

AnimaX: Animating the Inanimate in 3D with Joint Video-Pose Diffusion Models

ZEHUAN HUANG, Beihang University, China

HAORAN FENG, Tsinghua University, China

YANTIAN SUN, The University of Hong Kong, China

YUANCHEN GUO*, VAST, China

YANPEI CAO†, VAST, China

LU SHENG†, Beihang University, China



Fig. 1. Diverse articulated 3D models animated using *AnimaX*. The created animation, spanning various categories including humanoids, animals, and fictional models, demonstrates the versatility of our method. Selected models are visualized with keyframes of their predicted animations on the conveyor belts.

We present AnimaX, a feed-forward 3D animation framework that bridges the motion priors of video diffusion models with the controllable structure of skeleton-based animation. Traditional motion synthesis methods are either restricted to fixed skeletal topologies or require costly optimization in high-dimensional deformation spaces. In contrast, AnimaX effectively transfers video-based motion knowledge to the 3D domain, supporting diverse articulated meshes with arbitrary skeletons. Our method represents

3D motion as multi-view, multi-frame 2D pose maps, and enables joint video-pose diffusion conditioned on template renderings and a textual motion prompt. We introduce shared positional encodings and modality-aware embeddings to ensure spatial-temporal alignment between video and pose sequences, effectively transferring video priors to motion generation task. The resulting multi-view pose sequences are triangulated into 3D joint positions and converted into mesh animation via inverse kinematics. Trained on a newly curated dataset of 160,000 rigged sequences, AnimaX achieves state-of-the-art results on VBFench in generalization, motion fidelity, and efficiency, offering a scalable solution for category-agnostic 3D animation. Project page: <https://anima-x.github.io/>.

*Project leader.

†Corresponding author.

Authors' addresses: Zehuan Huang, Beihang University, China, huangzehuan@buaa.edu.cn; Haoran Feng, Tsinghua University, China, fenghr24@mails.tsinghua.edu.cn; Yangtian Sun, The University of Hong Kong, China, sunyangtian98@gmail.com; Yuanchen Guo, VAST, China, imbennguo@gmail.com; Yanpei Cao, VAST, China, caoyanpei@gmail.com; Lu Sheng, Beihang University, China, lsheng@buaa.edu.cn.

Additional Key Words and Phrases: 3D animation generation, generative model, 4D generation

Table 1. Comparison of *AnimaX* with existing generative work in 3D animation. Not all related methods are listed, but other approaches are generally similar to those included in the table.

Method	Multi Categories	Skeleton-Based	Generative Prior	Output Format	Cost Time
MotionDiffuse [Zhang et al. 2024a]	✗	✓	Motion Diffusion	Pose Sequence	20 s
MDM [Tevet et al. 2023]	✗	✓	Motion Diffusion	Pose Sequence	25 s
ATU [Millán et al. 2025]	✗	✓	Video Diffusion	Pose and Mesh Sequence	1.5 hours
Diffusion4D [Liang et al. 2024b]	✓	✗	MV Video Diffusion	NeRF	8 min
Animate3D [Jiang et al. 2024a]	✓	✗	MV Video Diffusion	GS	45 min
MotionDreamer [Uzolas et al. 2024]	✓	✗	Video Diffusion	Mesh Sequence	20 min
AKD [Li et al. 2025]	✓	✓	Score Distillation Sampling	GS	25 hours
AnimaX (Ours)	✓	✓	MV Video-Pose Diffusion	Pose and Mesh Sequence	6 min

1 INTRODUCTION

In traditional computer graphics, skeleton-based character animation typically involve binding a skeleton to a mesh and defining keyframes for motions. While this established technique affords high realism and fine-grained control over the resulting motion, it necessitates substantial manual effort from highly skilled artists, which is both time-consuming and expensive.

Recent advances in generative models [Ho et al. 2020; Peebles and Xie 2023; Radford et al. 2019] offer promising avenues for automating the character animation pipeline. Several studies have trained motion diffusion [Chen et al. 2023a; Tevet et al. 2023; Zhang et al. 2024a] or auto-regressive models [Zhang et al. 2023, 2024d] on collected motion capture data [Guo et al. 2022; Mahmood et al. 2019], enabling text-to-motion generation. However, these models can only be trained on datasets with a pre-defined, fixed skeleton system (*i.e.*, definition of joints with their connectivity), or rely on parametric 3D human models [Loper et al. 2015] for reconstruction. These methods primarily support motion synthesis for a single skeletal topology, such as humanoid motion, limiting their ability to generate animations for more diverse character categories.

Another series of work [Bahmani et al. 2024b; Jiang et al. 2024a; Liang et al. 2024b; Ren et al. 2023; Uzolas et al. 2024] explores 3D animation by leveraging advanced video generation models [Blattmann et al. 2023a; Guo et al. 2024; Wang et al. 2025; Yang et al. 2024b], distilling their learned generalized dynamic motion into consistent 4D sequences. As summarized in Tab. 1, these methods commonly leverage multi-view video diffusion models [Jiang et al. 2024a; Liang et al. 2024b] to guide the optimization of neural deformation fields [Pumarola et al. 2021; Wu et al. 2024], which predict displacements at each location within a 3D volume to deform a 3D shape. The resulting animation is a temporal sequence of these deformed shapes. While flexible, these approaches do not involve low-level skeleton-based motion representation, and instead introduce a large number of degrees of freedom (DoFs), making optimization challenging and often leading to in-consistent shapes and suboptimal quality. More recently, AKD [Li et al. 2025] distills articulated motion sequences from a pre-trained video diffusion model using Score Distillation Sampling (SDS) [Poole et al. 2022], simplifying the optimization by limiting the number of DoFs to that of a few joints. But it requires expensive optimization that takes even 25 hours.

We focus on efficiently animating articulated 3D meshes with arbitrary skeletal structures in a feed-forward manner, combining the diverse motion knowledge of video generation models [Kong et al. 2024; Wang et al. 2025; Yang et al. 2024b] with the low-DoF control of skeleton-based animation. Given an articulated mesh and a textual description, our goal is to generate 3D motion sequences that animate the mesh. The challenge lies in encoding and decoding sparse 3D poses for diverse skeletal topologies. While representing motion as graphs is straightforward, it hinders leveraging motion priors within video generation models, which are fundamental for category-agnostic and motion-diverse 3D animation.

Our system, *AnimaX*, addresses this by representing 3D motion as multi-view, multi-frame 2D pose maps, and adapting video diffusion models to generate such motion sequences in a feed-forward way. A straightforward baseline is to fine-tune a video diffusion model to generate pose sequences alone. However, the sparsity of pose representation and modality gap between RGB and pose maps make fine-tuning challenging, and disrupts the learned spatial-temporal priors, leading to distorted or nearly static pose outputs (Fig. 6). To better preserve and transfer video-based motion priors, we reveal the spatial alignment between video frames and pose frames at each timestep, and introduce a joint video-pose diffusion model that simultaneously predicts RGB videos and pose sequences. Crucially, we apply shared positional encoding across corresponding tokens in both modalities, ensuring coherence between video and pose streams. We find that this joint generation strategy—combined with shared positional encoding—allows the spatial-temporal priors learned from videos to be effectively grounded into the pose sequence generation process, resulting in more expressive motion outputs. Finally, we reconstruct the 3D animation by triangulating joint positions from multi-view poses and applying inverse kinematics to compute the joint angles.

We trained our multi-view video-pose diffusion models on a newly curated dataset of nearly 160,000 rigged 3D animation sequences, encompassing diverse categories such as humanoids, animals, and furniture. Evaluation on VBench [Huang et al. 2024b] demonstrates that *AnimaX* outperforms prior work in terms of generalizability across mesh categories, motion richness and naturalness, and efficiency. Our contributions are summarized as follows:

- We introduce *AnimaX*, an efficient feed-forward framework for animating diverse 3D articulated meshes with arbitrary

skeletal structures. AnimaX uniquely bridges rich motion priors from video diffusion models with the controllability of skeleton-based animation, overcoming key limitations of prior fixed-topology or category-specific approaches.

- We represent 3D motion as multi-view, multi-frame pose maps and design a joint multi-view video-pose diffusion model that simultaneously generates videos and corresponding 2D pose map sequences. This model incorporates novel shared positional encodings and modality-specific embeddings to ensure robust spatio-temporal alignment between video and pose, enabling a highly effective transfer of motion knowledge from video models to 3D animation task.
- We contribute a new, large-scale dataset of approximately 160,000 rigged 3D animation sequences. This dataset, encompassing diverse categories (e.g., humanoids, animals, articulated objects), is crucial for training generalizable, category-agnostic animation models like AnimaX and will serve as a valuable resource for future research.

2 RELATED WORK

Generative Models for 3D Animation. Recent advances in generative models [Esser et al. 2021; Ho et al. 2020; Peebles and Xie 2023; Radford et al. 2019] have spurred rapid progress in 3D animation. A significant body of work focuses on category-specific motion generation, such as text-driven human motion synthesis [Ahuja and Morency 2019; Azadi et al. 2023; Chen et al. 2023a; Jiang et al. 2023; Li et al. 2024; Liang et al. 2024a; Petrovich et al. 2022; Pi et al. 2024; Tevet et al. 2022, 2023; Zhang et al. 2023, 2024a,d]. For example, MDM [Tevet et al. 2023] successfully employ diffusion models [Ho et al. 2020] for this task, with subsequent work [Chen et al. 2023a] exploring latent diffusion models [Rombach et al. 2022]. Others [Jiang et al. 2023; Zhang et al. 2024d] leverage large language models [Brown et al. 2020; Radford et al. 2019] within the motion domain to support diverse motion-related tasks. However, these models only adapt to a pre-defined, fixed skeletal structures, or rely on parametric 3D human models [Bogo et al. 2016; Loper et al. 2015], hindering the generation of diverse character animations.

Another series of research [Bahmani et al. 2024a,b; Chen et al. 2025; Jiang et al. 2024b; Liang et al. 2024b; Ling et al. 2024; Liu et al. 2025; Ren et al. 2023; Shi et al. 2025; Sun et al. 2024; Uzolas et al. 2024; Wu et al. 2025; Yang et al. 2024a; Zeng et al. 2024; Zhang et al. 2024c; Zhao et al. 2024; Zhu et al. 2025] explores 3D animation [Azadi et al. 2023; Zhang et al. 2024b, 2021] by leveraging pre-trained image [Labs 2024; Podell et al. 2023; Rombach et al. 2022], video [Bao et al. 2024; Blattmann et al. 2023a; Guo et al. 2024; Ho et al. 2022; Kong et al. 2024; Singer et al. 2022; Wang et al. 2025; Yang et al. 2024b], multi-view image [Gao et al. 2024; Huang et al. 2024a,c; Liu et al. 2023; Shi et al. 2023; Wen et al. 2024; Zuo et al. 2024], multi-view video [Bai et al. 2024; Jiang et al. 2024a; Liang et al. 2024b; Xie et al. 2024; Yao et al. 2025; Zhang et al. 2024c] diffusion models, distilling their generalized motion or 3D priors into 4D sequences. Some approaches directly construct prior models, including diffusion [Gat et al. 2025; Jiang et al. 2024a; Liang et al. 2024b] and reconstruction [Ren et al. 2024] models, in the 4D domain. Others distill 4D motion from a

combination of generative models operating in lower dimensions, such as images, videos, and multi-view images.

Most relevant to *AnimaX* are Diffusion4D [Liang et al. 2024b], Animate3D [Jiang et al. 2024a], and MotionDreamer [Uzolas et al. 2024], which can generate animation from various 3D models. Diffusion4D and Animate3D train multi-view video diffusion models on 4D data, and distill their spatial-temporal prior into 4D generation. MotionDreamer extracts semantic motion priors from the deep features of video diffusion models to optimize deformation parameters in a zero-shot manner. However, these methods do not explicitly represent or generate 4D content in a physically grounded way.

A more recent work, Articulated Kinematics Distillation (AKD) [Li et al. 2025], combines traditional skeleton-based character animation pipelines with generative models, introducing a more physically plausible approach. Given a rigged 3D asset, AKD distills articulated motion sequences from video diffusion models using Score Distillation Sampling (SDS) [Poole et al. 2022]. However, AKD requires approximately 25 hours to optimize a single animation, motivating our focus on efficient, feed-forward articulated motion generation.

Video Diffusion Models. Video generation [Blattmann et al. 2023a,b; Chen et al. 2023b; Ho et al. 2022; Xing et al. 2024; Yu et al. 2023] has developed rapidly in recent years. Previous video diffusion models [Guo et al. 2024] usually build upon image diffusion models [Rombach et al. 2022], leveraging their pre-trained image knowledge by preserving spatial layers and inserting temporal layers to model motion dynamics. These methods typically utilize image VAEs, failing to compress temporal information effectively, thus limiting generation to very short videos. State-of-the-art video diffusion models [Bao et al. 2024; Kong et al. 2024; Wang et al. 2025; Yang et al. 2024b] incorporate 3D causal VAEs, compressing both spatial and temporal dimensions of videos, coupled with diffusion transformers (DiTs) [Peebles and Xie 2023] that denoise in the latent space. Trained on large-scale video datasets, these models demonstrated the capacity to generate diverse and realistic videos from textual prompts. These implicitly learned motion priors provide a foundation for category-agnostic and motion-diverse 3D animation.

3 METHODOLOGY

AnimaX generates 3D animations from a given articulated 3D mesh and a textual description of the motion. *AnimaX* has two stages: first, conditioned on the rendered views and pose maps from the mesh with a textual prompt, we generate multi-view consistent videos and corresponding pose sequences simultaneously using a joint video-pose diffusion model, and second, we reconstruct 3D motion from the generated multi-view pose sequences by multi-view triangulation and inverse kinematics (see Fig. 3). Below we describe our fundamental video diffusion models (Sec. 3.1), our multi-view video-pose diffusion model (Sec. 3.2), and how the generated pose sequences are recovered to a 3D animation (Sec. 3.3).

3.1 Preliminary: Video Diffusion Model

Video diffusion models [Kong et al. 2024; Wang et al. 2025] have demonstrated superior capabilities in generating generalized motion videos from textual descriptions, or animating arbitrary images

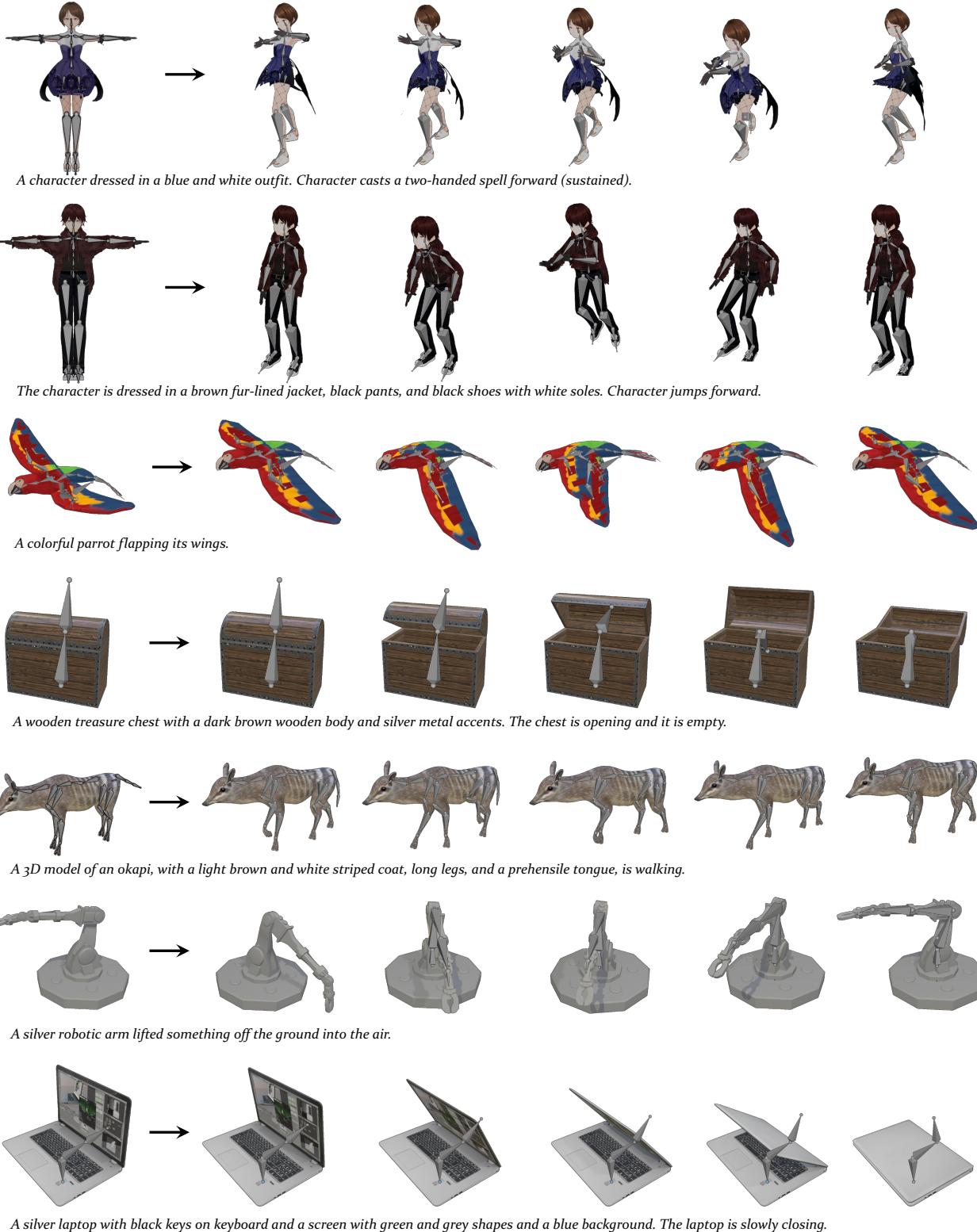


Fig. 2. Animation results on generalized 3D models, including biped 3D assets, animals, chests, robotic arms.

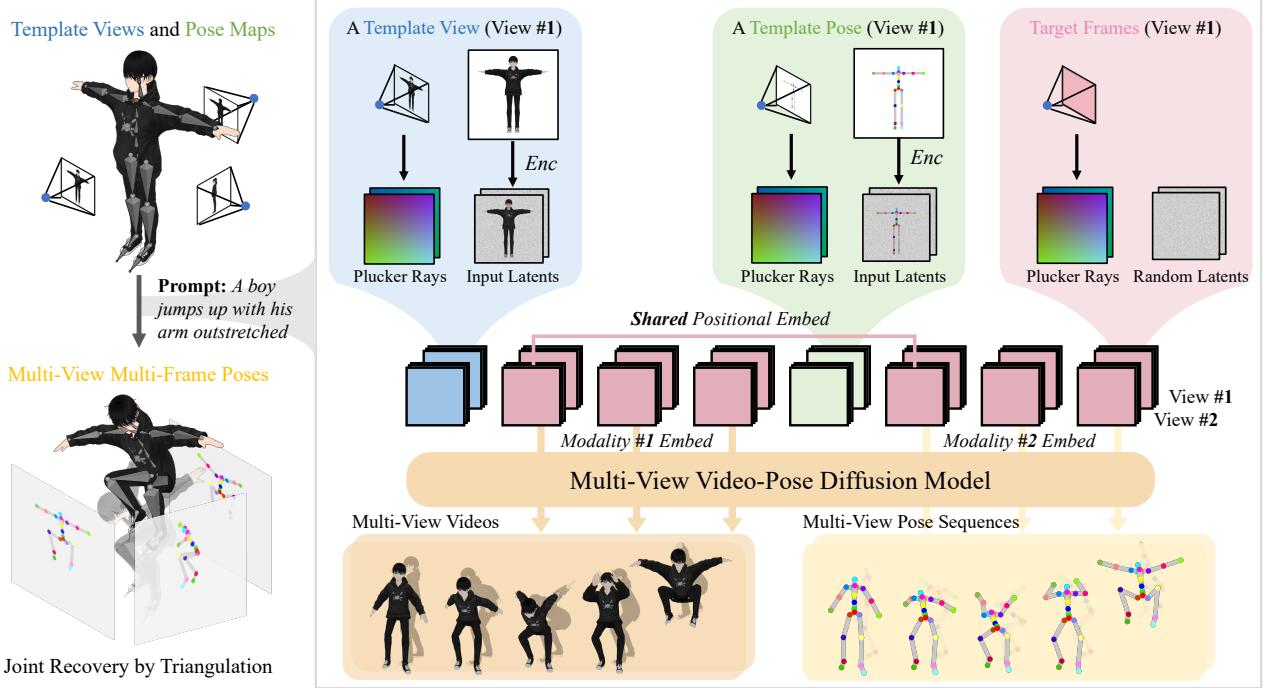


Fig. 3. Illustration of *AnimaX*. Given an articulated 3D mesh, *AnimaX* creates a sequence of 3D animation in minutes. *AnimaX* has two stages: (1) generating multi-view consistent videos and corresponding pose sequences simultaneously, conditioned on rendered template views and pose maps from the input mesh, with a textual description of the motion; and (2) recovering 3D joint positions per frame using multi-view triangulation [Hartley and Sturm 1997] and applying inverse kinematics to obtain the joint angles and animate the mesh.

guided by text. Typically, these models include a 3D causal variational auto-encoder (VAE) and a diffusion transformer (DiT) [Peebles and Xie 2023] ϵ_θ for denoising. The VAE compresses the video’s spatio-temporal dimensions; given a video $V \in \mathbb{R}^{(1+F) \times H \times W \times 3}$, the encodes maps it from pixel space to latent space $x \in \mathbb{R}^{(1+f) \times h \times w \times c}$. A diffusion transformer is then trained in this latent space, progressively denoising latent variables into video latents, which the VAE decoder reconstructs back into video frames.

In diffusion transformer, positional information of video latents x is typically encoded via RoPE [Su et al. 2024], applying rotation matrices based on each token’s coordinate (i, j, k) in a 3D grid:

$$\hat{x}^{i,j,k} = x^{i,j,k} \cdot R(i, j, k), \quad (1)$$

where $R(i, j, k)$ denotes the rotation matrix at position (i, j, k) with $0 \leq i < f$, $0 \leq j < w$, and $0 \leq k \leq h$. Subsequently, 3D attention is applied to these position-encoded tokens to capture both intra-frame and inter-frame relationships, while cross-attention incorporates textual conditioning information c^{txt} .

3.2 Multi-View Video-Pose Diffusion Model

We train a multi-view video-pose diffusion model that takes multi-view images and pose maps of a static 3D articulated mesh as input, and generates multi-view videos and pose sequences given a textual motion description, as depicted in Fig. 4. Specifically, given N template views containing N RGB images I^{rgb} , pose maps I^{pose} and their corresponding camera parameters C^{cam} , the model learns to

capture the joint distribution of N view RGB videos V^{rgb} and pose sequences V^{pose} with the guidance of a textual prompt C^{txt} :

$$p(V^{rgb}, V^{pose} | I^{rgb}, I^{pose}, C^{cam}, C^{txt}) \quad (2)$$

Pose Map Definition. We project the head positions of each bone in the skeletal animation onto the 2D image plane and assign a unique color to each joint for accurate localization during the subsequent recovery stage. In the rendered image, joints are visualized as circular markers, while the skeletal structure is depicted by connecting lines between parent and child joints.

Model Architecture. Our model architecture is initialized from video latent diffusion models [Wang et al. 2025], but with additional camera and modality embeddings as well as enhanced positional embeddings for joint video-pose generation. Given template RGB images, pose maps $I^{rgb}, I^{pose} \in \mathbb{R}^{H \times W \times 3}$ and target RGB videos and pose sequences $V^{rgb}, V^{pose} \in \mathbb{R}^{(1+F) \times H \times W \times 3}$ under N views, the model encodes each image and video into a latent representation through a 3D causal VAE, to obtain image latent tokens $c^{rgb}, c^{pose} \in \mathbb{R}^{1 \times h \times w \times c}$ and video latent tokens $x^{rgb}, x^{pose} \in \mathbb{R}^{(1+f) \times h \times w \times c}$ respectively. Then, a diffusion transformer [Peebles and Xie 2023], initialized from video DiT [Wang et al. 2025], is trained to estimate the joint distribution of the latent representations under multi-view capturing, given conditioning signals.

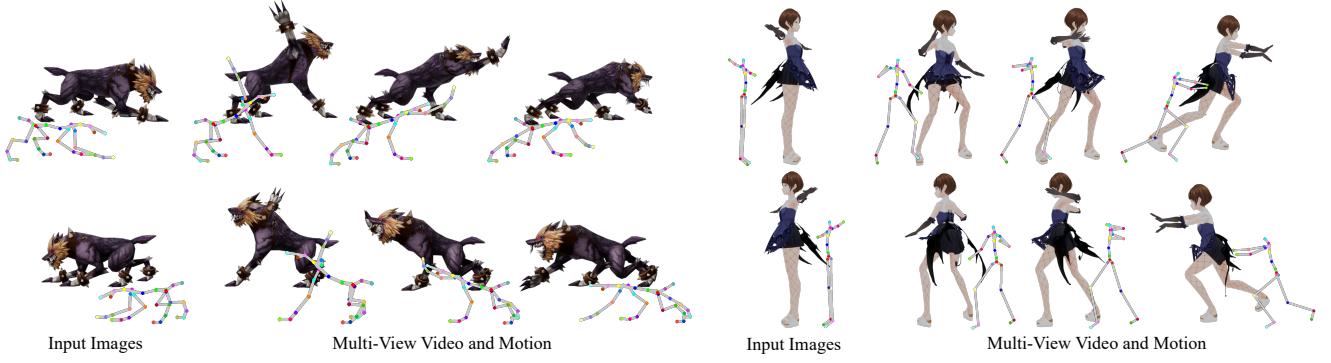


Fig. 4. Input and output examples of our multi-view video-pose diffusion models, with input prompts “A wolf is attacking something” and “A girl casting spell with hands forward”. We show two views here, but the model actually generates four views.

Image Conditioning and Temporal Modeling. We propose a simple yet unified approach that simultaneously enables image conditioning and consistent temporal modeling of two sequences of video latent tokens. Given the input image tokens c^{rgb} , c^{pose} and noisy video tokens x_t^{rgb} , x_t^{pose} , we concatenate them along the temporal dimension into a unified video token sequence $x^{total} \in \mathbb{R}^{(2f+4) \times h \times w \times c}$, where RGB video tokens are located in the first half while pose video tokens are in the second half. We then extend the processing scope of existing 3D self-attention layers to jointly attend over the entire temporally-padded sequence:

$$\text{Attention}([c^{rgb}, x_t^{rgb}, c^{pose}, x_t^{pose}]), \quad (3)$$

where c^{rgb} and c^{pose} denote the clean, non-noised conditioning tokens, while x_t^{rgb} and x_t^{pose} are noisy latent tokens at timestep t . This design leverages pre-trained 3D self-attention to seamlessly incorporate conditioning information from input images, efficiently utilizing the spatial-temporal priors of video diffusion models. As a result, it enables encoding of visual details with strong compatibility to the generative framework.

Cross-Modal Modeling. To distinguish between the two modalities—RGB and pose—in the unified video token sequence x^{total} , we introduce an additional modality embedding. We assign constant identifiers 0 and 1 to indicate the two modalities. These identifiers are further transformed via frequency encoding, followed by several linear layers, yielding embeddings that share the same dimensionality as the original timestep embeddings. The resulting embeddings are then added to the corresponding token’s timestep embedding.

Furthermore, considering that in x^{total} , the first half $[c^{rgb}, x_t^{rgb}]$ and the second half $[c^{pose}, x_t^{pose}]$ are spatially aligned, we introduce a shared positional encoding mechanism to enforce structural consistency. Specifically, for the sequence $x^{total} \in \mathbb{R}^{(2f+4) \times h \times w \times c}$, we define that tokens at positions (i, j, k) and $(i + f + 2, j, k)$ share the same positional encoding, formally expressed as:

$$\text{PE}^{i,j,k} = \text{PE}^{i+(f+2),j,k} = R(i, j, k), \quad (4)$$

where $R(i, j, k)$ denotes the rotation matrix used in RoPE [Su et al. 2024] at location (i, j, k) , with $0 \leq i < f + 2$, $0 \leq j < w$, and

$0 \leq k < h$. The design enables effective alignment and interaction between spatially corresponding tokens from different modalities.

Multi-View Consistency Modeling. To enable consistent multi-view video generation, we introduce additional camera conditioning and multi-view attention layers. Specifically, we adopt Plücker ray map to represent camera poses [Huang et al. 2024c; Sitzmann et al. 2021]. For each view, the corresponding ray map is concatenated channel-wise to the latent representations of both the input images and the generated video frames. To further enforce cross-view consistency, our multi-view layers operate on the multi-view video token sequence $x^{mv} \in \mathbb{R}^{N \times (2f+4) \times h \times w \times c}$, where N denotes the number of views. The token sequence is first inflated into $\hat{x}^{mv} \in \mathbb{R}^{(2f+4) \times (N \cdot h \cdot w) \times c}$, and self-attention is performed across the spatial dimension that aggregates all views. This formulation enables the model to directly learn spatial correspondences and enforce consistency across different camera viewpoints.

3.3 3D Motion Reconstruction and Animation

After obtaining multi-view pose sequences, we recover 3D poses and animate the 3D mesh through a three-stage process. 1) *2D joint localization*: For each frame, we first extract 2D joint positions $p^{1:v}$ by clustering [Arthur and Vassilvitskii 2006] the colors in the pose maps corresponding to each joint and taking the cluster centers as joint coordinates. 2) *3D joint optimization via triangulation*: We then estimate the 3D joint positions $P^{1:v}$ by solving a non-linear least-squares optimization problem [Hartley and Zisserman 2003]. The objective is to minimize the re-projection error between the projected 3D joints and the observed multi-view 2D joint positions $\{p^{1:v}\}_{v=1}^N$, while enforcing bone length consistency. 3) *Kinematic parameter estimation*: Based on the joint positions in both the template pose and the predicted pose, we apply the inverse process of forward kinematics to estimate the animation parameters [Aristidou and Lasenby 2011]. Traversing from the root node to the end-effectors, the rotation angle of each joint is estimated based on its positional deviation relative to the template pose. The resulting joint angles are then applied to animate the articulated 3D mesh.

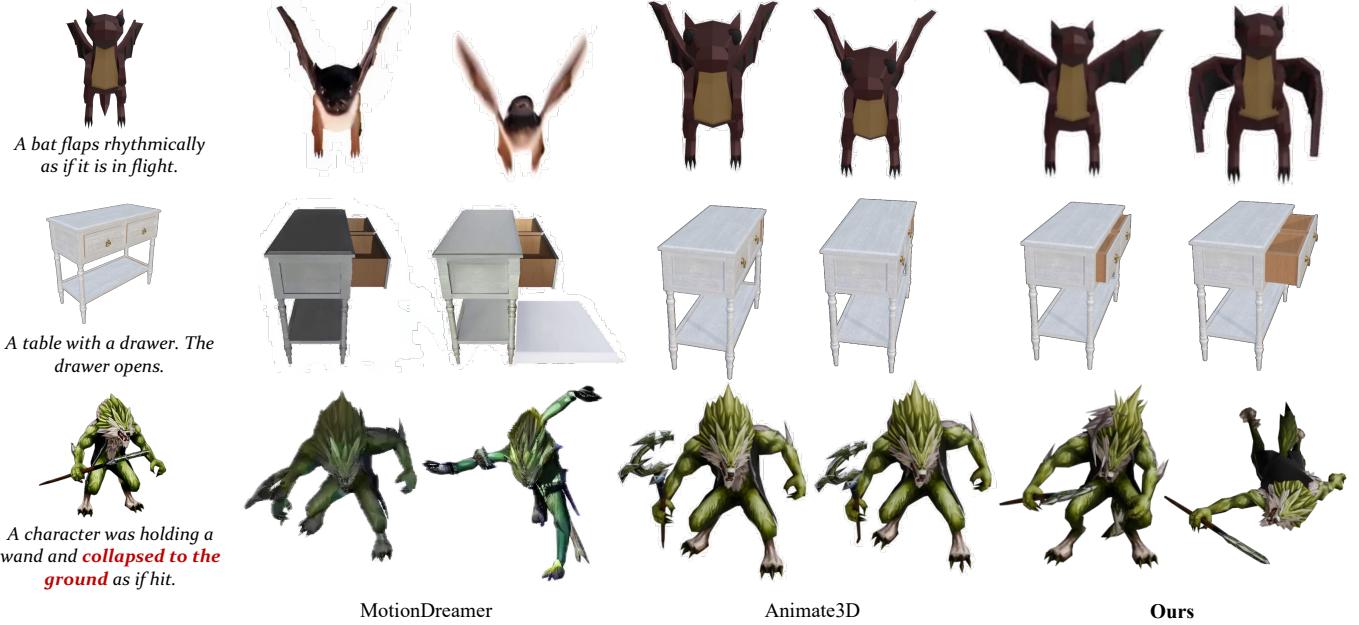


Fig. 5. Comparison with state-of-the-art generalizable 4D generation methods. We compare our method with representative 3D-to-4D methods, including MotionDreamer [Uzolas et al. 2024] and Animate3D [Jiang et al. 2024a]. Our model can synthesize more correct and authentic animation clips compared to these methods which rely on optimization of neural deformation fields and do not involve low-level skeleton-based representation.

4 EXPERIMENTS

We used Objaverse [Deitke et al. 2023a,b], Mixamo [Family 2022], VRoid [Hub 2022] as our raw data source, and extracted a total of 161,023 animation clips after processing and cleaning. Of these, we pick out 35 data pairs for evaluation, covering various categories such as humanoids, quadrupeds, birds, cabinets. We render multi-view videos from the generated animation and use VBench [Huang et al. 2024b] to evaluate on them.

We implement our multi-view video-pose diffusion model based on the Wan2.1 [Wang et al. 2025] text-to-video diffusion architecture, using the 1.3B parameter variant. A two-stage training strategy is employed. In the first stage, we fine-tune a single-view joint video-pose diffusion model using the LoRA [Hu et al. 2021] technique to efficiently adapt the pretrained backbone. In the second stage, we freeze all pretrained weights and train only the newly introduced camera embeddings and multi-view attention layers, thereby extending the model to support multi-view video-pose generation without disrupting the original learned priors.

4.1 Main Results and Comparisons

We present our primary results, which include the animated 3D conditioned on a 3D mesh and a text prompt, as illustrated in Fig. 1. In Fig. 4, we display examples of our multi-view video-pose diffusion models. Our results demonstrate that the model, adapted from pretrained video diffusion models, jointly generates consistent and high-fidelity videos and motion sequences.

Qualitative Comparisons. We perform a qualitative comparison with Animate3D [Jiang et al. 2024a] and MotionDreamer [Uzolas

Table 2. Quantitative comparisons. *I2V Subject*, *Smooth.*, *Dynamic Deg.*, *Quality* in VBench [Huang et al. 2024b] are used to evaluate the consistency with the given image, the motion smoothness, the motion degree, and the appearance quality, respectively. Values of all metrics are the higher, the better, except for *Dynamic Deg.*, since completely failed results (e.g., subject disappears) present an extremely high degree.

Methods	I2V Subject↑	Smooth.↑	Dynamic Deg.	Quality↑
Animate3D	0.943	0.986	0.446	0.481
MotionDreamer	0.817	0.977	0.827	0.439
Ours	0.962	0.990	0.661	0.517

Table 3. User study results on 3D animation. We collected user preference on motion-text alignment, 3D shape consistency, and overall motion quality from 30 participants. Our method receives the best preference on all metrics.

Methods	Motion-Text Align.↑	Shape Consist.↑	Overall Motion.↑
Animate3D	12.8%	26.7%	19.2%
MotionDreamer	4.3%	0.0%	2.9%
Ours	82.9%	73.3%	77.9%

et al. 2024], two representative 3D-to-4D methods that utilize multi-view video diffusion models to guide the optimization of neural deformation fields. As shown in Fig. 5, both baselines exhibit notable limitations. MotionDreamer relies on pretrained video diffusion models to supervise deformation optimization; however, the excessive degrees of freedom (DoFs) inherent in deformation fields lead to inconsistent geometry and unstable temporal behavior. Animate3D fine-tunes a multi-view video diffusion model to improve

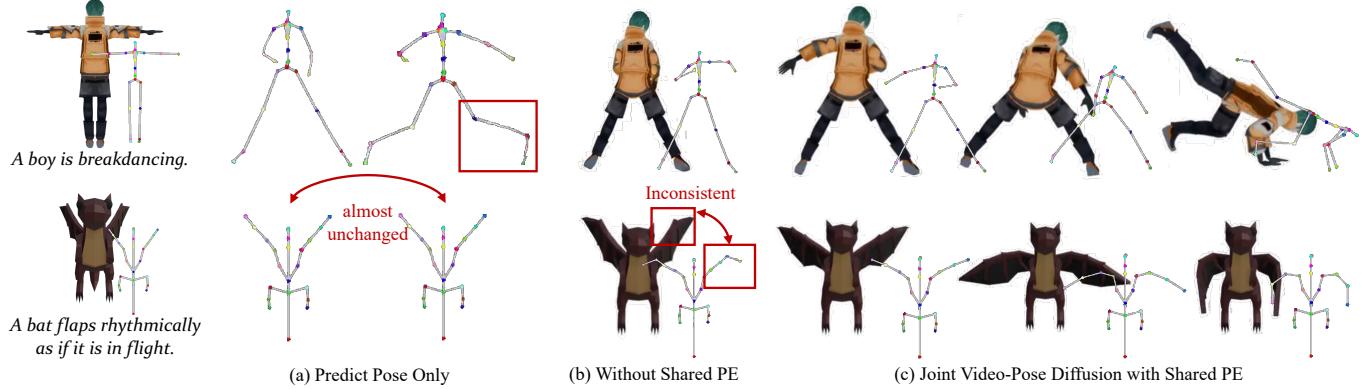


Fig. 6. Qualitative ablation results on the joint video-pose diffusion model.

cross-view consistency. While it reduces artifacts, the reconstruction remains challenging and often results in near-static outputs. In contrast, our method transfers video-based motion priors into the motion synthesis task via joint video-pose modeling. This enables the generation of temporally coherent and semantically aligned 3D animations, demonstrating superior consistency and motion expressiveness without requiring costly optimization procedures.

Quantitative Comparisons. Following [Jiang et al. 2024a], we render multi-view videos from animated 3D and evaluate them with Vbench [Huang et al. 2024b], a comprehensive video evaluation toolkit. We choose 4 image-to-video metrics, *i.e.*, *I2V Subject*, *Motion Smoothness*, *Dynamic Degree*, and *Aesthetic Quality*, measuring the consistency with the given image, the motion smoothness, the motion degree, and the appearance quality, respectively. Values of all metrics are the higher, the better. The *Dynamic Deg.* metric aims to capture motion richness; however, it can be less robust, as severe generation failures (*e.g.*, subject disappearing) may also produce misleadingly high scores. As shown in Tab. 2, our method outperforms other methods, especially in appearance quality, due to our low-DoF animation design.

User Study. We conducted a user study comparing our approach with baseline 3D-to-4D methods [Jiang et al. 2024a; Uzolas et al. 2024]. The study aimed to evaluate motion-text alignment, 3D shape consistency, and overall motion quality. A total of 30 participants were recruited to provide their preferences between the outputs of difference methods on test set. As shown in Tab. 3, our method receives the best preference on all metrics. This highlights the superior capability of our method in transferring the motion priors of video diffusion models to skeleton-based 3D animation.

4.2 Ablation Studies

We ablate the key design in our joint video-pose diffusion models. Specifically, based on video diffusion models, we examine three settings: (a) a pose diffusion model that is fine-tuned to generate pose sequences alone, (b) a video-pose diffusion model that generates videos and pose sequences simultaneously but do not share the same positional encoding in these two modalities, and (c) our full video-pose diffusion model with shared positional encoding mechanism.

Table 4. Quantitative ablation results on the video-pose diffusin model design. Based on the video diffusion model, we fine-tune a model that (a) only predicts pose sequences and (b) a joint video-pose diffusion model without shared positional encoding mechanism, and (c) our full setting. Vbench [Huang et al. 2024b] is used to evaluate on multi-view renderings.

Methods	I2V Subject↑	Smooth.↑	Dynamic Deg.	Quality↑
(a) Pose Only	0.960	0.982	0.402	0.448
(b) w/o Shared PE	0.954	0.988	0.660	0.429
(c) Full Setting	0.962	0.990	0.661	0.517

As illustrated in Fig. 6 (a), directly fine-tuning a pre-trained video diffusion model to generate only pose sequences often leads to degraded performance due to the sparsity of pose representation and the significant modality gap between RGB videos and pose maps, as well as the relatively sparse supervision available for the pose modality. This mismatch disrupts the learned spatial-temporal priors of the original model, frequently resulting in degenerate outputs such as distorted pose frames or nearly static pose sequences. Therefore, we adopt a joint video-pose diffusion framework that simultaneously generates both RGB video frames and pose sequences. Within this setting, we observe that compared to a baseline model (b), our full model—which shares positional encodings across the two modalities—significantly improves the spatial alignment between generated pose sequences and RGB videos. This architectural design facilitates more effective transfer of pre-trained video priors to the motion generation task, resulting in more coherent and realistic pose outputs. We also lifted the pose sequences generated by these three models to 3D animation clips, and used Vbench [Huang et al. 2024b] to evaluate on them. As shown in Tab. 4, model (c) performs the best, confirming the effectiveness of our model design.

5 CONCLUSION

We present AnimaX, a feed-forward framework for animating articulated 3D meshes with arbitrary skeletal structures by bridging the generalizable motion priors of video diffusion models with the structured controllability of skeleton-based animation. Unlike prior approaches that either rely on fixed skeletal topologies or

require costly optimization, our method enables efficient generation of temporally and spatially consistent multi-view pose and video sequences conditioned on a textual motion prompt. By introducing joint video-pose diffusion, shared positional encodings, and modality-aware embeddings, AnimaX effectively transfers video-based motion knowledge to the 3D domain and supports a broad spectrum of mesh categories. Extensive experiments on VBench validate the superiority of our method in terms of generalization, animation quality, and runtime efficiency. We believe this work opens new avenues for scalable, category-agnostic 3D animation driven by text and visual priors.

REFERENCES

- Chaitanya Ahuja and Louis-Philippe Morency. 2019. Language2Pose: Natural Language Grounded Pose Forecasting. In *2019 International Conference on 3D Vision, 3DV 2019, Québec City, QC, Canada, September 16–19, 2019*. IEEE, 719–728. <https://doi.org/10.1109/3DV400084>
- Andreas Aristidou and Joan Lasenby. 2011. FABRIK: A fast, iterative solver for the Inverse Kinematics problem. *Graphical Models* 73, 5 (2011), 243–260.
- David Arthur and Sergei Vassilvitskii. 2006. *k-means++: The advantages of careful seeding*. Technical Report. Stanford.
- Samaneh Azadi, Akbar Shah, Thomas Hayes, Devi Parikh, and Sonal Gupta. 2023. Make-An-Animation: Large-Scale Text-conditional 3D Human Motion Generation. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1–6, 2023*. IEEE, 14993–15002. <https://doi.org/10.1109/ICCV51070.2023.01381>
- Sherwin Bahmani, Xian Liu, Wang Yifan, Ivan Skorokhodov, Victor Rong, Ziwei Liu, Xihui Liu, Jeong Joon Park, Sergey Tulyakov, Gordon Wetzstein, Andrea Tagliasacchi, and David B. Lindell. 2024a. TC4D: Trajectory-Conditioned Text-to-4D Generation. In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29–October 4, 2024, Proceedings, Part XLVI (Lecture Notes in Computer Science, Vol. 15104)*, Ales Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gür Varol (Eds.). Springer, 53–72. https://doi.org/10.1007/978-3-031-72952-2_4
- Sherwin Bahmani, Ivan Skorokhodov, Victor Rong, Gordon Wetzstein, Leonidas J. Guibas, Peter Wonka, Sergey Tulyakov, Jeong Joon Park, Andrea Tagliasacchi, and David B. Lindell. 2024b. 4D-fy: Text-to-4D Generation Using Hybrid Score Distillation Sampling. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16–22, 2024*. IEEE, 7996–8006. <https://doi.org/10.1109/CVPR52733.2024.00764>
- Jianhong Bai, Menghan Xia, Xintao Wang, Ziyang Yuan, Xiao Fu, Zuozhu Liu, Haoji Hu, Pengfei Wan, and Di Zhang. 2024. SynCamMaster: Synchronizing Multi-Camera Video Generation from Diverse Viewpoints. *arXiv preprint arXiv:2412.07760* (2024).
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiaob Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. 2025. Qwen2.5-VL Technical Report. *arXiv preprint arXiv:2502.13923* (2025).
- Fan Bao, Chendong Xiang, Gang Yue, Guande He, Hongzhou Zhu, Kaiwen Zheng, Min Zhao, Shilong Liu, Yaole Wang, and Jun Zhu. 2024. Vidu: a highly consistent, dynamic and skilled text-to-video generator with diffusion models. *arXiv preprint arXiv:2405.04233* (2024).
- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. 2023a. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127* (2023).
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. 2023b. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 22563–22575.
- Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J. Black. 2016. Keep it SMPL: Automatic Estimation of 3D Human Pose and Shape from a Single Image. In *Computer Vision – ECCV 2016 (Lecture Notes in Computer Science)*. Springer International Publishing.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann. 2024. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in Neural Information Processing Systems* 37 (2024), 24081–24125.
- Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xiantao Wang, et al. 2023b. Videocrafter1: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512* (2023).
- Jianqi Chen, Biao Zhang, Xiangjun Tang, and Peter Wonka. 2025. V2M4: 4D Mesh Animation Reconstruction from a Single Monocular Video. *CoRR abs/2503.09631* (2025). <https://doi.org/10.48550/ARXIV.2503.09631> arXiv:2503.09631
- Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, and Gang Yu. 2023a. Executing your Commands via Motion Diffusion in Latent Space. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17–24, 2023*. IEEE, 18000–18010. <https://doi.org/10.1109/CVPR52729.2023.01726>
- Karan Dalal, Daniel Koceja, Gashon Hussein, Jiarui Xu, Yue Zhao, Youjin Song, Shihao Han, Ka Chun Cheung, Jan Kautz, Carlos Guestrin, et al. 2025. One-minute video generation with test-time training. *arXiv preprint arXiv:2504.05298* (2025).
- Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhabek Gadre, et al. 2023a. Objaverse-xl: A universe of 10m+ 3d objects. *Advances in Neural Information Processing Systems* 36 (2023), 35799–35813.
- Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. 2023b. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 13142–13153.
- Patrick Esser, Robin Rombach, and Björn Ommer. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 12873–12883.
- Adobe Family. 2022. Mixamo. <https://www.mixamo.com/> (2022).
- Ruiqi Gao, Aleksander Holynski, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla, Pratul Srinivasan, Jonathan T Barron, and Ben Poole. 2024. Cat3d: Create anything in 3d with multi-view diffusion models. *arXiv preprint arXiv:2405.10314* (2024).
- Inbar Gat, Sigal Raab, Guy Tevet, Yuval Reshef, Amit H Bermano, and Daniel Cohen-Or. 2025. AnyTop: Character Animation Diffusion with Any Topology. *arXiv preprint arXiv:2502.17327* (2025).
- Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. 2022. Generating Diverse and Natural 3D Human Motions from Text. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18–24, 2022*. IEEE, 5142–5151. <https://doi.org/10.1109/CVPR52688.2022.00509>
- Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahu Lin, and Bo Dai. 2024. AnimateDiff: Animate Your Personalized Text-to-Image Diffusion Models without Specific Tuning. In *ICLR*.
- Richard Hartley and Andrew Zisserman. 2003. *Multiple view geometry in computer vision*. Cambridge university press.
- Richard I Hartley and Peter Sturm. 1997. Triangulation. *Computer vision and image understanding* 68, 2 (1997), 146–157.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems* 33 (2020), 6840–6851.
- Jonathan Ho, Tim Salimans, Alexey Grigchenko, William Chan, Mohammad Norouzi, and David J Fleet. 2022. Video diffusion models. *Advances in Neural Information Processing Systems* 35 (2022), 8633–8646.
- 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).
- Zehuan Huang. 2025. Bpy-renderer: A Go-To Library for Rendering 3D Scenes and Animations. <https://github.com/huangzh/bpy-renderer>.
- Zehuan Huang, Yuan-Chen Guo, Haoran Wang, Ran Yi, Lizhuang Ma, Yan-Pei Cao, and Lu Sheng. 2024a. Mv-adapter: Multi-view consistent image generation made easy. *arXiv preprint arXiv:2412.03632* (2024).
- Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahu Lin, Yu Qiao, and Ziwei Liu. 2024b. VBench: Comprehensive Benchmark Suite for Video Generative Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Zehuan Huang, Hao Wen, Junting Dong, Yaohui Wang, Yangguang Li, Xinyuan Chen, Yan-Pei Cao, Ding Liang, Yu Qiao, Bo Dai, et al. 2024c. Epidiff: Enhancing multi-view synthesis via localized epipolar-constrained diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9784–9794.
- VRoid Hub. 2022. VRoid. <https://vroid.com/> (2022).
- Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. 2023. MotionGPT: Human Motion as a Foreign Language. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (Eds.). http://papers.nips.cc/paper_files/paper/2023/hash/3fbfc1ea0716c03dea93bb6e78dd6f-Abstract-Conference.html
- Yanqin Jiang, Chaohui Yu, Chenjie Cao, Fan Wang, Weiming Hu, and Jin Gao. 2024a. Animate3D: Animating Any 3D Model with Multi-view Video Diffusion. In *Advances*

- in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024, Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (Eds.).* http://papers.nips.cc/paper_files/paper/2024/hash/e3b53f89136b1bc69a5714ea4a65f01be-Abstract-Conference.html
- Yanqin Jiang, Li Zhang, Jin Gao, Weiming Hu, and Yao Yao. 2024b. Consistent4D: Consistent 360° Dynamic Object Generation from Monocular Video. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net. <https://openreview.net/forum?id=sPUrdGepF>
- Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangueng Xiong, Xin Lin, Bo Wu, Jianwei Zhang, et al. 2024. Hunyuanyideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603* (2024).
- Black Forest Labs. 2024. FLUX. <https://github.com/black-forest-labs/flux>.
- Jiaman Li, C Karen Liu, and Jiajun Wu. 2024. Lifting Motion to the 3D World via 2D Diffusion. *arXiv preprint arXiv:2411.18808* (2024).
- Xuan Li, Qianli Ma, Tsung-Yi Lin, Yongxin Chen, Chenfanfu Jiang, Ming-Yu Liu, and Donglai Xiang. 2025. Articulated Kinematics Distillation from Video Diffusion Models. *arXiv preprint arXiv:2504.01204* (2025).
- Han Liang, Jiacheng Bao, Ruichi Zhang, Sihan Ren, Yuecheng Xu, Sibei Yang, Xin Chen, Jingyi Yu, and Lan Xu. 2024a. OMG: Towards Open-vocabulary Motion Generation via Mixture of Controllers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 482–493. <https://doi.org/10.1109/CVPR52733.2024.00053>
- Hanwen Liang, Yuyang Yin, Dejia Xu, Hanxue Liang, Zhangyang Wang, Konstantinos N. Plataniotis, Yao Zhao, and Yunchao Wei. 2024b. Diffusion4D: Fast Spatial-temporal Consistent 4D generation via Video Diffusion Models. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024, Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (Eds.).* http://papers.nips.cc/paper_files/paper/2024/hash/c7f4dbbf3739b36029ba71a47844696-Abstract-Conference.html
- Huan Ling, Seung Wook Kim, Antonio Torralba, Sanja Fidler, and Karsten Kreis. 2024. Align Your Gaussians: Text-to-4D with Dynamic 3D Gaussians and Composed Diffusion Models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*. IEEE, 8576–8588. <https://doi.org/10.1109/CVPR52733.2024.00819>
- Tianqi Liu, Zihao Huang, Zhaoxi Chen, Guangcong Wang, Shoukang Hu, Liao Shen, Huiqiang Sun, Zhiguo Cao, Wei Li, and Ziwei Liu. 2025. Free4D: Tuning-free 4D Scene Generation with Spatial-Temporal Consistency. *CoRR abs/2503.20785* (2025). [https://doi.org/10.48550/ARXIV.2503.20785 arXiv:2503.20785](https://doi.org/10.48550/ARXIV.2503.20785)
- Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. 2023. Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453* (2023).
- Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. 2015. SMPL: a skinned multi-person linear model. *ACM Trans. Graph.* 34, 6 (2015), 248:1–248:16. <https://doi.org/10.1145/2816795.2818013>
- Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. 2019. AMASS: Archive of Motion Capture As Surface Shapes. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 5441–5450. <https://doi.org/10.1109/ICCV.2019.000554>
- Marc Benedi San Millán, Angelia Dai, and Matthias Nießner. 2025. Animating the Uncaptured: Humanoid Mesh Animation with Video Diffusion Models. *CoRR abs/2503.15996* (2025). <https://doi.org/10.48550/ARXIV.2503.15996 arXiv:2503.15996>
- OpenAI. 2023. GPT-4o. <https://openai.com/index/hello-gpt-4o/> (2023).
- William Peebles and Saining Xie. 2023. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 4195–4205.
- Mathis Petrovich, Michael J. Black, and Gülcin Varol. 2022. TEMOS: Generating Diverse Human Motions from Textual Descriptions. In *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXII (Lecture Notes in Computer Science, Vol. 13682)*, Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (Eds.). Springer, 480–497. https://doi.org/10.1007/978-3-031-20047-2_28
- Huaijin Pi, Ruoxi Guo, Zehong Shen, Qing Shuai, Zechen Hu, Zhumei Wang, Yajiao Dong, Ruizhen Hu, Taku Komura, Sida Peng, and Xiaowei Zhou. 2024. Motion-2-to-3: Leveraging 2D Motion Data to Boost 3D Motion Generation. *CoRR abs/2412.13111* (2024). <https://doi.org/10.48550/ARXIV.2412.13111 arXiv:2412.13111>
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Jon Penna, and Robin Rombach. 2023. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952* (2023).
- Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988* (2022).
- Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. 2021. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10318–10327.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- Jiawei Ren, Liang Pan, Jiaxiang Tang, Chi Zhang, Ang Cao, Gang Zeng, and Ziwei Liu. 2023. DreamGaussian4D: Generative 4D Gaussian Splatting. *CoRR abs/2312.17142* (2023). <https://doi.org/10.48550/ARXIV.2312.17142 arXiv:2312.17142>
- Jiawei Ren, Cheng Xie, Ashkan Mirzaei, Hanxue Liang, Xiaohui Zeng, Karsten Kreis, Ziwei Liu, Antonio Torralba, Sanja Fidler, Seung Wook Kim, and Huan Ling. 2024. L4GM: Large 4D Gaussian Reconstruction Model. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (Eds.). http://papers.nips.cc/paper_files/paper/2024/hash/6808f2c57d9564a2639a4710e3bb9b9-Abstract-Conference.html
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.
- Yahao Shi, Yang Liu, Yanmin Wu, Xing Liu, Chen Zhao, Jie Luo, and Bin Zhou. 2025. Drive Any Mesh: 4D Latent Diffusion for Mesh Deformation from Video. *arXiv preprint arXiv:2506.07489* (2025).
- Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. 2023. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512* (2023).
- Uriel Singer, Adam Polyk, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. 2022. Make-A-Video: Text-to-Video Generation without Text-Video Data. In *The Eleventh International Conference on Learning Representations*.
- Vincent Sitzmann, Semon Rezchikov, Bill Freeman, Josh Tenenbaum, and Fredo Durand. 2021. Light field networks: Neural scene representations with single-evaluation rendering. *Advances in Neural Information Processing Systems* 34 (2021), 19313–19325.
- Chaooyue Song, Jianfeng Zhang, Xiu Li, Fan Yang, Yiwen Chen, Zhongcong Xu, Jun Hao Lieuw, Xiaoyang Guo, Fayao Liu, Jiaoshi Feng, et al. 2025. MagicArticulate: Make Your 3D Models Articulation-Ready. *arXiv preprint arXiv:2502.12135* (2025).
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing* 568 (2024), 127063.
- Qi Sun, Zhiyang Guo, Ziyu Wan, Jing Nathan Yan, Shengming Yin, Wengang Zhou, Jing Liao, and Houqiang Li. 2024. EG4D: Explicit Generation of 4D Object without Score Distillation. *CoRR abs/2405.18132* (2024). <https://doi.org/10.48550/ARXIV.2405.18132 arXiv:2405.18132>
- Guy Tevet, Brian Gordon, Amir Hertz, Amit H. Bermano, and Daniel Cohen-Or. 2022. MotionCLIP: Exposing Human Motion Generation to CLIP Space. In *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXII (Lecture Notes in Computer Science, Vol. 13682)*, Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (Eds.). Springer, 358–374. https://doi.org/10.1007/978-3-031-20047-2_21
- Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafit, Daniel Cohen-Or, and Amit Haim Bermano. 2023. Human Motion Diffusion Model. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. <https://openreview.net/forum?id=SJ1kSyOjwu>
- Lukas Uzelas, Elmar Eisemann, and Petr Kellnhofer. 2024. MotionDreamer: Zero-Shot 3D Mesh Animation from Video Diffusion Models. *CoRR abs/2405.20155* (2024). <https://doi.org/10.48550/ARXIV.2405.20155 arXiv:2405.20155>
- Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuena, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas Wolf. 2022. Diffusers: State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>.
- Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, et al. 2025. Wan: Open and Advanced Large-Scale Video Generative Models. *arXiv preprint arXiv:2503.20314* (2025).
- Hao Wen, Zehuan Huang, Yaohui Wang, Xinyuan Chen, Yu Qiao, and Lu Sheng. 2024. Ouroboros3d: Image-to-3d generation via 3d-aware recursive diffusion. *arXiv preprint arXiv:2406.03184* (2024).
- Guanjun Wu, Taoran Yi, Jiemi Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 2024. 4d gaussian splatting for real-time dynamic scene rendering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 20310–20320.
- Zijie Wu, Chaohui Yu, Fan Wang, and Xiang Bai. 2025. AnimateAnyMesh: A Feed-Forward 4D Foundation Model for Text-Driven Universal Mesh Animation. *arXiv preprint arXiv:2506.09982* (2025).
- Yiming Xie, Chun-Han Yao, Vikram Voleti, Huaizu Jiang, and Varun Jampani. 2024. SV4D: Dynamic 3D Content Generation with Multi-Frame and Multi-View Consistency. *CoRR abs/2407.17470* (2024). <https://doi.org/10.48550/ARXIV.2407.17470 arXiv:2407.17470>
- Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Wangbo Yu, Hanyuan Liu, Gongye Liu, Xintao Wang, Ying Shan, and Tien-Tsin Wong. 2024. Dynamicroft: Animating open-domain images with video diffusion priors. In *European Conference on Computer Vision*. Springer, 399–417.

- Zeyu Yang, Zijie Pan, Chun Gu, and Li Zhang. 2024a. Diffusion²: Dynamic 3D Content Generation via Score Composition of Orthogonal Diffusion Models. *CoRR* abs/2404.02148 (2024). <https://doi.org/10.48550/ARXIV.2404.02148> arXiv:2404.02148
- Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. 2024b. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072* (2024).
- Chun-Han Yao, Yiming Xie, Vikram Voleti, Huaizu Jiang, and Varun Jampani. 2025. SV4D 2.0: Enhancing Spatio-Temporal Consistency in Multi-View Video Diffusion for High-Quality 4D Generation. *arXiv preprint arXiv:2503.16396* (2025).
- Jiwen Yu, Xiaodong Cun, Chenyang Qi, Yong Zhang, Xintao Wang, Ying Shan, and Jian Zhang. 2023. Animatezero: Video diffusion models are zero-shot image animators. *arXiv preprint arXiv:2312.03793* (2023).
- Yifei Zeng, Yanqin Jiang, Siyu Zhu, Yuanxun Lu, Youtian Lin, Hao Zhu, Weiming Hu, Xun Cao, and Yao Yao. 2024. STAG4D: Spatial-Temporal Anchored Generative 4D Gaussians. In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29–October 4, 2024, Proceedings, Part XXXVI (Lecture Notes in Computer Science, Vol. 15094)*, Ales Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gülcü Varol (Eds.). Springer, 163–179. https://doi.org/10.1007/978-3-031-72764-1_10
- Hao Zhang, Di Chang, Fang Li, Mohammad Soleymani, and Narendra Ahuja. 2024b. Magicpose4d: Crafting articulated models with appearance and motion control. *arXiv preprint arXiv:2405.14017* (2024).
- Haiyu Zhang, Xinyuan Chen, Yaohui Wang, Xihui Liu, Yunhong Wang, and Yu Qiao. 2024c. 4Diffusion: Multi-view Video Diffusion Model for 4D Generation. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 – 15, 2024*, Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (Eds.). http://papers.nips.cc/paper_files/paper/2024/hash/1bbfea488a8968e2d3c6565639b08e5e-Abstract-Conference.html
- Jia-Qi Zhang, Xiang Xu, Zhi-Meng Shen, Zehuan Huang, Yang Zhao, Yan-Pei Cao, Pengfei Wan, and Miao Wang. 2021. Write-An-Animation: High-level Text-based Animation Editing with Character-Scene Interaction. *Comput. Graph. Forum* 40, 7 (2021), 217–228. <https://doi.org/10.1111/CGF.14415>
- Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Yong Zhang, Hongwei Zhao, Hongtao Lu, Xi Shen, and Shan Ying. 2023. Generating Human Motion from Textual Descriptions with Discrete Representations. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17–24, 2023*. IEEE, 14730–14740. <https://doi.org/10.1109/CVPR52729.2023.01415>
- Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. 2024a. MotionDiffuse: Text-Driven Human Motion Generation With Diffusion Model. *IEEE Trans. Pattern Anal. Mach. Intell.* 46, 6 (2024), 4115–4128. <https://doi.org/10.1109/TPAMI.2024.3355414>
- Yaqi Zhang, Di Huang, Bin Liu, Shixiang Tang, Yan Lu, Lu Chen, Lei Bai, Qi Chu, Nenghai Yu, and Wanli Ouyang. 2024d. MotionGPT: Finetuned LLMs Are General-Purpose Motion Generators. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20–27, 2024, Vancouver, Canada*, Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (Eds.). AAAI Press, 7368–7376. <https://doi.org/10.1609/AAAI.V38I7.28567>
- Yuyang Zhao, Chung-Ching Lin, Kevin Lin, Zhiwen Yan, Linjie Li, Zhengyuan Yang, Jianfeng Wang, Gim Hee Lee, and Lijuan Wang. 2024. GenXD: Generating Any 3D and 4D Scenes. *CoRR* abs/2411.02319 (2024). <https://doi.org/10.48550/ARXIV.2411.02319> arXiv:2411.02319
- Hanxin Zhu, Tianyu He, Xiqian Yu, Junliang Guo, Zhibo Chen, and Jiang Bian. 2025. AR4D: Autoregressive 4D Generation from Monocular Videos. *CoRR* abs/2501.01722 (2025). <https://doi.org/10.48550/ARXIV.2501.01722> arXiv:2501.01722
- Qi Zuo, Xiaodong Gu, Lingteng Qiu, Yuan Dong, Zhengyi Zhao, Weihao Yuan, Rui Peng, Siyu Zhu, Zilong Dong, Liefeng Bo, et al. 2024. VideoMV: Consistent Multi-View Generation Based on Large Video Generative Model. *arXiv preprint arXiv:2403.12010* (2024).

A APPENDIX

A.1 Dataset

Data Curation. We use Objaverse [Deitke et al. 2023a,b], Mixamo [Family 2022], VRoid [Hub 2022] as our raw data source. For Objaverse and Objaverse-XL, we curate high-quality animation clips by applying the following filtering criteria: 1) The asset must contain both rigging and animation data. 2) Geometry and texture quality meet a fidelity threshold. 3) The animation sequence must contain more than 16 frames. 4) We discard clips with negligible motion

based on the mean optical flow magnitude. After filtering, we obtain a total of 48,020 animation clips from Objaverse and Objaverse-XL. For Mixamo and VRoid, we retarget motion-only animation clips from Mixamo humanoid characters onto both Mixamo and VRoid static characters. We then apply the same filtering strategy as used for Objaverse. This process yields 55,222 animation clips from Mixamo characters and 57,781 stylized animation clips from VRoid anime-style characters.

In total, we collect 161,023 high-quality animation clips across multiple categories. For evaluation, we sample 35 representative sequences covering a wide range of object categories, including humanoids, quadrupeds, flying animals, and articulated furniture.

Category Annotation. Following [Song et al. 2025], we render each 3D model from four predefined viewpoints and arrange the resulting images into a 2×2 grid. We then utilize GPT-4o [OpenAI 2023] to perform automatic category labeling based on the composed images.

Rendering. To support multi-view supervision, we render four-view videos and their corresponding poses from animation clips and their models in rest pose using Blender scripts [Huang 2025]. The camera setup consists of a fixed elevation angle of 0° , with azimuth angles set to 0° , 90° , 190° , and 270° , respectively. Additionally, we generate video captions using the vision-language model Qwen2.5-VL [Bai et al. 2025]. These paired multi-view videos and pose maps form the training data for our multi-view video-pose diffusion model.

A.2 Implementation Details

Training. We implement our multi-view video-pose diffusion model based on the Wan2.1 [Wang et al. 2025] text-to-video diffusion architecture, using the 1.3B parameter variant and diffusers codebase [von Platen et al. 2022]. A two-stage training strategy is employed. In the first stage, we fine-tune a single-view joint video-pose diffusion model using the LoRA [Hu et al. 2021] technique to efficiently adapt the pretrained backbone. In the second stage, we freeze all pretrained weights and train only the newly introduced camera embeddings and multi-view attention layers, thereby extending the model to support multi-view video-pose generation without disrupting the original learned priors.

The video-pose diffusion model is trained to generate videos at a resolution of 480p, with a sequence length of up to 81 frames. Since Mixamo and VRoid datasets primarily consist of humanoid characters, we assign sampling weights during training to balance contributions from different data sources. Specifically, training samples are drawn from Mixamo and VRoid with probabilities of 0.25 each, and from Objaverse [Deitke et al. 2023a,b] with a probability of 0.5. We randomly selected either a frame from the video or a rendered view of the mesh in its rest pose as the image condition. This image condition was dropped with a probability of 0.2. We train the model for 5 epochs using a learning rate of 5×10^{-5} on 16 NVIDIA A100 GPUs.

Inference. At inference time, given an articulated 3D mesh, we first render four-view images and corresponding pose maps as templates. These are fed into our multi-view video-pose diffusion model, along with a textual motion description. For the denoising process, we set

Table 5. Animation counts for each category in the curated dataset. Each column is sorted in descending order of animation count.

Category	# Animation	Category	# Animation	Category	# Animation	Category	# Animation
character	140,000	sculpture	1132	furniture	487	miscellaneous	196
anthropomorphic	22881	vehicle	1085	scanned data	425	planet	191
toy	12725	anatomy	1070	electronic device	410	sporting goods	140
animal	8603	household item	679	architecture	378	paper	77
mythical creature	5428	plant	606	clothing	304	jewelry	62
weapon	3221	accessory	570	food	270	musical instrument	31
tool	1297						

the image condition guidance scale to 3.0, and use 50 denoising steps. The model simultaneously generates four-view videos and pose sequences in approximately 5 minutes. Subsequently, we extract 2D joint positions from the generated pose maps and lift them to 3D joint angles via multi-view triangulation followed by kinematic parameter estimation. The total inference time, including both video generation and 3D reconstruction, is approximately 6 minutes.

A.3 Limitations and Discussions

Currently, our video-pose diffusion model generates videos from a fixed set of camera viewpoints with limited and static fields of view, making it challenging to synthesize animations with large spatial

motion. However, we believe that our network architecture could potentially handle such scenarios if trained on videos captured with more flexible and dynamic camera movements. Under this setting, it would be possible to specify arbitrary camera trajectories during inference, enabling the generation of multi-view videos that capture wider motion ranges and larger 3D spaces.

In addition, our model inherits the limitation of pretrained video diffusion backbones, which typically constrain the maximum video length. As a result, generating temporally continuous long-form animations remains difficult. We anticipate that incorporating test-time training [Dalal et al. 2025] or autoregressive denoising [Chen et al. 2024] extensions may offer a promising direction to support the generation of longer and more coherent video sequences.