4. Results of transferring "turning" action to the target image using the pose-based method and the animation version of the global motion method. dynamically adjusts the generated video's attention to Frequency-aware Embedding according to denoising timesteps, transferring actions from the reference video to the target image.

## 4 EXPERIMENT

#### 4.1 Implementation Details

**Evaluation Dataset.** We conduct experiments on a evaluation dataset of 250 video-image pairs, featuring 25 distinct action categories. Each action is transferred to 10 different target images, covering a wide range of human motions (e.g., yoga, fitness exercises) and animal motions (e.g., jumping, running). The target images include real humans, animals, animated, and game characters. This diversity ensures our dataset encompasses a broad spectrum of scenarios, allowing for a comprehensive evaluation of our method's generalization capabilities.

**Comparison Methods.** Existing methods for action transfer include those based on predefined signals and global motion. Predefined signal methods are ineffective for non-human entities or subjects with significant skeletal differences, as shown in Figure 4. Therefore, we use the recent global motion transfer method, MotionDirector [Tu et al. 2024a], as our baseline. For a fair comparison, we reimplement MotionDirector on the stronger CogVideoX-I2V backbone (referred to as MD-I2V) with identical training settings to our methods. Additionally, we implement a base model that learns actions directly through standard learnable action embeddings, without using RefAdapter and FAE (referred to as BaseModel).

**Training Details.** For RefAdapter's training, we conduct a one-time 40,000-step training on Miradata [Ju et al. 2024] with a learning rate of 1e-5 and a batch size of 8 with the AdamW optimizer. RefAdapter introduces 66M parameters, constituting 5% of CogVideoX-I2V's total parameters. Frequency-aware embeddings require 1,500 to 3,000 training steps on each reference video, depending on action complexity. In comparison, the Motion Director needs 3,000 steps for temporal LoRA and 300 for spatial LoRA.

#### 4.2 Quantitative Evaluation

Automatic Evaluations. Following [Jeong et al. 2023; Wang et al. 2024b; Yatim et al. 2023], we employ *Text Similarity, Motion Fidelity,* and *Temporal Consistency* to evaluate the semantic accuracy of the



Fig. 5. Qualitative comparison of action transfer from reference video (Ref Video) to target images with varying spatial structures. Red boxes highlight regions where the appearance deviates from the target image. Our method demonstrates superior performance in maintaining appearance consistency with the target image and motion fidelity to the reference video compared to other approaches.

Method	Automatic Evaluations				Human Evaluations		
	Text Similarity <sup>↑</sup>	Motion Fidelity ↑	Temporal Consistency <sup>↑</sup>	Appearance Consistency ↑		Motion Consistency	Appearance Consistency
MD-I2V [Zhao et al. 2023] Base Model	0.2446 0.2541	0.3496 0.3562	0.9276 0.9283	0.8963 0.8951	v.s. Base Model -	47.2 v.s. 52.8 -	53.1 v.s. 46.9 -
w/o FAE w/o RefAdapter	0.2675 0.2640	0.3614 0.3856	0.9255 0.9217	0.9134 0.9021	v.s. Base Model v.s. Base Model	59.7 v.s. 40.3 68.6 v.s. 31.4	76.4 v.s. 23.6 52.2 v.s. 47.8
Ours	0.2732	0.4103	0.9342	0.9162	v.s. Base Model	79.5 v.s. 20.5	78.3 v.s. 21.7

Table 1. Quantitative comparisons and human evaluations. We train an I2V version of Motion Director (MD-I2V) based on CogVideoX. The Base Model trains a set of learnable embeddings without incorporating both RefAdapter and FAE. " $p_1$  v.s.  $p_2$ " means  $p_1$ % results of the first method are preferred.

generated videos, the degree of motion alignment with the reference videos, and the temporal coherence, respectively. Furthermore, we introduce *Appearance Consistency* to assess the consistency in appearance between the generated videos and the target image. Below, we provide a brief overview of these metrics.

*Text Similarity.* It is calculated with CLIP [Radford et al. 2021] frame-to-text similarity, reflecting the semantic alignment degree between the output video and the prompt.

*Motion Fidelity.* Introduced by [Yatim et al. 2023], it utilizes tracklets computed by a tracking model [Karaev et al. 2023], measuring the similarity between the motion trajectories in unaligned videos.

*Temporal Consistency.* It measures the smoothness and coherence of a video sequence [Chen et al. 2023a; Jeong et al. 2023; Wu et al. 2023a; Zhao et al. 2023], quantified by the average similarity between the CLIP image features of all frame pairs within the output video.

Appearance Consistency. It reflects the appearance consistency between the output video and the target image, calculated as the average CLIP similarity between the first frame and the remaining frames of the output video. **Human Evaluation.** Following [Zhao et al. 2023], we conduct a human evaluation with 5 raters who assessed each generated video for appearance consistency with the target image and motion consistency with the reference video. Each participant compares 50 randomly selected video pairs, each containing one video generated by a random method and one by the Base Model. Following [Zhao et al. 2023], all methods are compared against the Base Model, serving as a solid baseline due to its comparable performance to MD-I2V. In Table 1, " $p_1$  v.s.  $p_2$ " indicates that  $p_1$ % of the first method's results are prefer over  $p_2$ % of the second method's results.

**Results.** As shown in Table 1, our method significantly outperforms baseline approaches in both motion fidelity and appearance consistency. This underscores the challenges of action transfer in heterogeneous scenarios and demonstrates our approach's effectiveness in balancing action accuracy with appearance consistency. Notably, our Motion Fidelity scores are generally lower than those in global motion tasks, as they are affected by layout consistency, whereas global motion tasks involve transferring motion to scenarios with identical layouts, making them not directly comparable to action transfer in heterogeneous scenarios.



Fig. 6. Qualitative results of ablation study. We ablate Frequency-aware Action Extraction (FAE) and RefAdapter, comparing the action transfer results from reference videos (Ref Video) to different subjects. Ablating FAE reduces action accuracy, demonstrating its effectiveness in action extraction. Ablating RefAdapter degrades both appearance consistency and action precision, proving its capability in spatial structure adaptation for cross-subject action transfer.

#### 4.3 Qualitative Evaluation

Figure 5 shows a qualitative comparison with the baseline method. MD-I2V struggles to replicate the reference video's motion accurately. In the first example, the man fails to stand up after squatting, and his arm movements do not match the original. In the second, he does not lift his leg as in the reference, and one eye closes in later frames. The Base Model also suffers from motion accuracy and appearance consistency issues. In the first example, the man puts on clothes, deviating from the original image, and his final motion differs from the reference. In the second, his leg lift is exaggerated, and clothing-like folds appear in the final frames. In contrast, our method excels in both motion accuracy and appearance consistency.

#### 4.4 Ablation Study

Table 1 and Figure 6 present our Ablation Study results. Quantitative data show that removing FAE significantly decreases Motion Fidelity, highlighting its role in enhancing motion generation quality. This is corroborated by qualitative results, where two distinct actions transferred to different characters exhibit inconsistencies without FAE. For example, in a stretching motion, the character merely raises their hand without proper bending or stretching, deviating from the reference video. Similar mismatches in the second example further emphasize FAE's importance for motion consistency.

We also examined the impact of removing RefAdapter. Quantitative results indicate noticeable declines in both appearance consistency and motion fidelity, as RefAdapter ensures adaptability to varying spatial structures. Without it, the model struggles to adapt motion to target images with different spatial layouts, weakening appearance consistency. Qualitative results in Figure 6 support this: in the first example, discrepancies in the character's face and clothing are resolved with RefAdapter. In the second example, without RefAdapter, the output video fails to extend arms fully, maintaining them bent, and shows noticeable differences in facial details, underscoring RefAdapter's role in maintaining both motion and appearance consistency.

# 5 DISCUSSION

In this paper, we tackle the action transfer in heterogeneous scenarios, where the main difficulty is achieving precise action transfer for subjects with different spatial structures while maintaining appearance consistency. We introduce FlexiAct, a flexible and versatile approach that surpasses existing methods. Our RefAdapter adapts to various spatial structures and ensures appearance consistency, while Frequency-aware Action Extraction allows for precisely extracting action during the denoising process. Extensive experiments show that FlexiAct effectively balances action accuracy and appearance consistency across diverse spatial structures and domains.

Despite achieving precise action and appearance consistency, like [Jeong et al. 2023; Wang et al. 2024b; Yatim et al. 2023; Zhao et al. 2023], our method requires optimization for each reference video. Developing feed-forward motion transfer methods for heterogeneous scenarios is a key direction for future work.

Acknowledgments. This work was supported by Guangdong Natural Science Funds for Distinguished Young Scholar (No. 2025B1515020012).



Fig. 7. FlexiAct can transfer actions to diverse subjects while maintaining both appearance consistency with the target subject and action consistency with the reference video.



Fig. 8. Examples of human action transfer using FlexiAct.





Fig. 11. Ablation of bias transition.

SIGGRAPH Conference Papers '25, August 10-14, 2025, Vancouver, BC, Canada

Shiyi Zhang\*, Junhao Zhuang, Zhaoyang Zhang, Ying Shan, and Yansong Tang

#### REFERENCES

- Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat, Junhwa Hur, Guanghui Liu, Amit Raj, et al. 2024. Lumiere: A space-time diffusion model for video generation. arXiv preprint arXiv:2401.12945 (2024).
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. 2023. Align your Latents: High-Resolution Video Synthesis with Latent Diffusion Models. In Proc. CVPR.
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. 2024. Video generation models as world simulators. (2024). https: //openai.com/research/video-generation-models-as-world-simulators
- Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. 2024. Genie: Generative interactive environments. In Forty-first International Conference on Machine Learning.
- Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. 2023b. VideoCrafter1: Open Diffusion Models for High-Quality Video Generation. arXiv preprint arXiv:2310.19512 (2023).
- Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying Shan. 2024. Videocrafter2: Overcoming data limitations for high-quality video diffusion models. In Proc. CVPR. 7310–7320.
- Weifeng Chen, Jie Wu, Pan Xie, Hefeng Wu, Jiashi Li, Xin Xia, Xuefeng Xiao, and Liang Lin. 2023a. Control-A-Video: Controllable Text-to-Video Generation with Diffusion Models. arXiv preprint arXiv:2305.13840 (2023).
- Tri Dao and Albert Gu. 2024. Transformers are SSMs: Generalized Models and Efficient Algorithms Through Structured State Space Duality. In *ICML*.
- Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. 2023. Structure and content-guided video synthesis with diffusion models. In Proc. ICCV. 7346–7356.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. 2024. Scaling rectified flow transformers for high-resolution image synthesis. In Forty-first international conference on machine learning.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing text-toimage generation using textual inversion. arXiv preprint arXiv:2208.01618 (2022).
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. Commun. ACM (2020).
- Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. 2023. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. arXiv preprint arXiv:2307.04725 (2023).
- Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. 2022. Latent Video Diffusion Models for High-Fidelity Long Video Generation. (2022). arXiv:2211.13221 [cs.CV]
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. Proc. NeurIPS 33 (2020), 6840–6851.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021).
- Li Hu. 2024. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. In Proc. CVPR. 8153–8163.
- Zhichao Huang, Xintong Han, Jia Xu, and Tong Zhang. 2021. Few-shot human motion transfer by personalized geometry and texture modeling. In CVPR.
- Hyeonho Jeong, Jinho Chang, Geon Yeong Park, and Jong Chul Ye. 2024. DreamMotion: Space-Time Self-Similarity Score Distillation for Zero-Shot Video Editing. arXiv preprint arXiv:2403.12002 (2024).
- Hyeonho Jeong, Geon Yeong Park, and Jong Chul Ye. 2023. VMC: Video Motion Customization using Temporal Attention Adaption for Text-to-Video Diffusion Models. arXiv preprint arXiv:2312.00845 (2023).
- Xuan Ju, Yiming Gao, Zhaoyang Zhang, Ziyang Yuan, Xintao Wang, Ailing Zeng, Yu Xiong, Qiang Xu, and Ying Shan. 2024. Miradata: A large-scale video dataset with long durations and structured captions. arXiv preprint arXiv:2407.06358 (2024).
- Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. 2023. Cotracker: It is better to track together. arXiv preprint arXiv:2307.07635 (2023).
- Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. 2024. Common diffusion noise schedules and sample steps are flawed. In Proceedings of the IEEE/CVF winter conference on applications of computer vision. 5404–5411.
- Pengyang Ling, Jiazi Bu, Pan Zhang, Xiaoyi Dong, Yuhang Zang, Tong Wu, Huaian Chen, Jiaqi Wang, and Yi Jin. 2024. MotionClone: Training-Free Motion Cloning for Controllable Video Generation. arXiv preprint arXiv:2406.05338 (2024).
- Yu Lu, Yuanzhi Liang, Linchao Zhu, and Yi Yang. 2024. Freelong: Training-free long video generation with spectralblend temporal attention. arXiv preprint arXiv:2407.19918 (2024).

- Xin Ma, Yaohui Wang, Gengyun Jia, Xinyuan Chen, Ziwei Liu, Yuan-Fang Li, Cunjian Chen, and Yu Qiao. 2024. Latte: Latent diffusion transformer for video generation. arXiv preprint arXiv:2401.03048 (2024).
- Bohao Peng, Jian Wang, Yuechen Zhang, Wenbo Li, Ming-Chang Yang, and Jiaya Jia. 2024. ControlNeXt: Powerful and Efficient Control for Image and Video Generation. arXiv preprint arXiv:2408.06070 (2024).
- Haonan Qiu, Menghan Xia, Yong Zhang, Yingqing He, Xintao Wang, Ying Shan, and Ziwei Liu. 2023. Freenoise: Tuning-free longer video diffusion via noise rescheduling. arXiv preprint arXiv:2310.15169 (2023).
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *Proc. ICML*. PMLR, 8748–8763.
- Chenyang Si, Ziqi Huang, Yuming Jiang, and Ziwei Liu. 2024. Freeu: Free lunch in diffusion u-net. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 4733–4743.
- Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019. First Order Motion Model for Image Animation. In *NeurIPS*.
- Aliaksandr Siarohin, Oliver J Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. 2021. Motion representations for articulated animation. In CVPR.
- Shuyuan Tu, Qi Dai, Zhi-Qi Cheng, Han Hu, Xintong Han, Zuxuan Wu, and Yu-Gang Jiang. 2024a. Motioneditor: Editing video motion via content-aware diffusion. In *CVPR*.
- Shuyuan Tu, Zhen Xing, Xintong Han, Zhi-Qi Cheng, Qi Dai, Chong Luo, and Zuxuan Wu. 2024b. StableAnimator: High-Quality Identity-Preserving Human Image Animation. arXiv preprint arXiv:2411.17697 (2024).
- Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. 2023c. Modelscope text-to-video technical report. arXiv preprint arXiv:2308.06571 (2023).
- Luozhou Wang, Ziyang Mai, Guibao Shen, Yixun Liang, Xin Tao, Pengfei Wan, Di Zhang, Yijun Li, and Yingcong Chen. 2024b. Motion inversion for video customization. arXiv preprint arXiv:2403.20193 (2024).
- Tan Wang, Linjie Li, Kevin Lin, Yuanhao Zhai, Chung-Ching Lin, Zhengyuan Yang, Hanwang Zhang, Zicheng Liu, and Lijuan Wang. 2024a. Disco: Disentangled control for realistic human dance generation. In CVPR.
- Wenjing Wang, Huan Yang, Zixi Tuo, Huiguo He, Junchen Zhu, Jianlong Fu, and Jiaying Liu. 2023b. VideoFactory: Swap Attention in Spatiotemporal Diffusions for Text-to-Video Generation. arXiv preprint arXiv:2305.10874 (2023).
- Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen, Deli Zhao, and Jingren Zhou. 2023d. VideoComposer: Compositional Video Synthesis with Motion Controllability. arXiv preprint arXiv:2306.02018 (2023).
- Xiang Wang, Shiwei Zhang, Changxin Gao, Jiayu Wang, Xiaoqiang Zhou, Yingya Zhang, Luxin Yan, and Nong Sang. 2024c. UniAnimate: Taming Unified Video Diffusion Models for Consistent Human Image Animation. arXiv preprint arXiv:2406.01188 (2024).
- Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. 2023a. LAVIE: High-Quality Video Generation with Cascaded Latent Diffusion Models. arXiv preprint arXiv:2309.15103 (2023).
- Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. 2023a. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 7623–7633.
- Tianxing Wu, Chenyang Si, Yuming Jiang, Ziqi Huang, and Ziwei Liu. 2023b. Freeinit: Bridging initialization gap in video diffusion models. arXiv preprint arXiv:2312.07537 (2023).
- Zhongcong Xu, Jianfeng Zhang, Jun Hao Liew, Hanshu Yan, Jia-Wei Liu, Chenxu Zhang, Jiashi Feng, and Mike Zheng Shou. 2024. Magicanimate: Temporally consistent human image animation using diffusion model. In CVPR.
- Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. 2024. CogVideoX: Text-to-Video Diffusion Models with An Expert Transformer. arXiv preprint arXiv:2408.06072 (2024).
- Danah Yatim, Rafail Fridman, Omer Bar Tal, Yoni Kasten, and Tali Dekel. 2023. Space-Time Diffusion Features for Zero-Shot Text-Driven Motion Transfer. arXiv preprint arXiv:2311.17009 (2023).
- Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. 2023a. IP-Adapter: Text Compatible Image Prompt Adapter for Text-to-Image Diffusion Models. arXiv preprint arxiv:2308.06721 (2023).
- Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. 2023b. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721* (2023).
- Hangjie Yuan, Shiwei Zhang, Xiang Wang, Yujie Wei, Tao Feng, Yining Pan, Yingya Zhang, Ziwei Liu, Samuel Albanie, and Dong Ni. 2024. InstructVideo: instructing video diffusion models with human feedback. In Proc. CVPR. 6463–6474.

FlexiAct: Towards Flexible Action Control in Heterogeneous Scenarios

- David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. 2024b. Show-1: Marrying pixel and latent diffusion models for text-to-video generation. *Int. J. Comput. Vis.* (2024), 1–15.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. In *Proc. ICCV*. 3836–3847.
  Yuang Zhang, Jiaxi Gu, Li-Wen Wang, Han Wang, Junqi Cheng, Yuefeng Zhu, and
- Fangy uan Zou. 2024. Alimicmotion: High-quality human motion video generation with confidence-aware pose guidance. *arXiv preprint arXiv:2406.19680* (2024).
- Rui Zhao, Yuchao Gu, Jay Zhangjie Wu, David Junhao Zhang, Jiawei Liu, Weijia Wu, Jussi Keppo, and Mike Zheng Shou. 2023. MotionDirector: Motion Customization of Text-to-Video Diffusion Models. arXiv:2310.08465 [cs.CV]
- Yuan Zhou, Qiuyue Wang, Yuxuan Cai, and Huan Yang. 2024. Allegro: Open the Black Box of Commercial-Level Video Generation Model. arXiv preprint arXiv:2410.15458 (2024).
- Shenhao Zhu, Junming Leo Chen, Zuozhuo Dai, Yinghui Xu, Xun Cao, Yao Yao, Hao Zhu, and Siyu Zhu. 2024. Champ: Controllable and Consistent Human Image Animation with 3D Parametric Guidance. In *EECV*.