

AnimateAnyMesh: A Feed-Forward 4D Foundation Model for Text-Driven Universal Mesh Animation

Zijie Wu^{1,2*}, Chaohui Yu^{2,3}, Fan Wang², Xiang Bai^{1†}

¹ Huazhong University of Science and Technology ² DAMO Academy, Alibaba Group ³ Hupan Lab

{zjw1031,xbai}@hust.edu.cn, {huakun.ych,fan.w}@alibaba-inc.com

<https://animateanymesh.github.io/AnimateAnyMesh/>

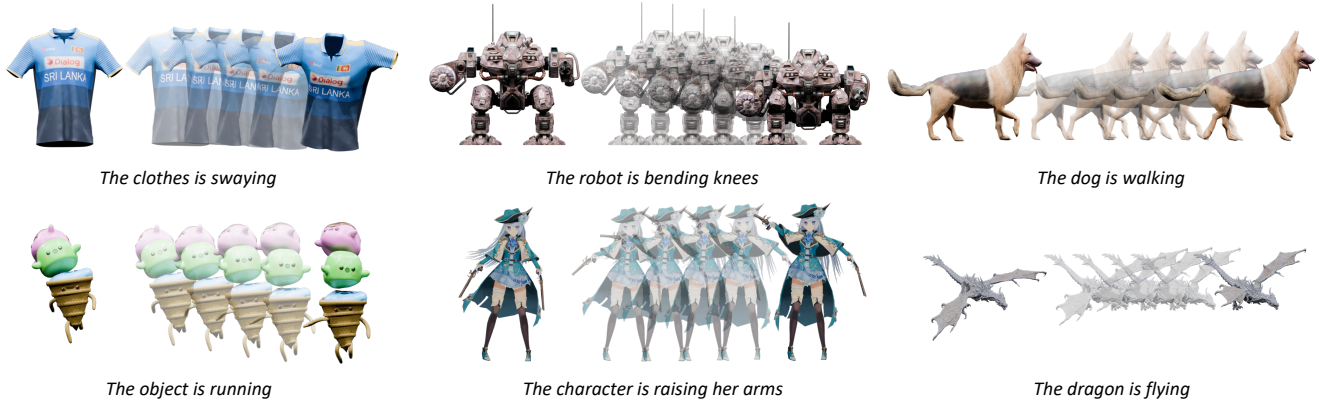


Figure 1. We present **AnimateAnyMesh**: the first feed-forward universal mesh animation framework that enables efficient motion generation for arbitrary 3D meshes. Given a static mesh and prompt, our method generates high-quality animations in only **a few seconds**.

Abstract

Recent advances in 4D content generation have attracted increasing attention, yet creating high-quality animated 3D models remains challenging due to the complexity of modeling spatio-temporal distributions and the scarcity of 4D training data. In this paper, we present **AnimateAnyMesh**, the first feed-forward framework that enables efficient text-driven animation of arbitrary 3D meshes. Our approach leverages a novel **DyMeshVAE** architecture that effectively compresses and reconstructs dynamic mesh sequences by disentangling spatial and temporal features while preserving local topological structures. To enable high-quality text-conditional generation, we employ a **Rectified Flow**-based training strategy in the compressed latent space. Additionally, we contribute the **DyMesh Dataset**, containing over 4M diverse dynamic mesh sequences with text annotations. Experimental results demonstrate that our method generates semantically accurate and temporally coherent

mesh animations in a few seconds, significantly outperforming existing approaches in both quality and efficiency. Our work marks a substantial step forward in making 4D content creation more accessible and practical. All the data, code, and models will be open-released.

1. Introduction

The revolution in 3D content creation has transformed various domains like VR/AR and gaming. While recent generative models [22, 32, 39, 40, 73, 84, 88, 89, 98, 100] excel at producing high-quality 3D assets, extending these advances to 4D content generation remains challenging due to the complexity of spatio-temporal modeling and scarcity of high-quality 4D assets.

Existing 4D generation approaches fall into two categories: per-scene optimization methods [4, 25, 58, 85, 101] and multi-view dynamic video generation methods [26, 36, 97]. The former employs SDS [54] with pre-trained generative models but suffers from high computational costs

*Work done during internship at DAMO Academy, Alibaba Group

† Corresponding author

and inconsistency. The latter fine-tunes video generation models [12, 20, 64, 81] on multi-view renderings of 4D assets [15, 16], followed by per-scene 4D reconstruction. Although inference efficiency improves, the need for post-processing impedes real-time application. Moreover, these methods, which typically adopt dynamic 3DGS [27] or NeRF [50] as 4D representations, suffer from view discrepancies due to the lack of ground-truth 4D data, relying solely on multi-view rendering supervision.

Given these limitations, we argue that dynamic meshes serve as an ideal representation for 4D content creation. As the de facto standard in modern graphics pipelines, mesh representations not only offer superior rendering efficiency but also enable natural decoupling of geometry and motion. Moreover, leveraging existing high-quality 3D meshes for animation, rather than pursuing direct 4D generation, enjoys several advantages: First, there exists a vast repository of high-quality 3D meshes, whether crafted by artists [15, 67] or synthesized by generative methods [84, 88], while high-quality 4D assets remain scarce. Second, decomposing 4D generation into geometry creation and motion modeling allows for higher fidelity and better control in both aspects. These observations motivate us to focus on the fundamental task of text-driven mesh animation.

To this end, we propose **AnimateAnyMesh**, the first feed-forward framework for text-driven universal mesh animation. At its core, we introduce **DyMeshVAE**, a novel VAE [29] architecture tailored for dynamic mesh sequences, aiming to compress and reconstruct the trajectory of each vertex. Specifically, DyMeshVAE first decomposes each trajectory into the initial vertex position and the relative trajectory. These components are then encoded using distinct positional encoding mechanisms, mapping spatial and temporal features into high-dimensional spaces. To enhance the quality of compression and reconstruction, we leverage the mesh topology by constructing the vertex connectivity matrix from the face information. This matrix serves as an attention mask to encode connectivity information into vertex features, which helps preserve topological structures and prevents trajectory entanglement during reconstruction (check Sec. 3.1 for details). During decoding, the enhanced vertex features of the initial mesh serve as queries to retrieve corresponding relative trajectories through cross-attention, enabling accurate trajectory reconstruction. Through this hierarchical design, DyMeshVAE effectively compresses meshes with varying numbers of vertices/faces into a fixed number of tokens (512 as default) during training while maintaining high-quality reconstruction. Moreover, being entirely attention-based, DyMeshVAE can dynamically adjust the number of encoded tokens during inference to accommodate meshes of different complexities, enabling animation of more intricate geometries.

To bridge the gap between text descriptions and mesh an-

imations, we propose **Shape-Guided Text-to-Trajectory Model**, which leverages a Rectified Flow-based [42] training strategy in the compressed latent space. By learning the conditional distribution of relative trajectories given text prompts and initial mesh features, our approach enables the generation of smooth and realistic animations. To facilitate this learning process, we introduce the **DyMesh Dataset**, a large-scale collection of over **4M** dynamic mesh sequences comprising about **2.6M** 16-frame and **1.6M** 32-frame sequences. This comprehensive dataset is carefully curated from diverse 4D assets through rigorous collection, filtering, processing, and annotation procedures, providing a robust foundation for training and evaluation.

Training with our DyMesh Dataset, AnimateAnyMesh achieves remarkable capabilities in generating high-quality animations that faithfully align with text descriptions for meshes of arbitrary topology. The feed-forward nature of our framework enables real-time inference, establishing AnimateAnyMesh as the first solution that combines universal mesh animation with both versatility and efficiency.

We highlight the contribution of this paper as follows:

1. We propose AnimateAnyMesh, the first feed-forward framework for text-driven universal mesh animation, enabling generation of high-quality animations for meshes of arbitrary topology within a few seconds.
2. We introduce DyMeshVAE, a novel compression-reconstruction architecture that effectively handles dynamic mesh sequences through trajectory decomposition and topology-aware attention mechanisms, naturally scaling to meshes of varying complexities.
3. We curate DyMesh Dataset, a large-scale collection of over 4M dynamic mesh sequences with text annotations, providing a robust foundation for 4D generation.
4. Extensive experiments demonstrate that our approach achieves state-of-the-art performance in text-driven mesh animation, combining high fidelity, versatility, and computational efficiency.

2. Related Works

3D Generation. Early approaches [14, 37, 54, 69, 71, 72, 80, 83, 87, 91, 93] for 3D generation leverage CLIP [57] score or Score Distillation Sampling (SDS) [54] to distill geometric priors from pre-trained 2D generative models [61, 62]. However, due to the inherent lack of 3D information, these methods often suffer from view discrepancy issues and require time-consuming per-scene optimization, significantly limiting their practical applications. To address these limitations, some methods [41, 43, 56, 63, 82] try to fine-tune 2D generative models using multi-view renderings of 3D assets [15, 16], thereby enhancing view consistency in the generated 3D content. Nevertheless, these approaches still require per-instance reconstruction after obtaining multi-view outputs. In contrast, more recent

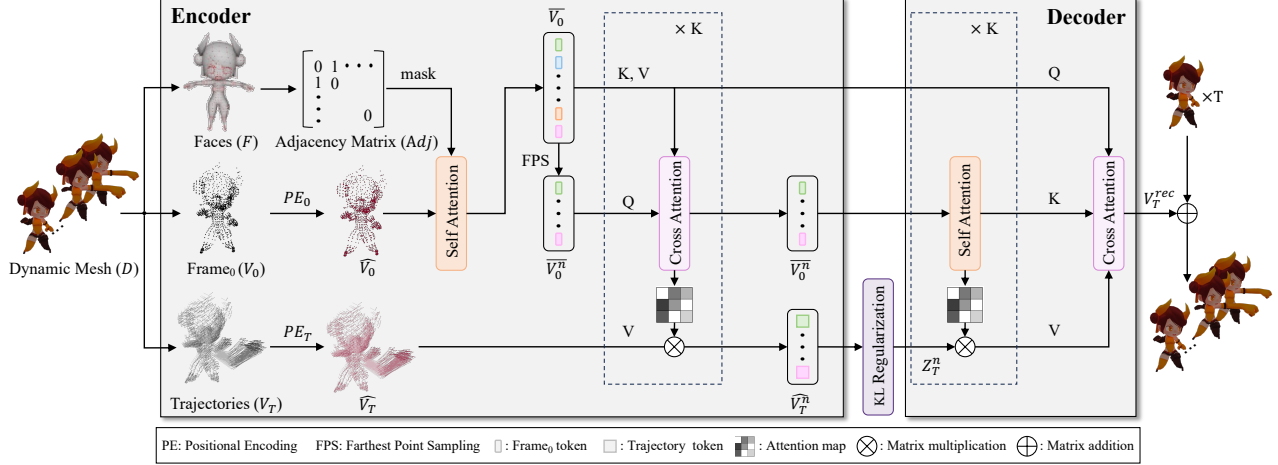


Figure 2. Illustration of our proposed **DyMeshVAE**. Given a dynamic mesh D , we first extract the initial frame vertex V_0 , the connectivity information from faces F , and the relative trajectories V_T . These information are then encoded to a decoupled latent space $\{\bar{V}_0^n, \bar{V}_T^n\}$ via the **Encoder**, which features trajectory decomposition and topology-aware attention mechanisms. Then the relative trajectories V_T^{rec} are reconstructed from the latent space via the **Decoder**. Finally, we add V_T^{rec} and V_0 to get the reconstructed dynamic mesh.

methods [22, 32, 39, 40, 73, 84, 88, 89, 98, 100] directly generate 3D representations [27, 50] and optimize the networks through multi-view rendering supervision, enabling rapid 3D asset generation within seconds for given prompts. This end-to-end approach circumvents the need for expensive post-processing optimization. Inspired by the evolution of 3D generation techniques, we pioneer a feed-forward 4D architecture for universal mesh animation. Our approach enables rapid mesh animation in a few seconds without per-scene optimization or reconstruction.

4D Generation. Following the evolution path of 3D generation, early 4D generation approaches [4, 25, 38, 52, 58, 66, 85, 92, 95, 101] attempt to distill spatio-temporal priors from 2D/3D/video generative models. Compared to 3D generation, these distillation approaches for 4D content creation not only demand significantly more computing power and longer optimization time to generate a single scene, but also tend to produce more noticeable spatio-temporal artifacts. Some recent works [26, 31, 36, 59, 97] have explored finetuning 3D/video generative models [12, 20, 64, 81] with 4D data to synthesize multi-view dynamic videos, aiming to accelerate 4D generation and improve spatio-temporal consistency. Nevertheless, these approaches still rely heavily on per-scene 4D reconstruction, and the generated objects exhibit spatial discrepancies due to the lack of true 4D training data. While some methods [9, 60, 74, 94, 99] have achieved efficient 4D generation in specific categories (e.g., human bodies) through parametric models [44] and modality-specific data [1–3, 7, 8, 10, 13, 18, 19, 21, 23, 30, 33, 34, 45–49, 51, 65, 70, 75–78], to the best of our knowledge, direct feed-forward 4D generation for general categories remains unexplored. In this work, we propose a novel framework that enables efficient feed-forward mesh

animation during inference by disentangling and compressing dynamic mesh shapes and motions. Then, we leverage Rectified Flow [42] to effectively model the posterior distribution of object motions conditioned on text and shape latents. This approach marks a significant advancement toward general-purpose feed-forward 4D content generation.

3. Method

Existing approaches to mesh animation typically rely on computationally intensive per-instance optimization [26, 58, 79] or are constrained to specific object categories [60, 74, 94, 99], limiting their practical applications. To address these limitations, we propose AnimateAnyMesh, a feed-forward framework that enables text-driven universal mesh animation. Our framework consists of two key components: DyMeshVAE and Shape-Guided Text-to-Trajectory Model. The former effectively handles meshes of arbitrary topology by decomposing them into initial frames and relative trajectories, which are then compressed into a structured latent space, while the latter learns to generate trajectory features conditioned on both the initial mesh latent and text embeddings. In this section, we elaborate on the architecture of DyMeshVAE and Shape-Guided Text-to-Trajectory Model in Sec. 3.1 and Sec. 3.2, respectively, followed by the training and inference pipeline of in Sec. 3.3.

3.1. DyMeshVAE

DyMeshVAE Encoder. As shown in Fig. 2, given a dynamic mesh sequence $D \subset \{F \in \mathbb{R}^{M \times 3}, V \in \mathbb{R}^{T \times N \times 3}\}$, we first disentangle the vertices sequence V into the initial frame vertices $V_0 \in \mathbb{R}^{N \times 3}$ and the relative trajectories $V_T \in \mathbb{R}^{N \times (T-3)}$, which satisfy the following:

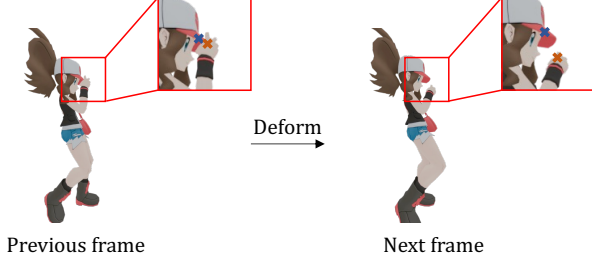


Figure 3. Demonstration of divergent trajectories for nearby mesh vertices in the initial frame.

$$V^t = V_0^t + V_T^t, \quad \text{for } t = 1, 2, \dots, T \quad (1)$$

where t stands for the index of the time sequence.

We decompose the vertex sequence V into initial positions V_0 and relative trajectories V_T based on our empirical observation that such decomposition leads to the disentanglement of shape and motion, while yielding a motion distribution that better approximates a zero-mean normal distribution. Following this decomposition, each trajectory is represented as a combination of its initial vertex position and the subsequent temporal offsets. The initial position serves as a spatial identifier, while the temporal offsets become our modeling target. To enhance trajectory reconstruction stability and prevent adhesion effects (as demonstrated in Fig. 8), we get inspirations from [50, 96] and employ distinct positional encoding schemes for V_0 and V_T , resulting in encoded features \widehat{V}_0 and \widehat{V}_T .

Due to the inherent local rigidity of objects, vertices that are spatially proximate in the initial frame typically exhibit similar motion trajectories. However, this spatial proximity assumption breaks down when considering vertices that belong to different structural components of the object, as illustrated in Fig. 3 (compare the trajectories of the blue and orange crosses). To address this limitation, we exploit the mesh’s topological structure and design a topology-aware attention mechanism by incorporating connectivity information into vertex features, which serve as more reliable trajectory identifiers. Specifically, we first construct an adjacency matrix Adj from the face information F of the input mesh. This adjacency matrix then acts as an attention mask in a self-attention layer, enabling each vertex to aggregate information from its connected neighbors:

$$\overline{V}_0 = \text{Softmax} \left(\frac{\widehat{V}_0 \cdot \widehat{V}_0^T \odot Adj}{\sqrt{d_k}} \right) \widehat{V}_0 + \widehat{V}_0, \quad (2)$$

where d_k denotes the channel dimension of the projected space, and \odot denotes Hadamard product.

Building upon our observation that vertices with similar neighborhood-enhanced features tend to exhibit similar

motion patterns, we leverage this correlation by applying Farthest Point Sampling (FPS) [55] to the topology-aware vertex features \overline{V}_0 to obtain \overline{V}_0^n . The same sampling indices are used to sample \widehat{V}_T , resulting in \widehat{V}_T^n . Subsequently, we employ a cross-attention mechanism where \overline{V}_0^n serves as the query, while the complete set of topology-aware features \overline{V}_0 acts as both key and value. The resulting attention map is then used to project \widehat{V}_T^n , as formulated below:

$$\overline{V}_0^n = \text{Softmax} \left(\frac{\overline{V}_0^n \cdot \overline{V}_0^T}{\sqrt{d_k}} \right) \overline{V}_0 + \overline{V}_0^n, \quad (3)$$

$$\widehat{V}_T^n = \text{Softmax} \left(\frac{\overline{V}_0^n \cdot \overline{V}_0^T}{\sqrt{d_k}} \right) \widehat{V}_T + \widehat{V}_T^n. \quad (4)$$

We employ a stack of cross-attention layers, as described above, and iteratively apply Equ. (3), (4) to progressively gather global information into \overline{V}_0^n and \widehat{V}_T^n .

KL Regularization. After encoding, a dynamic mesh sequence is compressed to a pair of latents $\{\overline{V}_0^n, \widehat{V}_T^n\}$. Since we focus on the task of animating existing meshes, we only have to model the distribution of the relative trajectories \widehat{V}_T^n . In accordance with the latent diffusion framework [61, 96], our model employs KL-regularization in the latent space for the purpose of modulating feature diversity. Specifically, we linearly project \widehat{V}_T^n by two fully-connected layers to predict the mean μ_T^n and standard variance σ_T^n of the distribution of \widehat{V}_T^n . Then, the VAE latent is sampled as:

$$Z_T^n = \mu_T^n + \sigma_T^n \cdot \epsilon, \quad (5)$$

where $\epsilon \sim \mathcal{N}(0, 1)$. The KL regularization is as follows:

$$L_{kl} = \frac{1}{2n \cdot C} \sum_{i=1}^n \sum_{j=1}^C ((\mu_T^{i,j})^2 + (\sigma_T^{i,j})^2 - \log(\sigma_T^{i,j})^2), \quad (6)$$

where C denotes the number of channels.

DyMeshVAE Decoder. Similarly, to further exploit the correlation between motion patterns of vertices with similar neighborhood-enhanced features, we perform self-attention on the sampled topology-aware features \overline{V}_0^n . The resulting attention map is then used to project Z_T^n , capturing the motion relationships among vertices with similar local structural characteristics, which is formulated as:

$$\overline{V}_0^n = \text{Softmax} \left(\frac{\overline{V}_0^n \cdot \overline{V}_0^T}{\sqrt{d_k}} \right) \overline{V}_0^n + \overline{V}_0^n, \quad (7)$$

$$Z_T^n = \text{Softmax} \left(\frac{\overline{V}_0^n \cdot \overline{V}_0^T}{\sqrt{d_k}} \right) Z_T^n + Z_T^n. \quad (8)$$

To enhance the model capacity, we stack K identical blocks in a cascaded manner, where each block comprises

the aforementioned attention mechanism. Next, we perform cross attention where the encoded vertex features of the initial mesh \bar{V}_0 serve as the query, while \bar{V}_0^n and Z_T^n act as the key and value respectively, which can be formulated as:

$$V_T^{rec} = \text{Softmax} \left(\frac{\bar{V}_0 \cdot \bar{V}_0^{nT}}{\sqrt{d_k}} \right) Z_T^n. \quad (9)$$

The reconstructed trajectories are obtained by projecting V_T^{rec} through a fully connected layer to map the channel dimension to $T \times 3$. The reconstruction loss is the mean square error between V_T^{rec} and V_T :

$$L_{rec} = \frac{1}{N} \sum_{i=1}^N \|(V_T^{rec})^i - V_T^i\|_2^2. \quad (10)$$

The overall loss is a weighted combination of L_{rec} and L_{kl} : $L_{dvae} = L_{rec} + \gamma \cdot L_{kl}$, where we set $\gamma = 0.001$ as default. Finally, we transpose the obtained relative trajectories V_T^{rec} and add the initial vertex positions V_0 to get the reconstructed dynamic mesh.

3.2. Shape-Guided Text-to-Trajectory Model

For text-driven mesh animation, our model learns to estimate the posterior distribution of relative trajectories conditioned on both the initial mesh and textual prompts. As shown in Fig. 4, we get inspirations from [17, 53, 90] and propose the Shape-Guided Text-to-Trajectory Model, which builds upon the MMDiT [17] architecture.

As described in Sec. 3.1, each dynamic mesh sequence D is encoded into a pair of latent representations $\{\bar{V}_0^n, Z_T^n\}$. To facilitate effective multi-modal learning, we first normalize these latents using their respective global statistics (mean μ_0, μ_T and standard deviation σ_0, σ_T) to eliminate numerical disparities. The normalized features are then concatenated along the channel dimension to form a comprehensive trajectory embedding. Meanwhile, text prompts are encoded through a pre-trained CLIP [57] text encoder to obtain text embeddings. Following [17, 90], we apply separate Adaptive Layer Normalization (AdaLN) parameters conditioned on timestep t to modulate the trajectory and text embeddings, respectively. These rescaled features are subsequently concatenated for self-attention computation. This modality-specific normalization strategy effectively bridges the distributional gap between varied representations, enabling more robust multi-modal learning. After processing through multiple such attention blocks, we decompose the output features back into and restore their original scales using the preserved statistics μ_0, μ_T and σ_0, σ_T .

3.3. Diffusion Pipeline

Training. Following the diffusion paradigm in Rectified Flow (RF) [42], we aim to minimize the mean square error between the predicted and ground truth flow. To be noted,

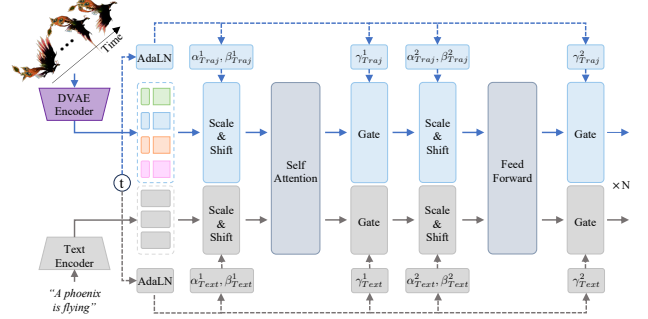


Figure 4. The architecture of the Shape-Guided Text-to-Trajectory Model. DVAE stands for the proposed DyMeshVAE.

since \bar{V}_0^n is settled, our model focuses solely on learning the distribution of the relative trajectory features Z_T^n . Specifically, we apply the diffusion process to Z_T^n to obtain its noisy version \tilde{Z}_T^n , as formulated below:

$$\tilde{Z}_T^n = (1 - t)Z_T^n + t\epsilon, \quad t = 1 - \frac{1}{\tan(\frac{\pi}{2}u) + 1}, \quad (11)$$

where $\epsilon \sim \mathcal{N}(0, 1)$ and $u \sim \mathcal{U}(0, 1)$. Then, the optimization process can be formulated as:

$$L_{rf} = \mathbb{E}_{\bar{V}_0^n, C_{text}} \left\| v_\theta(\tilde{Z}_T^n; \bar{V}_0^n, C_{text}) - u_T^n \right\|_2^2, \quad (12)$$

where v_θ represents the backbone in Sec. 3.2. And $u_T^n = Z_T^n - \epsilon$, where ϵ is the sampled noise.

Inference. We implement the sampling procedure for Rectified Flow following the flow-based ODE formulation. The sampling process can be formulated as:

$$\frac{d\tilde{Z}_T^n}{dt} = f_\theta(\tilde{Z}_T^n(t), t), \quad Z_T^n(1) \sim \mathcal{N}(0, 1), \quad (13)$$

where f_θ denotes the learned velocity field. For classifier-free guidance (CFG), we modify the velocity field as:

$$f_\theta^{cfg} = f_\theta^{uncond} + \gamma \cdot (f_\theta^{cond} - f_\theta^{uncond}), \quad (14)$$

where we set the scale $\gamma = 3.0$ as default.

After sampling, we feed both the sampled features Z_T^n and the topology-aware vertex features \bar{V}_0^n into the decoder to generate the relative trajectories. The final animated mesh sequence is generated by applying the decoded vertex displacements to the vertices of the given mesh.

4. Experiments

4.1. Settings

Data Curation. As mentioned in Sec. 1, learning generalizable 4D generative models requires large-scale, high-quality 4D datasets. Recent works [26, 31, 36, 59, 97] have

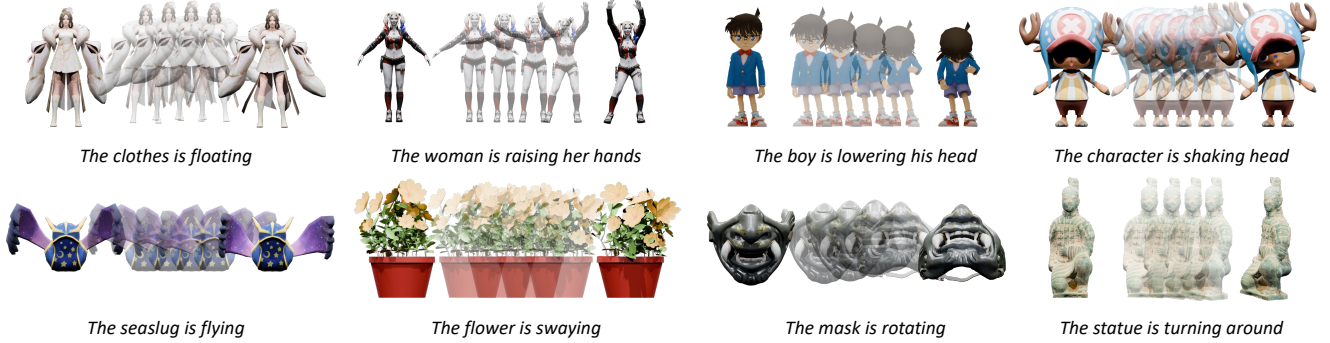


Figure 5. Animation examples of AnimateAnyMesh. Our model demonstrates the capability to generate high-quality and semantically plausible mesh animations for arbitrary input meshes based on text prompts. Best viewed when zoomed in.

All Animations	Frame number	Vertex number		
		$\leq 50,000$	≤ 8192	≤ 4096
66,209	16	2,604,053	1,237,824	330,658
	32	1,572,944	848,239	182,172

Table 1. Statistical information of our DyMesh dataset.

attempted to construct such datasets by filtering and rendering multi-view video sequences from Objaverse [15, 16]. However, as discussed in Sec. 1, the lack of genuine 4D data in these datasets makes single-stage 4D generation impractical. To address this limitation, we convert various 4D assets into sequential dynamic meshes and propose DyMesh Dataset. Our dataset integrates multiple sources, including dynamic objects/scenes from Objaverse ($\sim 55k$ instances, counting multiple animations per object separately), dynamic human SMPL [44] models from AMASS [47] ($\sim 8k$ instances), and dynamic human and animal models from DeformingThings4D [33] ($\sim 2k$ instances), yielding a total of 66k complete animations. The final dataset statistics after sequence segmentation and filtering are presented in Tab. 1. For each sequence, we render front-view videos at 256×256 resolution using Blender and generate captions using the Qwen-2.5-VL [5] model. Detailed information about our dataset can be found in Sec. A.1 of the appendix.

Implementation Details. We conduct extensive experiments to evaluate our proposed framework by training DyMeshVAE and the corresponding rectified flow model on four subsets of our dataset, varying the maximum vertex count (4096/8192) and sequence length (16/32 frames). For comparison with existing methods, we adopt the model variant trained on sequences with 4096 vertices and 16 frames. In this configuration, our DyMeshVAE architecture comprises an encoder with 8 cross-attention layers and a decoder with 8 self-attention layers, where the encoder samples 512 latent tokens via Farthest Point Sampling (FPS). The model is optimized using the Adam optimizer [28] with a learning rate of $1e-4$ for 1000 epochs on 8 H20 GPUs. For the rectified flow model, we leverage a pre-trained CLIP [57] ViT-L/14 as the text encoder with a maximum sequence length of 77 tokens. The training

is performed on 32 H20 GPUs using the Adam optimizer with a learning rate of $2e-4$ for 1000 epochs. At inference time, our framework demonstrates strong generalization capability by adapting to varying mesh resolutions through dynamic FPS sampling density adjustment. Specifically, processing a mesh with 20K vertices takes approximately 6 seconds on a single Nvidia A800 GPU for arbitrary mesh animation. We provide comprehensive implementation details in Sec. A.2 of the appendix.

Evaluation Metrics. To quantitatively evaluate DyMeshVAE’s reconstruction quality, we compute the frame-wise average L2 distance between the reconstructed and ground-truth vertex trajectories. For evaluating mesh animations without references, we adopt the evaluation protocol adopted in Animate3D [26]. Specifically, we render the animated meshes from specified viewpoints and leverage the view-aligned renderings as references to compute multiple perceptual metrics from VBench [24], including I2V Subject Similarity, Motion Smoothness, and Aesthetic Quality (denoted as I2V, M.sm, and Aest.Q respectively). Note that we exclude the Dynamic Degree metric from VBench due to the presence of shape degeneration in some comparative methods under drastic motions. To further assess the perceptual quality, we conduct a user study with 20 randomly recruited participants. Each participant is asked to rate the generated animations on a 5-point Likert scale (1-lowest, 5-highest) considering three aspects: text-to-motion alignment (User. Ta), motion naturalness (User. Mn), and shape preservation (User. Sp). The final scores are computed by averaging all ratings for each metric. Additionally, we report the computational efficiency by measuring the average time consumption for generating a single mesh animation. Comprehensive descriptions of evaluation metrics are provided in Sec. A.3 of the appendix.

4.2. Results

Our proposed AnimateAnyMesh enables text-driven mesh animation for generic objects, generating dynamic mesh sequences that demonstrate sophisticated performance in se-

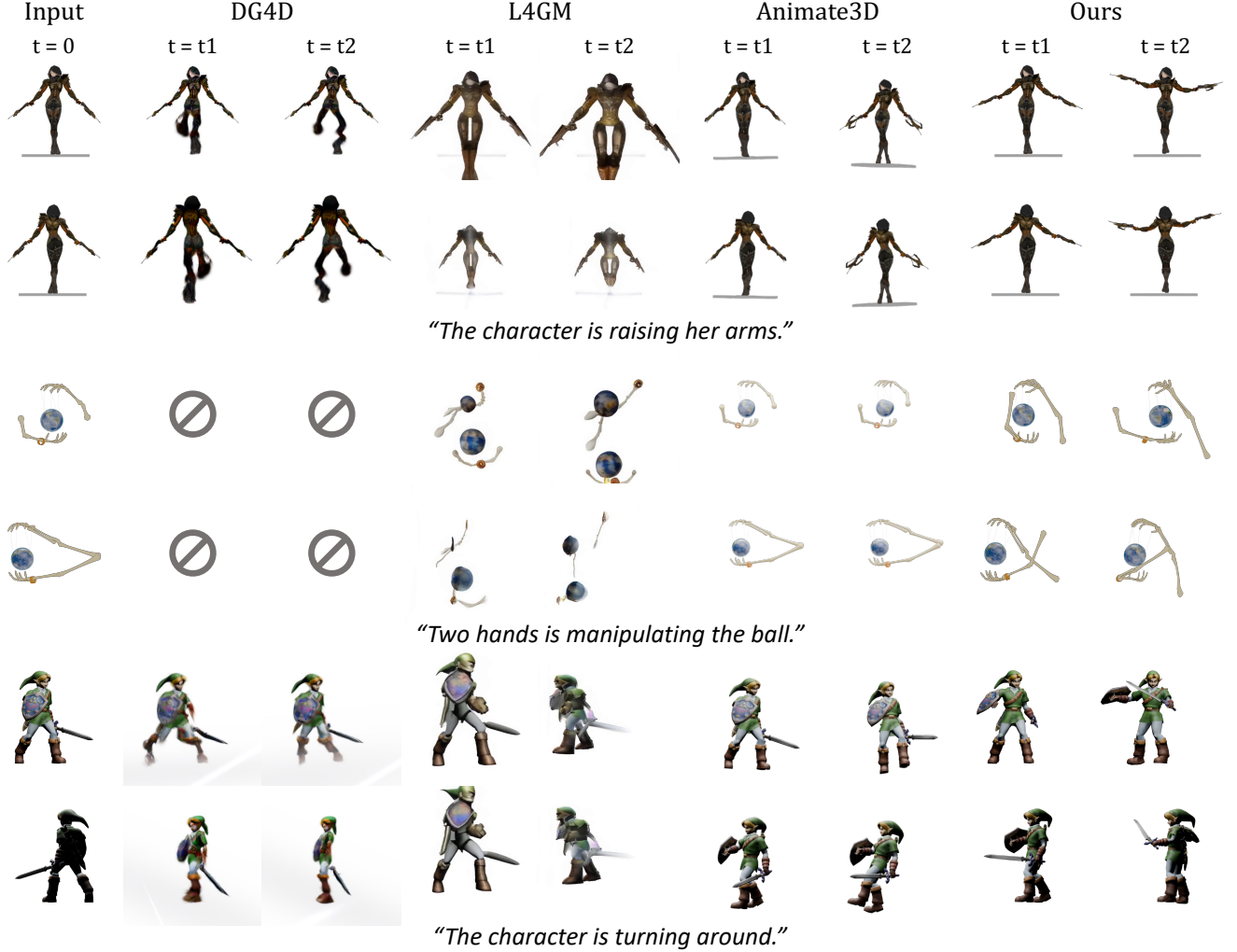


Figure 6. Qualitative comparison with state-of-the-art methods. \otimes represents the failure cases. Best viewed when zoomed in.

mantic alignment, motion naturalness, and geometric consistency. Representative examples are illustrated in Fig. 5. Additional examples can be found in Fig. 11 of the appendix.

4.3. Comparison

We evaluate our approach through comprehensive comparisons with three state-of-the-art methods, each employing distinct paradigms for multi-stage mesh animation. Specifically, we compare against DG4D [58], L4GM [59], and Animate3D [26], all comparative evaluations are obtained using their official code and settings. We adopt DynamiCrafter [86] as the video generator for L4GM.

Qualitative Comparison. As illustrated in Fig. 6, existing methods exhibit various limitations in text-driven mesh animation. DG4D [58], which distills knowledge from video diffusion models [6] to optimize 3D Gaussians through SDS-based optimization, demonstrates limited shape preservation capability and suffers from severe

object drifting, as evidenced in rows 3_{rd} and 4_{th} of Fig. 6. L4GM [59], which aims to reconstruct 4D representations directly from single-view videos, is inherently limited by the performance of the video generator, particularly when handling rendered sequences of isolated 3D objects without natural backgrounds. While Animate3D [26] proposes a theoretically sound approach by converting initial meshes to Gaussian Splatting representations and optimizing motion trajectories through multi-view dynamic video generation with ARAP [68] regularization, its multi-stage pipeline suffers from error accumulation, leading to degraded performance in complex scenarios. In contrast, our method achieves superior results through direct vertex trajectory prediction, demonstrating two key advantages: (1). efficient mesh animation through a feed-forward architecture, and (2). enhanced preservation of local geometric details while generating prompt-aligned realistic motions. Please refer to Sec. C of the appendix for more visualization results.

Quantitative Comparison. For quantitative comparison,

Method	VBench			User study			Time ↓
	I2V ↑	M.Sm ↑	Aest.Q ↑	User.Ta ↑	User.Mn ↑	User.Sp ↑	
DG4D [58]	0.811	0.926	0.476	2.130	2.460	2.755	~10min
L4GM [59]	0.844	0.992	0.464	2.885	2.865	2.835	~30s
Animate3D [26]	0.936	0.992	0.526	2.850	3.195	3.405	~14min
Ours	0.954	0.995	0.539	4.505	4.700	4.790	~6s

Table 2. Quantitative comparison with state-of-the-art methods. All the evaluation is performed on a Nvidia A800 GPU.

<i>Adj</i>	PE_0	PE_t	<i>Sep Attn</i>	<i>Emb FPS</i>	Rec Error ↓
×	✓	✓	✓	✓	0.500
✓	×	✓	✓	✓	0.443
✓	✓	×	✓	✓	0.441
✓	✓	✓	×	✓	0.478
✓	✓	✓	✓	×	0.291
✓	✓	✓	✓	✓	0.223

Table 3. Ablation studies of technical components of DyMeshVAE. Rec Error denotes the average L2 sum error per-instance.

we curate a test set comprising 10 randomly selected objects across multiple categories (human, animals, weapons, etc.). Using identical text prompts as conditions, we compute the mentioned metrics (detailed in Sec. 4.1) across all baseline methods. The comparative results are presented in Tab. 2.

As shown in Tab. 2, our approach achieves superior performance across all VBench metrics (I2V, M.Sm, Aest.Q), indicating its effectiveness in both shape preservation and temporal motion coherence. The user study results further validate our method’s advantages, showing substantial improvements in both text-motion alignment accuracy and motion naturalness. These quantitative evaluations align well with our qualitative observations. Moreover, AnimateAnyMesh significantly reduces the computational overhead during inference compared to existing methods, making it particularly promising for practical applications.

4.4. Ablation Studies

We conduct comprehensive ablation studies to validate the effectiveness of key components in our DyMeshVAE architecture. Specifically, we examine the incorporation of mesh adjacency information (*Adj*), positional encoding for both initial frame vertex features (PE_0) and relative trajectory features (PE_T), separate attention computation using initial frame features instead of complete trajectory features (*Sep Attn*), and Farthest Point Sampling on neighborhood-enhanced vertex features (*Emb FPS*). Further ablation studies are provided in Sec. D of the appendix.

As shown in Tab. 3, each architectural component in our framework plays a crucial role in ensuring high-quality mesh animation. The incorporation of mesh connectivity information (*Adj*) is essential for distinguishing vertices from different semantic regions. Without this topology information, vertices that are spatially adjacent but semantically distinct (e.g., hand and waist regions in Fig. 7) tend to be encoded similarly, leading to undesired adhesion artifacts during decoding. The mesh topology enables better feature discrimination through distinct neigh-

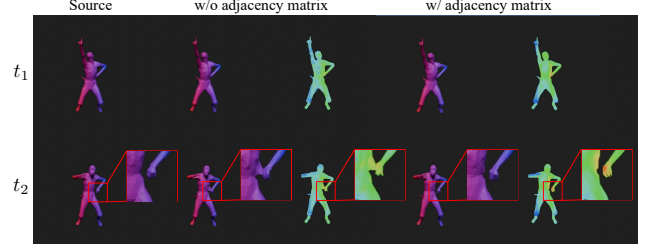


Figure 7. Ablation study on mesh adjacency information. Visualization shows the cosine similarity between feature embeddings of a query vertex (in hand region) and all mesh vertices, with similarity intensity ranging from blue (low) to red (high).

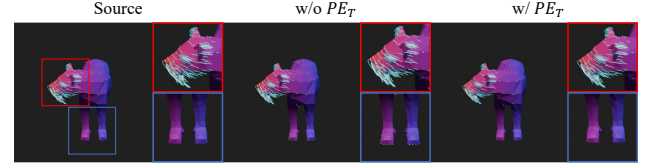


Figure 8. Ablation study on trajectory positional encoding (PE_T). Vertex trajectories shown in teal. Please zoom in for a better view.

borhood patterns, resulting in more accurate trajectory predictions. Positional encoding schemes for both initial frame features (PE_0) and trajectory features (PE_T) project the low-dimensional trajectory information into a higher-dimensional space. This enhancement not only improves the discriminative power of queries but also ensures robust trajectory reconstruction, even under VAE sampling and noise perturbation. As demonstrated in Fig. 8, the absence of trajectory feature positional encoding leads to unstable motion patterns. Furthermore, our design choice of computing attention maps using only initial frame features (*Sep Attn*) effectively leverages the prior knowledge that vertices with similar neighborhood patterns tend to share similar trajectories. The implementation of Farthest Point Sampling (*Emb FPS*) on topology-enhanced initial vertex features, rather than raw vertex coordinates, proves more effective in selecting representative trajectories, thereby enabling better compression and reconstruction capabilities.

5. Conclusion

In this paper, we introduce AnimateAnyMesh, the first feed-forward framework for universal text-driven mesh animation. At its core lies DyMeshVAE, a novel VAE architecture specifically designed for dynamic mesh sequences. By decomposing dynamic meshes into initial states and relative trajectories while incorporating mesh topology information, DyMeshVAE achieves efficient variable-length compression and reconstruction of dynamic mesh sequences through carefully designed attention mechanisms. To enable modeling the motion distribution in the compressed latent space, we propose the Shape-Guided Text-to-Trajectory Model, which adopts a Rectified Flow-based training strategy that effectively learns the distribution of trajectories conditioned on text

and initial mesh features. To facilitate the training of our framework, we construct DyMesh Dataset, a large-scale dynamic mesh dataset containing over 4M dynamic mesh sequences. Extensive experiments demonstrate that AnimateAnyMesh can generate high-quality language-driven animations for arbitrary meshes within seconds, marking a significant advancement in feed-forward 4D generation.

References

- [1] Advanced Computing Center for the Arts and Design. AC-CAD MoCap Dataset. 3
- [2] Ijaz Akhter and Michael J. Black. Pose-conditioned joint angle limits for 3D human pose reconstruction. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) 2015*, 2015.
- [3] Andreas Aristidou, Ariel Shamir, and Yiorgos Chrysanthou. Digital dance ethnography: Organizing large dance collections. *J. Comput. Cult. Herit.*, 12(4), 2019. 3
- [4] Sherwin Bahmani, Ivan Skorokhodov, Victor Rong, Gordon Wetzstein, Leonidas Guibas, Peter Wonka, Sergey Tulyakov, Jeong Joon Park, Andrea Tagliasacchi, and David B Lindell. 4d-fy: Text-to-4d generation using hybrid score distillation sampling. *arXiv preprint arXiv:2311.17984*, 2023. 1, 3
- [5] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 6, 13
- [6] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 7
- [7] Federica Bogo, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Dynamic FAUST: Registering human bodies in motion. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3
- [8] Samarth Brahmabhatt, Cusuh Ham, Charles C. Kemp, and James Hays. ContactDB: Analyzing and predicting grasp contact via thermal imaging. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3
- [9] Wei Cao, Chang Luo, Biao Zhang, Matthias Nießner, and Jiapeng Tang. Motion2vecsets: 4d latent vector set diffusion for non-rigid shape reconstruction and tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20496–20506, 2024. 3
- [10] Carnegie Mellon University. CMU MoCap Dataset. 3
- [11] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 13
- [12] Cersense. Zeroscope text-to-video model. https://huggingface.co/cersense/zeroscope_v2_576w, 2023. Accessed: 2023-10-31. 2, 3
- [13] Anargyros Chatzitofis, Leonidas Saroglou, Prodromos Boutis, Petros Drakoulis, Nikolaos Zioulis, Shishir Subramanyam, Bart Kevelham, Caecilia Charbonnier, Pablo Cesar, Dimitrios Zarpalas, et al. Human4d: A human-centric multimodal dataset for motions and immersive media. *IEEE Access*, 8:176241–176262, 2020. 3
- [14] Zilong Chen, Feng Wang, and Huaping Liu. Text-to-3d using gaussian splatting. *arXiv preprint arXiv:2309.16585*, 2023. 2
- [15] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. 2, 6, 13
- [16] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of 10m+ 3d objects. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 6
- [17] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 5
- [18] Nima Ghorbani and Michael J. Black. SOMA: Solving optical marker-based mocap automatically. In *Proc. International Conference on Computer Vision (ICCV)*, pages 11117–11126, 2021. 3
- [19] Saeed Ghorbani, Kimia Mahdavian, Anne Thaler, Konrad Kording, Douglas James Cook, Gunnar Blohm, and Nikolaus F. Troje. MoVi: A large multipurpose motion and video dataset. *arXiv preprint arXiv: 2003.01888*, 2020. 3
- [20] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 2, 3
- [21] Fabian Helm, Nikolaus Troje, Mathias Reiser, and Jörn Munzert. Bewegungsanalyse geträuschter und nicht-geträuschter 7m-wärfe im handball. 2015. 3
- [22] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint arXiv:2311.04400*, 2023. 1, 3
- [23] Ludovic Hoyet, Kenneth Ryall, Rachel McDonnell, and Carol O’Sullivan. Sleight of hand: Perception of finger motion from reduced marker sets. In *Proceedings of the ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games*, pages 79–86, 2012. 3
- [24] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 6, 13, 16

- [25] Yanqin Jiang, Li Zhang, Jin Gao, Weimin Hu, and Yao Yao. Consistent4d: Consistent 360 $\{\deg\}$ dynamic object generation from monocular video. *arXiv preprint arXiv:2311.02848*, 2023. 1, 3
- [26] Yanqin Jiang, Chaohui Yu, Chenjie Cao, Fan Wang, Weiming Hu, and Jin Gao. Animate3d: Animating any 3d model with multi-view video diffusion. *arXiv preprint arXiv:2407.11398*, 2024. 1, 3, 5, 6, 7, 8
- [27] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 2, 3
- [28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [29] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2
- [30] Franziska Krebs, Andre Meixner, Isabel Patzer, and Tamim Asfour. The KIT bimanual manipulation dataset. In *IEEE/RAS International Conference on Humanoid Robots (Humanoids)*, pages 499–506, 2021. 3
- [31] Bing Li, Cheng Zheng, Wenxuan Zhu, Jinjie Mai, Biao Zhang, Peter Wonka, and Bernard Ghanem. Vivid-zoo: Multi-view video generation with diffusion model. *Advances in Neural Information Processing Systems*, 37: 62189–62222, 2025. 3, 5
- [32] Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model. *arXiv preprint arXiv:2311.06214*, 2023. 1, 3
- [33] Yang Li, Hikari Takehara, Takafumi Taketomi, Bo Zheng, and Matthias Nießner. 4dcomplete: Non-rigid motion estimation beyond the observable surface. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12706–12716, 2021. 3, 6, 13
- [34] Yunzhi Li, Vimal Mollyn, Kuang Yuan, and Patrick Carrington. Wheelposer: Sparse-imu based body pose estimation for wheelchair users. In *Proceedings of the 26th International ACM SIGACCESS Conference on Computers and Accessibility*, pages 1–17, 2024. 3
- [35] Zhen Li, Zuo-Liang Zhu, Ling-Hao Han, Qibin Hou, Chun-Le Guo, and Ming-Ming Cheng. Amt: All-pairs multi-field transforms for efficient frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9801–9810, 2023. 13
- [36] Hanwen Liang, Yuyang Yin, Dejia Xu, Hanxue Liang, Zhangyang Wang, Konstantinos N Plataniotis, Yao Zhao, and Yunchao Wei. Diffusion4d: Fast spatial-temporal consistent 4d generation via video diffusion models. *arXiv preprint arXiv:2405.16645*, 2024. 1, 3, 5
- [37] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 300–309, 2023. 2
- [38] Huan Ling, Seung Wook Kim, Antonio Torralba, Sanja Fidler, and Karsten Kreis. Align your gaussians: Text-to-4d with dynamic 3d gaussians and composed diffusion models. *arXiv preprint arXiv:2312.13763*, 2023. 3
- [39] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *Advances in Neural Information Processing Systems*, 36:22226–22246, 2023. 1, 3
- [40] Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen, Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3d objects with consistent multi-view generation and 3d diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10072–10083, 2024. 1, 3
- [41] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9298–9309, 2023. 2
- [42] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022. 2, 3, 5
- [43] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 2
- [44] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 851–866, 2023. 3, 6, 13
- [45] Matthew M. Loper, Naureen Mahmood, and Michael J. Black. MoSh: Motion and shape capture from sparse markers. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 33(6):220:1–220:13, 2014. 3
- [46] Eyes JAPAN Co. Ltd. Eyes Japan MoCap Dataset.
- [47] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS: Archive of motion capture as surface shapes. In *International Conference on Computer Vision*, pages 5442–5451, 2019. 6, 13
- [48] Christian Mandery, Ömer Terlemez, Martin Do, Nikolaus Vahrenkamp, and Tamim Asfour. The KIT whole-body human motion database. In *International Conference on Advanced Robotics (ICAR)*, pages 329–336, 2015.
- [49] Christian Mandery, Ömer Terlemez, Martin Do, Nikolaus Vahrenkamp, and Tamim Asfour. Unifying representations and large-scale whole-body motion databases for studying human motion. *IEEE Transactions on Robotics*, 32(4):796–809, 2016. 3
- [50] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2, 3, 4

- [51] M. Müller, T. Röder, M. Clausen, B. Eberhardt, B. Krüger, and A. Weber. Documentation mocap database hdm05. Technical Report CG-2007-2, Universität Bonn, 2007. 3
- [52] Zijie Pan, Zeyu Yang, Xiatian Zhu, and Li Zhang. Fast dynamic 3d object generation from a single-view video. *arXiv preprint arXiv:2401.08742*, 2024. 3
- [53] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 5
- [54] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 1, 2
- [55] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 4
- [56] Lingteng Qiu, Guanying Chen, Xiaodong Gu, Qi Zuo, Mutian Xu, Yushuang Wu, Weihao Yuan, Zilong Dong, Liefeng Bo, and Xiaoguang Han. Richdreamer: A generalizable normal-depth diffusion model for detail richness in text-to-3d. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9914–9925, 2024. 2
- [57] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. Pmlr, 2021. 2, 5, 6
- [58] Jiawei Ren, Liang Pan, Jiaxiang Tang, Chi Zhang, Ang Cao, Gang Zeng, and Ziwei Liu. Dreamgaussian4d: Generative 4d gaussian splatting. *arXiv preprint arXiv:2312.17142*, 2023. 1, 3, 7, 8
- [59] Jiawei Ren, Cheng Xie, Ashkan Mirzaei, Karsten Kreis, Ziwei Liu, Antonio Torralba, Sanja Fidler, Seung Wook Kim, Huan Ling, et al. L4gm: Large 4d gaussian reconstruction model. *Advances in Neural Information Processing Systems*, 37:56828–56858, 2025. 3, 5, 7, 8
- [60] Zeping Ren, Shaoli Huang, and Xiu Li. Realistic human motion generation with cross-diffusion models. In *European Conference on Computer Vision*, pages 345–362. Springer, 2024. 3
- [61] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 4
- [62] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022. 2
- [63] Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen, Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base model. *arXiv preprint arXiv:2310.15110*, 2023. 2
- [64] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023. 2, 3
- [65] L. Sigal, A. Balan, and M. J. Black. HumanEva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. *International Journal of Computer Vision*, 87(1):4–27, 2010. 3
- [66] Uriel Singer, Shelly Sheynin, Adam Polyak, Oron Ashual, Iurii Makarov, Filippos Kokkinos, Naman Goyal, Andrea Vedaldi, Devi Parikh, Justin Johnson, et al. Text-to-4d dynamic scene generation. *arXiv preprint arXiv:2301.11280*, 2023. 3
- [67] Sketchfab. Sketchfab. <https://sketchfab.com/>, 2024. Accessed: 2024-05-21. 2, 13, 15
- [68] Olga Sorkine and Marc Alexa. As-rigid-as-possible surface modeling. In *Symposium on Geometry processing*, pages 109–116. Citeseer, 2007. 7
- [69] Jingxiang Sun, Bo Zhang, Ruizhi Shao, Lizhen Wang, Wen Liu, Zhenda Xie, and Yebin Liu. Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior. *arXiv preprint arXiv:2310.16818*, 2023. 2
- [70] Omid Taheri, Nima Ghorbani, Michael J. Black, and Dimitrios Tzionas. GRAB: A dataset of whole-body human grasping of objects. In *European Conference on Computer Vision (ECCV)*, 2020. 3
- [71] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023. 2
- [72] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. *arXiv preprint arXiv:2303.14184*, 2023. 2
- [73] Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm: Large multi-view gaussian model for high-resolution 3d content creation. In *European Conference on Computer Vision*, pages 1–18. Springer, 2024. 1, 3
- [74] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano. Human motion diffusion model. *arXiv preprint arXiv:2209.14916*, 2022. 3
- [75] Shashank Tripathi, Lea Müller, Chun-Hao P. Huang, Taheri Omid, Michael J. Black, and Dimitrios Tzionas. 3D human pose estimation via intuitive physics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [76] Nikolaus F. Troje. Decomposing biological motion: A framework for analysis and synthesis of human gait patterns. *Journal of Vision*, 2(5):2–2, 2002.
- [77] Matt Trumble, Andrew Gilbert, Charles Malleson, Adrian Hilton, and John Collomosse. Total Capture: 3d human pose estimation fusing video and inertial sensors. In *2017 British Machine Vision Conference (BMVC)*, 2017.
- [78] Simon Fraser University and National University of Singapore. SFU Motion Capture Database. 3

- [79] Lukas Uzolas, Elmar Eisemann, and Petr Kellnhofer. Motiondreamer: Zero-shot 3d mesh animation from video diffusion models. *arXiv preprint arXiv:2405.20155*, 2024. 3
- [80] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12619–12629, 2023. 2
- [81] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023. 2, 3
- [82] Peng Wang and Yichun Shi. Imagedream: Image-prompt multi-view diffusion for 3d generation. *arXiv preprint arXiv:2312.02201*, 2023. 2
- [83] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *Advances in Neural Information Processing Systems*, 36:8406–8441, 2023. 2
- [84] Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan, Valentin Deschaintre, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Meshlrn: Large reconstruction model for high-quality meshes. *arXiv preprint arXiv:2404.12385*, 2024. 1, 2, 3
- [85] Zijie Wu, Chaohui Yu, Yanqin Jiang, Chenjie Cao, Fan Wang, and Xiang Bai. Sc4d: Sparse-controlled video-to-4d generation and motion transfer. In *European Conference on Computer Vision*, pages 361–379. Springer, 2024. 1, 3
- [86] Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Wangbo Yu, Hanyuan Liu, Gongye Liu, Xintao Wang, Ying Shan, and Tien-Tsin Wong. Dynamicrafter: Animating open-domain images with video diffusion priors. In *European Conference on Computer Vision*, pages 399–417. Springer, 2024. 7
- [87] Jiale Xu, Xintao Wang, Weihao Cheng, Yan-Pei Cao, Ying Shan, Xiaohu Qie, and Shenghua Gao. Dream3d: Zero-shot text-to-3d synthesis using 3d shape prior and text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20908–20918, 2023. 2
- [88] Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. Instantmesh: Efficient 3d mesh generation from a single image with sparse-view large reconstruction models. *arXiv preprint arXiv:2404.07191*, 2024. 1, 2, 3
- [89] Yinghao Xu, Zifan Shi, Wang Yifan, Hansheng Chen, Ceyuan Yang, Sida Peng, Yujun Shen, and Gordon Wetstein. Grm: Large gaussian reconstruction model for efficient 3d reconstruction and generation. In *European Conference on Computer Vision*, pages 1–20. Springer, 2024. 1, 3
- [90] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 5
- [91] Taoran Yi, Jiemin Fang, Guanjun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and Xinggang Wang. Gaussiandreamer: Fast generation from text to 3d gaussian splatting with point cloud priors. *arXiv preprint arXiv:2310.08529*, 2023. 2
- [92] Yuyang Yin, Dejia Xu, Zhangyang Wang, Yao Zhao, and Yunchao Wei. 4dgen: Grounded 4d content generation with spatial-temporal consistency. *arXiv preprint arXiv:2312.17225*, 2023. 3
- [93] Chaohui Yu, Qiang Zhou, Jingliang Li, Zhe Zhang, Zhibin Wang, and Fan Wang. Points-to-3d: Bridging the gap between sparse points and shape-controllable text-to-3d generation. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 6841–6850, 2023. 2
- [94] Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. Physdiff: Physics-guided human motion diffusion model. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 16010–16021, 2023. 3
- [95] Yifei Zeng, Yanqin Jiang, Siyu Zhu, Yuanxun Lu, Youtian Lin, Hao Zhu, Weiming Hu, Xun Cao, and Yao Yao. Stag4d: Spatial-temporal anchored generative 4d gaussians. In *European Conference on Computer Vision*, pages 163–179. Springer, 2024. 3
- [96] Biao Zhang, Jiapeng Tang, Matthias Niessner, and Peter Wonka. 3dshape2vecset: A 3d shape representation for neural fields and generative diffusion models. *ACM Transactions On Graphics (TOG)*, 42(4):1–16, 2023. 4
- [97] Haiyu Zhang, Xinyuan Chen, Yaohui Wang, Xihui Liu, Yunhong Wang, and Yu Qiao. 4diffusion: Multi-view video diffusion model for 4d generation. *Advances in Neural Information Processing Systems*, 37:15272–15295, 2025. 1, 3, 5
- [98] Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang Xu. Gs-lrm: Large reconstruction model for 3d gaussian splatting. In *European Conference on Computer Vision*, pages 1–19. Springer, 2024. 1, 3
- [99] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. *IEEE transactions on pattern analysis and machine intelligence*, 46(6):4115–4128, 2024. 3
- [100] Yisu Zhang, Chaohui Yu, Fan Wang, and Jianke Zhu. Lpm: Efficient 3d content creation from single image by large-scale partial 3d modeling. *IEEE Transactions on Circuits and Systems for Video Technology*, 2025. 1, 3
- [101] Yuyang Zhao, Zhiwen Yan, Enze Xie, Lanqing Hong, Zhenguo Li, and Gim Hee Lee. Animate124: Animating one image to 4d dynamic scene. *arXiv preprint arXiv:2311.14603*, 2023. 1, 3

A. Settings

In this section, we elaborate on the technical details of AnimateAnyMesh. The chapter is structured into three main components: data curation [A.1](#), implementation details [A.2](#), and evaluation metric [4.1](#).

A.1. Data Curation

As described in the main paper, our curated 4D data is sourced from three primary sources: Objaverse [15, 67], AMASS [47], and DT4D [33]. Initially, we filter and extract all .glb files containing animation sequences from Objaverse. Using Blender’s Python API (bpy), we convert each animation into a mesh sequence. Animations with fewer than 16 frames are discarded, and each sequence is capped at 200 frames. Post-conversion, each animation is encoded into a .bin file comprising $D \subset \{F \in \mathbb{R}^{M \times 3}, V \in \mathbb{R}^{T \times N \times 3}\}$, where M denotes the face count, T represents the temporal length of the dynamic mesh sequence, and N indicates the vertex count. Similarly, we develop scripts to convert SMPL [44] models from AMASS and .anime files from DT4D into the identical .bin format.

Subsequently, we traverse all stored files, implementing vertex merging operations for duplicate vertices while updating the corresponding face information. This serves two purposes: data optimization and, crucially, supporting DyMeshVAE’s encoding process, which embeds vertex connectivity information to prevent trajectory inconsistencies during decoding.

The processed dynamic mesh files undergo temporal slicing with window sizes $T = 16$ and $T = 32$. To maximize data utilization, we initiate slicing from both frame 0 and $T//2$, storing T-frame segments sequentially. We also preserve reverse-ordered sequences as independent files, effectively augmenting the dataset by 3-4 \times . Each new sequence undergoes center normalization, positioning the initial frame at the origin with maximum vertex absolute values normalized to 1.0.

Post-slicing, we implement motion-based filtering, eliminating sequences with inter-frame maximum absolute differences outside the range [0.01, 0.5]. We also filter out the instances whose faces/vertices ratio exceeds 2.5. The cleaned data is then rendered using bpy scripts to generate frontal video sequences. We apply uniform gradient coloring (purple-red) and consistent top-down point lighting, utilizing the Cycles engine for 256×256 resolution rendering on CPU clusters.

For caption generation, we employ Qwen-2.5-VL [5] as our annotation model with the prompt: “Describe the motion of the object in a sentence.” The generated captions are stored alongside their corresponding .bin files.

Finally, we validate all processed files, removing examples with anomalous vertex or face shapes. The resulting DyMesh dataset is partitioned into subsets based on maximum vertex counts (4,096/8,192/50,000) to facilitate training and testing across different configurations.

A.2. Implementation Details

In the DyMeshVAE architecture, both encoder and decoder utilize attention mechanisms with a hidden dimension of

512. For temporal settings of $T=16$ and $T=32$, we employ latent dimensions of 32 and 64 channels respectively for VAE sampling, with both configurations containing approximately 25M parameters. The Shape-Guided Text-to-Trajectory Model consists of 12 stacked transformer blocks as shown in Fig. 4 of the main paper, where each block incorporates 8-head attention layers with features projected to 512 dimensions, totaling approximately 200M parameters. We also conduct scaling experiments with an enhanced architecture of 740M parameters, comprising 24 transformer blocks, 16 attention heads per layer, and a 1024-dimensional latent space. During training, we implement an efficient batchify strategy where each sample’s vertices and faces are padded to maintain uniform tensor dimensions across the batch: vertex tensors are padded with zero vectors (0.0, 0.0, 0.0) up to the dataset’s maximum vertex count, while face indices are padded with invalid indices (-1, -1, -1) to reach the maximum face count (defined as 2.5 times the maximum vertex count), enabling consistent batch processing while preserving mesh topology integrity. During inference, we perform rectified flow sampling using 64 uniformly sampled timesteps within the [0,1] interval.

A.3. Evaluation Metrics

For quantitative evaluation, we employ three standardized metrics from VBench [24] to assess the performance of comparative methods: I2V Subject Consistency, Motion Smoothness, and Aesthetic Quality. Specifically, I2V Subject Consistency is computed by measuring the frame-wise similarity of DINO [11] features, quantifying the visual coherence between the generated video and the reference image. Motion Smoothness is evaluated through the AMT [35] video interpolation framework, which assesses the temporal continuity and fluidity of the generated motion sequences. The Aesthetic Quality metric leverages the LAION aesthetic predictor to quantify the perceptual quality and artistic value of individual frames from a human-centric perspective. These complementary metrics provide a comprehensive evaluation of both temporal consistency and visual fidelity of the generated results.

For our user study protocol, we recruit 20 participants from diverse backgrounds and age groups to evaluate comparative methods through a controlled assessment. We randomly selected 10 diverse test cases and generated motion sequences for each using all comparative methods based on text prompts, rendering each result from four orthogonal viewpoints and concatenating them temporally ($16 \times 4 = 64$ frames). The results from all methods were randomly shuffled and horizontally concatenated into 64-frame GIFs, with participants rating each result on a 5-point Likert scale (5: excellent, 1: poor) across three criteria: text-motion alignment, motion plausibility, and shape preservation fidelity. The final evaluation scores were computed by aggregat-

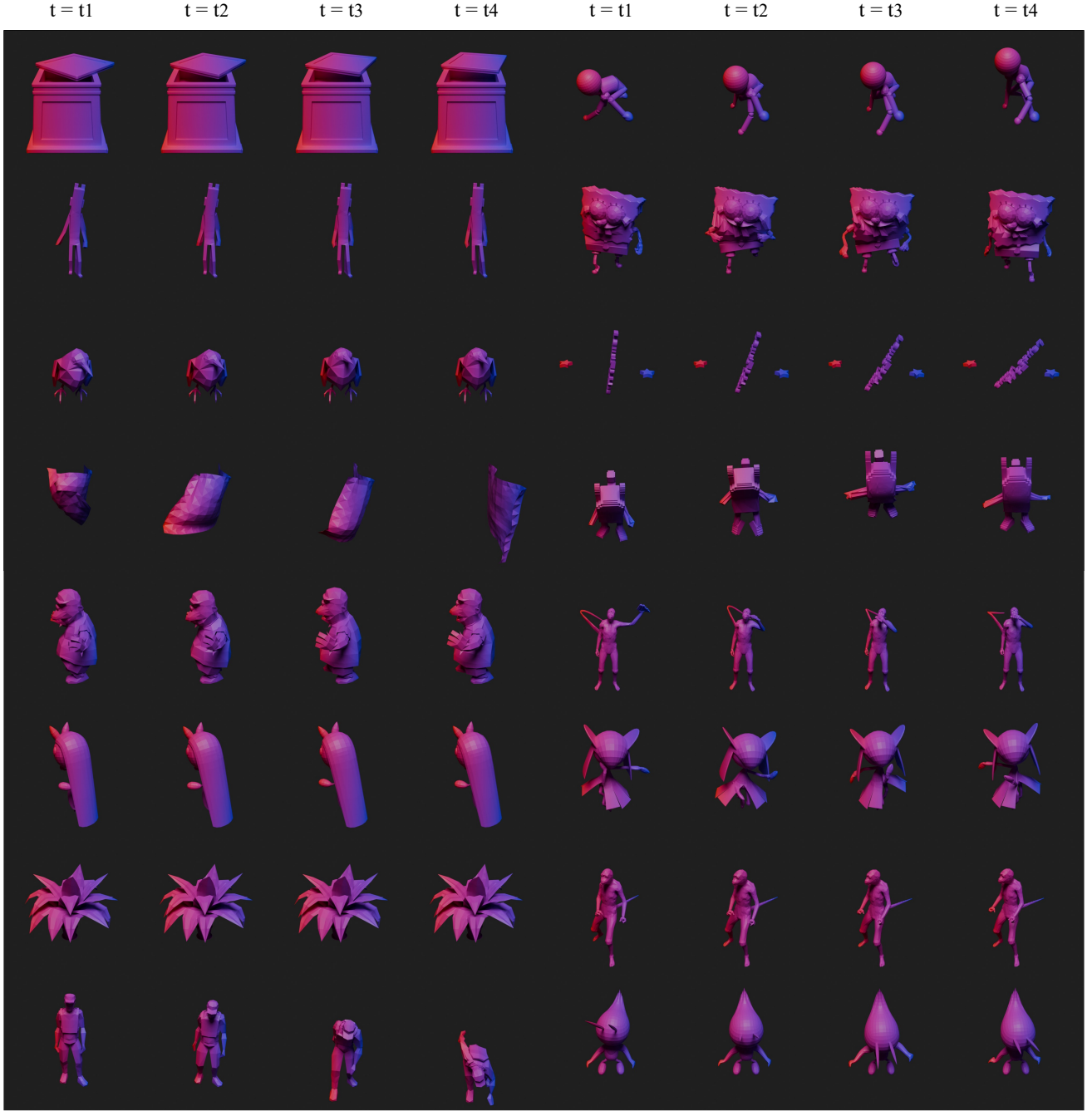


Figure 9. Examples of dynamic mesh sequences in our DyMesh dataset.

ing and de-shuffling ratings across all participants, with failed generations being handled by computing means from successful cases only, ensuring a comprehensive and unbiased assessment of perceptual quality and semantic accuracy. All the User Study source videos can be found in the *User_Study* folder of the supplementary materials.

The inference efficiency of our framework scales with

both mesh complexity (vertex/face count) and FPS feature sampling density. Our empirical studies indicate that an 8:1~4:1 ratio between vertex count and FPS samples achieves optimal performance-efficiency trade-off. The corresponding inference times across different mesh resolutions under this sampling configuration are presented in Tab. 4.

num_v	4,096	8,192	16,384	32,768
t (s)	3.95	5.99	10.68	21.86

Table 4. Inference time evaluation. num_v represents the number of mesh vertices. We sample 1/8 number of vertices in the FPS sampling procedure as default. All these testing is conducted on a single Nvidia A800 GPU.

B. Animation Results of AnimateAnyMesh

We curated a diverse collection of high-fidelity static meshes from Sketchfab [67], encompassing various categories including humanoid figures, animals, weapons, clothes, and environmental assets. These meshes were animated using our proposed AnimateAnyMesh framework through text-driven synthesis. Fig. 11 demonstrates representative results, showcasing our framework’s capability to generate high-fidelity, versatile animations across arbitrary mesh topologies. The qualitative results validate the effectiveness of our approach in achieving generalized mesh animation with exceptional geometric fidelity, motion naturalness, and semantic flexibility.

Moreover, given identical input prompts and initial mesh as condition, our AnimateAnyMesh framework demonstrates robust multi-modal synthesis capabilities through different random seeds, generating diverse yet plausible high-fidelity mesh animations. Fig. 12 illustrates this generative flexibility through exemplar results, highlighting our framework’s ability to explore varied motion manifestations while maintaining semantic consistency and geometric integrity.

C. Additional Qualitative Comparison

For a comprehensive comparison of mesh animation approaches, we present additional comparative examples against baseline methods in Fig. 13. The results consistently support our main findings, demonstrating that our approach outperforms existing methods in terms of text-motion alignment, motion naturalness, and shape preservation fidelity.

D. Additional Ablation Studies

num_v	FPS Sampling Ratio			
	1/32	1/16	1/8	1/4
369	136.57	58.78	22.45	8.10
2,567	17.74	3.64	1.16	0.54
6,890	1.57	0.66	0.51	0.50

Table 5. Reconstruction error for dynamic mesh sequences with varying vertex counts under different FPS sampling ratios.

FPS Sampling Ratio. We conduct ablation study on the sampling ratio in the DyMeshVAE encoding procedure to

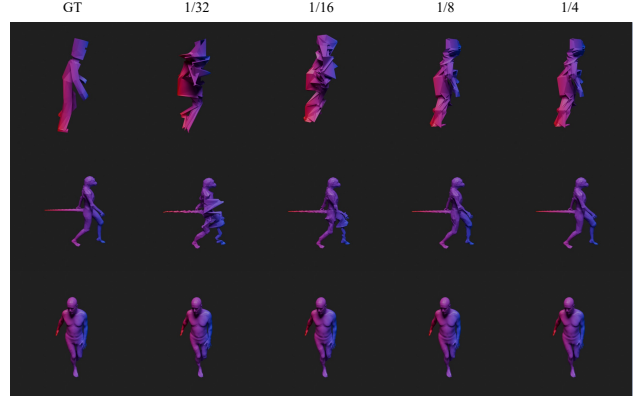


Figure 10. Ablation study on FPS Sampling Ratio. The numerical annotations above each image indicate the FPS ratio employed in the DyMeshVAE Encoder. Please zoom in for a better view.

see how it influences the reconstruction quality. We randomly sampled three dynamic mesh sequences from the DyMesh dataset with vertex counts of 369, 2,567, and 6,890 (top to bottom). For each mesh, we conducted experiments using FPS (Farthest Point Sampling) ratios of 1/32, 1/16, 1/8, and 1/4 for feature sampling and reconstruction. The qualitative visualizations and quantitative metrics are presented in Fig. 10 and Tab. 5, respectively, demonstrating the impact of sampling density on reconstruction fidelity.

As demonstrated in Fig. 10, meshes with low vertex counts (row 1_{st}) exhibit poor reconstruction quality even with a 1/4 sampling ratio, while high-vertex-count meshes (row 3_{rd}) maintain satisfactory reconstruction fidelity even at a 1/32 sampling ratio, as validated by quantitative metrics in Table 3. We attribute this phenomenon to spatial distribution characteristics: low-vertex meshes typically exhibit sparse spatial distribution, leading to significant geometric information loss in regions surrounding unsampled vertices, whereas high-vertex-count meshes maintain dense surface coverage even with lower sampling ratios, enabling better preservation of local geometric features during encoding. Based on these observations, we set the feature sampling count to 512 during training to facilitate efficient batch processing while achieving an optimal balance between performance and computational efficiency across meshes with diverse vertex counts. During inference, we adopt an adaptive sampling strategy where $n = \min(512, num_v // 8)$ for inference, where num_v represents the number of mesh vertices. We empirically find the robust performance across the majority of test cases with this setting.

E. Scaling Experiments

In this section, we conduct scaling experiments across three dimensions: dataset scale, temporal resolution (frame count), and model capacity. We evaluate four configu-

	Experiment	I2V \uparrow	M.Sm \uparrow	Aest.Q \uparrow	Dy.Dg \uparrow
A	rf_4096v_16f_200Mp	0.954	0.995	0.539	0.693
B	rf_8192v_16f_200Mp	0.985	0.996	0.550	0.605
C	rf_4096v_32f_200Mp	0.948	0.993	0.532	0.737
D	rf_4096v_16f_740Mp	0.968	0.997	0.545	0.705

Table 6. Scaling experiments of AnimateAnyMesh. We name the experiments as rf_< 1 > v_< 2 > f_< 3 > p, where 1, 2, 3 represents the maximum number of vertices of the dataset, number of frames per instance, and the number of parameters of the Shape-Guided Text-to-Trajectory Model.

rations, each denoted as rf_< 1 > v_< 2 > f_< 3 > p, where 1, 2, 3 represents the maximum number of vertices of the dataset, number of frames per instance, and the number of parameters of the Shape-Guided Text-to-Trajectory Model. All models were trained for 600,000 iterations on corresponding DyMesh subsets with a batch size of 2048 and a learning rate of $2e-4$. For evaluation, we generate mesh animations using the same 10 mesh-prompt pairs from our qualitative benchmark, maintaining consistent random seeds across all models. Front-view renderings were produced to compute the VBench [24] metrics (I2V, M.Sm, Aest.Q) discussed in the main text. Additionally, given that all four trained models demonstrate the capability to generate high-quality and semantically plausible mesh animations, we incorporate the Dynamic Degree metric (abbreviated as Dy.Dg) from VBench to quantitatively assess motion intensity. The comprehensive results are presented in Tab. 6. The results indicate that: (1) Increasing the maximum number of vertices leads to better performance on most metrics (B vs. A). (2) Increasing the number of frames will improve the output dynamic, improving Dy.Dg while maintaining promising results on other metrics. (C vs. A). (3) Scaling the model’s parameter size leads to better performance on all metrics, demonstrating good scalability of our method. (D vs. A).

F. Limitation

Our work exhibits three limitations: First, regarding dataset scale, while the proposed DyMesh Dataset encompasses over 4M dynamic mesh sequences, the number of unique mesh identities remains below 100k, potentially limiting model generalization across specialized categories. We plan to address this by creating and curating additional high-quality, diverse 4D datasets. Second, concerning annotation quality, we observe that current video captioning models demonstrate suboptimal performance when annotating 3D rendered sequences without natural backgrounds, compared to their performance on natural videos, particularly in motion description granularity. Enhancing caption fidelity for synthetic 3D content remains a key research direction. Third, in terms of model capabilities, the current implemen-

tation of AnimateAnyMesh is confined to 16/32-frame sequence generation, and extending the model’s capability to generate longer-duration mesh animations represents a significant future research objective.

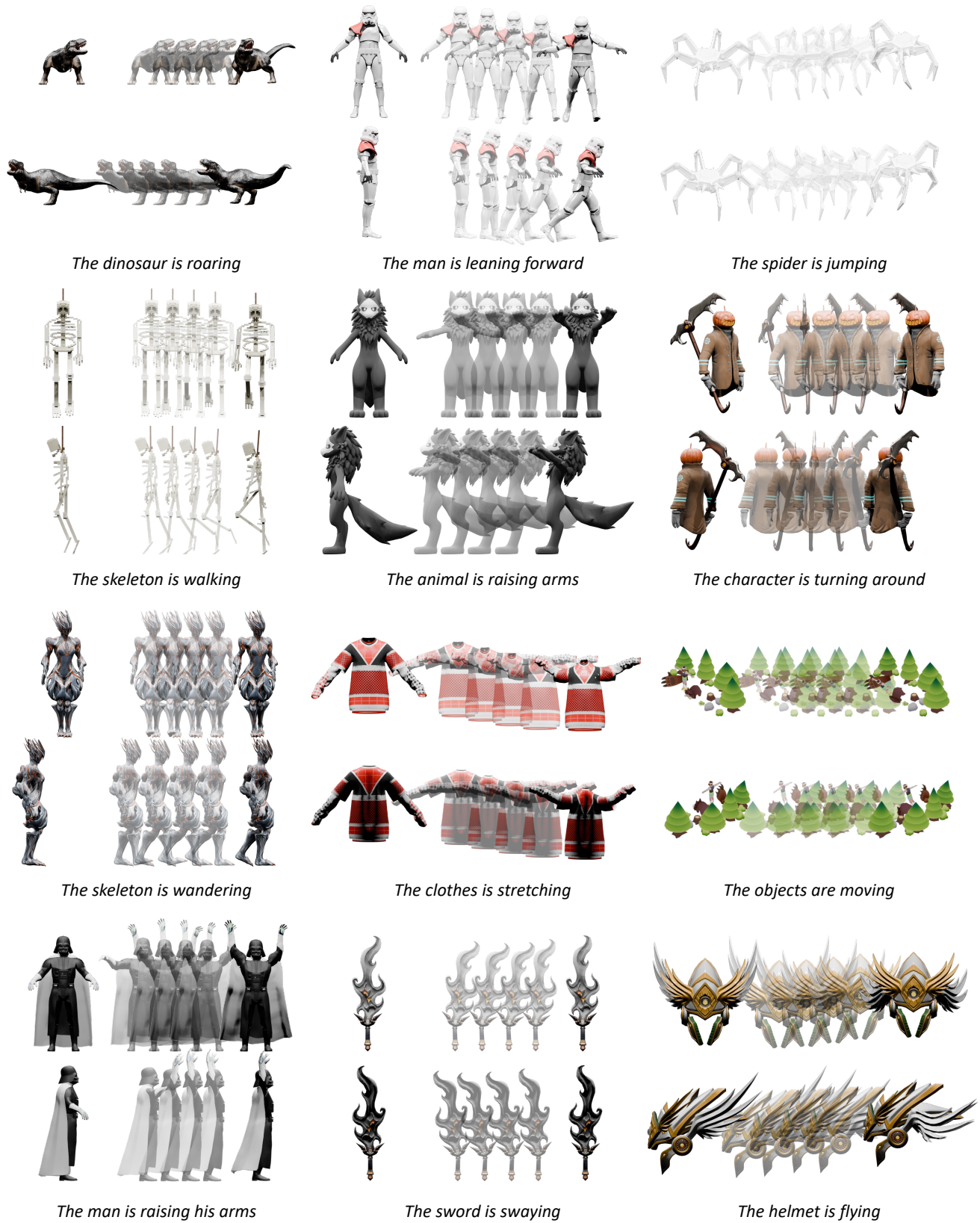


Figure 11. Examples of text-driven mesh animation results of the proposed AnimateAnyMesh. We render two random views for each example. Please zoom in for a better view.



Figure 12. Diversity Demonstration of AnimateAnyMesh Generations. Given identical text-prompt and initial mesh conditions, AnimateAnyMesh demonstrates the capability to generate diverse, high-quality mesh animations through random seed variation.



Figure 13. Additional qualitative comparison with state-of-the-art mesh animation methods.