

Make-It-Animatable: An Efficient Framework for Authoring Animation-Ready 3D Characters

Zhiyang Guo^{1*}Jinxu Xiang²Kai Ma²Wengang Zhou¹Houqiang Li¹Ran Zhang²

¹CAS Key Laboratory of Technology in GIPAS, EEIS Department,
University of Science and Technology of China

²Tencent PCG

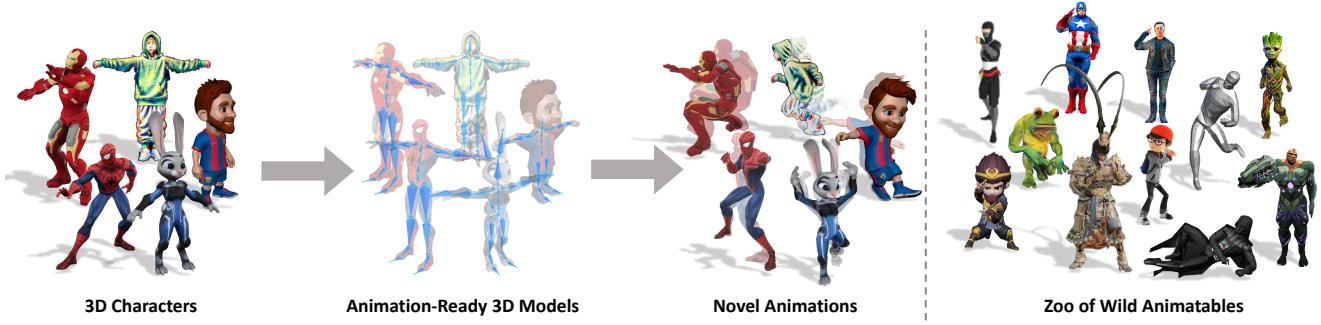


Figure 1. Given a 3D character represented by mesh or 3D Gaussian Splats with arbitrary pose and shape, our framework can produce high-quality results of rigging, skinning, and pose resetting for it within one second. The output 3D model is fully animatable with a fine-grained skeleton and optional bone topology of extra body structures.

Abstract

3D characters are essential to modern creative industries, but making them animatable often demands extensive manual work in tasks like rigging and skinning. Existing automatic rigging tools face several limitations, including the necessity for manual annotations, rigid skeleton topologies, and limited generalization across diverse shapes and poses. An alternative approach is to generate animatable avatars pre-bound to a rigged template mesh. However, this method often lacks flexibility and is typically limited to realistic human shapes. To address these issues, we present **Make-It-Animatable**, a novel data-driven method to make any 3D humanoid model ready for character animation in less than one second, regardless of its shapes and poses. Our unified framework generates high-quality blend weights, bones, and pose transformations. By incorporating a particle-based shape autoencoder, our approach supports various 3D representations, including meshes and 3D Gaussian splats. Additionally, we employ a coarse-to-fine representation and a structure-aware modeling strategy to

ensure both accuracy and robustness, even for characters with non-standard skeleton structures. We conducted extensive experiments to validate our framework's effectiveness. Compared to existing methods, our approach demonstrates significant improvements in both quality and speed. The source code will be made publicly available.

1. Introduction

3D characters, as the principal subjects of the digital world, play an essential role in various fields of modern creative industries, e.g., video games, 3D animations, films, mixed reality, etc. Bringing 3D characters to life requires making them animatable, a crucial yet labor-intensive task that involves substantial effort in rigging and skinning. Additionally, animation-ready character models must be positioned in predefined rest poses to support effective retargeting and deformation. Traditionally, rigging and skinning processes are either manually executed by artists, demanding considerable time and effort for each task, or managed using existing automatic rigging tools that often lack the robustness and generalizability to accommodate diverse character types and animation requirements. Recent advancements have sought to address these challenges by employing pa-

*Work is done during internship at Tencent.

¹Permission Pending: Some of the images in this paper are temporarily masked due to copyright issues. We will update them as soon as possible.

Categories	Methods	Mesh	3DGs	Template Free	Alterable Skeleton	Pose to Rest	Hand Animation	Rigging ³ Time
Text/Image to Animatable	Meshy [†] [36] Tripo [‡] [53]	✓ ✓	✗ ✗	✓ ✓	✗ ✗	✓ ✓	✗ ✗	~3 min ~3 min
Auto Rigging	Mixamo [§] [1] Anything World [†] [33] RigNet [57] TARig [38]	✓ ✓ ✓ ✓	✗ ✗ ✗ ✗	✓ ✓ ✓ ✓	✗ ✓ ✓ ✓	✓ [¶] ✓ [¶] ✗ ✗	✓ ✓ ✗ ✗	~2 min ~4 min ~10 min ~0.6 min
Template-based Animatable	TADA [32] HumanGaussian [35]	✓ ✗	✗ ✓	✗ ✓	✗ ✗	— —	✓ ✗	— —
	Ours	✓	✓	✓	✓	✓	✓	~0.5 s

[†] Commercial software.

[‡] Mixamo and Anything World only support simple poses like T- or A-pose, while our methods allow for arbitrary poses.

[§] Meshy and Mixamo need manual annotation of joint positions.

[¶] Tested on a typical input mesh with 8k vertices.

Table 1. Comparison between our method and existing approaches in terms of key features.

parameterized templates to directly generate animatable characters as meshes or 3D Gaussian splats. However, these methods are frequently constrained by their dependence on realistic humanoid models and lack the flexibility to handle non-standard poses or shapes, limiting their applicability to the dynamic and varied 3D character designs prevalent in modern applications.

In this work, we present a comprehensive data-driven framework that enables instant rigging and skinning of 3D characters, regardless of their initial shapes, poses, or structural complexities. Our solution is exceptionally fast, processing each character in approximately one second while generating high-quality, animation-ready models with accurate bones, blend weights, and rest pose transformations. Unlike other automatic rigging approaches, our system adeptly handles challenging cases, including exaggerated characters with unconventional head and body proportions, as well as non-standard poses featuring additional bone structures such as hands, ears, and tails.

Tab. 1 compares our method with existing approaches across key features, including the ability to handle both mesh and 3D Gaussian splatting inputs, flexibility in adapting to various poses, and support for advanced animation features such as hand and finger articulation.

Unlike commercial auto-rigging tools like [1, 33], our model is designed to work with any predefined skeleton structure, providing greater control over joint movements and significantly improving rigging speed and flexibility. Our framework’s rapid and precise rigging and skinning of diverse 3D characters unlocks new possibilities for dynamic character animations. This capability is particularly beneficial in applications requiring swift responses and high customization, such as virtual reality, gaming, and real-time simulations. Additionally, it serves as a valuable enhancement to existing static 3D model generation systems, enabling the creation of dynamic 3D models.

2. Related Works

2.1. Modeling Dynamics in 3D Vision and Graphics

Modeling dynamics is a fundamental task in computer vision and computer graphics. It involves both the spatial

representation of 3D models and the formulation of their temporal behaviors. Recent approaches in this area can be categorized into mesh-based, mesh-free, reduced-order, and proxy-based methods.

Mesh-based elasticity models have been a cornerstone in modeling dynamics within computer graphics [52]. One of the most commonly employed variants is the linear tetrahedral finite element method [7, 9, 16, 49], which requires a robust surface-to-tetrahedron algorithm like TetWild [20].

However, both the linear finite element method and tetrahedral generation are sensitive to the geometry of the model, necessitating special handling for intricate shapes, such as thin features [2] and rod-like geometries [8]. For complex geometries, including non-manifold shapes, mesh-free methods, first introduced by [19], provide an effective solution for addressing the challenges of mesh-based dynamic modeling. For example, smoothed-particle hydrodynamics (SPH) [41, 44] and material point method (MPM) [24, 42, 50], which use particles as the main spatial representation, overcomes the difficulties of modeling different materials and complex dynamics in mesh-based methods, yielding high-quality results.

One major drawback of particle-based methods is their high computational complexity, resulting from the dense sampling required. To address this, researchers introduced reduced-order dynamics [4–6], which decrease the degrees of freedom by projecting motions onto a limited set of deformation bases. This approach significantly accelerates computations and is applicable to both mesh-based and mesh-free methods. Recently, methods such as CROM [13] and LiCROM [12] have advanced continuous reduced-order modeling using neural networks.

Reduced-order dynamics simplify temporal behaviors by reducing their dimensionality. Similarly, spatial representations can be simplified through proxy-based methods, which employ low-dimensional proxies to efficiently compute dynamics and subsequently apply deformations to the original geometry. Notable proxy-based techniques include Free-Form Deformation (FFD) [48] and cage-based methods [15, 25, 26, 51]. These approaches utilize simplified spatial representations to expedite dynamic calculations, which are then used to deform the detailed geometry.

The most widely used proxy-based method is Linear Blend Skinning (LBS). It deforms a mesh by blending transformations from an underlying skeletal structure, allowing for efficient and intuitive animation of complex models. This technique was introduced by Magnenat-Thalmann et al. [39] and further developed by Lewis et al. [29]. Recent advancements with neural networks have further enhanced LBS, making it suitable for geometry-agnostic simulation [40] and enabling more natural deformations through neural blend shapes [30].

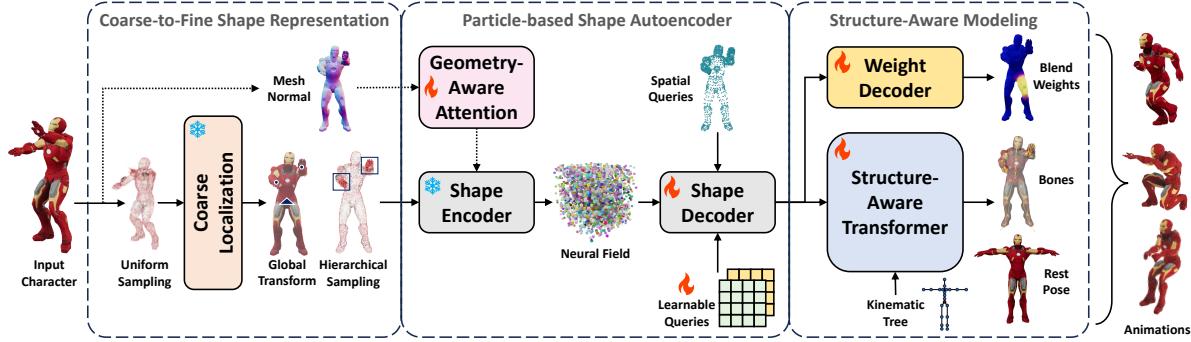


Figure 2. Pipeline of the proposed framework. Given an input 3D character, we produce high-quality blend weights, bones, and pose-to-rest transformation for it, so that any animation is within easy reach. First, we coarsely localize the joints with a pre-trained lite version of this framework, which enables a finer shape representation. Then the shape is encoded into a neural field with a particle-based autoencoder. The decoding process involves spatial and learnable queries for different animation assets. Finally, the structure-aware modeling of bones is proposed to better align the predictions with skeleton topology priors.

2.2. Authoring Animation-Ready 3D Models

In the 3D animation and gaming industry, Linear Blend Skinning (LBS) is the standard for character animation. Creating animation-ready 3D models with LBS involves two key processes: constructing rigs and assigning skinning weights to the input models. Traditionally, these tasks are performed manually in 3D modeling software like Autodesk Maya, where artists place rigs and paint skinning weights to achieve the desired deformations.

Pinocchio [3] is a pioneering system for automatic rigging and skinning of 3D characters. It automates the embedding of a skeleton into a 3D mesh and assigns skinning weights based on vertex proximity to bones, enabling smooth and natural deformations during animation. Recent advancements in automatic rigging have leveraged deep learning techniques to enhance flexibility and accuracy. RigNet [57] utilizes neural networks trained on extensive datasets of animated characters to predict rigging parameters, enabling the generation of custom rigs suitable for diverse character models. Similarly, TARig [38] employs a template-aware approach, combining a humanoid skeleton with a shared graph neural network backbone to ensure precise joint positioning and skin weight estimation, thereby improving rigging efficiency. To enhance skinning quality, Li et al. [30] proposed neural blend shapes in addition to predicting rigs, achieving better deformation compared to static skinning weights. Another advancement in automatic rigging involves relaxing the requirement for the input pose to be in a standard position, such as A-pose or T-pose; Theisel et al. [27] introduced a method that generates rigs for characters in arbitrary poses. Nonetheless, the inadequate quality, the limited robustness against complex inputs, and the unsatisfying time cost, all hinder the practical application of these auto-rigging methods.

An alternative approach to creating animation-ready models involves generating 3D models that are pre-bound to a rigged template. For instance, TADA [32] produces

textured human meshes by deforming a SMPL-X [43] template through score distillation sampling (SDS) optimization [46]. Similarly, HumanGaussian [35] generates human models in Gaussian splats from a SMPL [37] mesh using a structure-aware SDS algorithm. Additionally, DreamWaltz-G [21] employs a skeleton-guided distillation method combined with a hybrid 3D Gaussian avatar representation to achieve realistic generation and expressive animation. Although the aforementioned works have introduced various strategies to improve the shape diversity of generation, they remain constrained by the preset body ratio of SMPL mesh and are limited to avatars resembling realistic humans. Moreover, the fixed template skeleton topology also prevents their extension to non-standard bone structures.

3. Method

3.1. Preliminaries

Intuitively, modeling the dynamics of an object requires a per-timestamp deformation of all the particles (*i.e.*, vertices for meshes or splats for 3D Gaussians) that make up its geometry. Suppose that we have an object composed of N particles and a desired dynamic sequence of T timestamps. The temporal deformations of all the particles can be denoted by a matrix $\mathbf{D} \in \mathbb{R}^{T \times N \times d}$, where d is the degrees of freedom (*e.g.*, $d = 6$ for rigid transformations). Generally, \mathbf{D} has a lot of redundancy when representing real-world dynamics, so we would like to seek a low-rank approximation expressed by $\mathbf{D}^{T \times N \times d} \approx \mathbf{B}_t^{T \times K \times d} \mathbf{W}_s^{K \times N}$, where K is the desired rank, \mathbf{B}_t and \mathbf{W}_s are *temporal basis* and *spatial weight* matrices, respectively. The practical applications of this low-rank form of dynamics include the sparse control paradigms [22, 28] and the linear blend skinning (LBS) [37, 43] algorithm, where \mathbf{B}_t is interpreted as the transformations/poses of control nodes/body joints, and \mathbf{W}_s is the blend weights of the nodes or joints. Since the motions under the same joint setting can be reused, solving

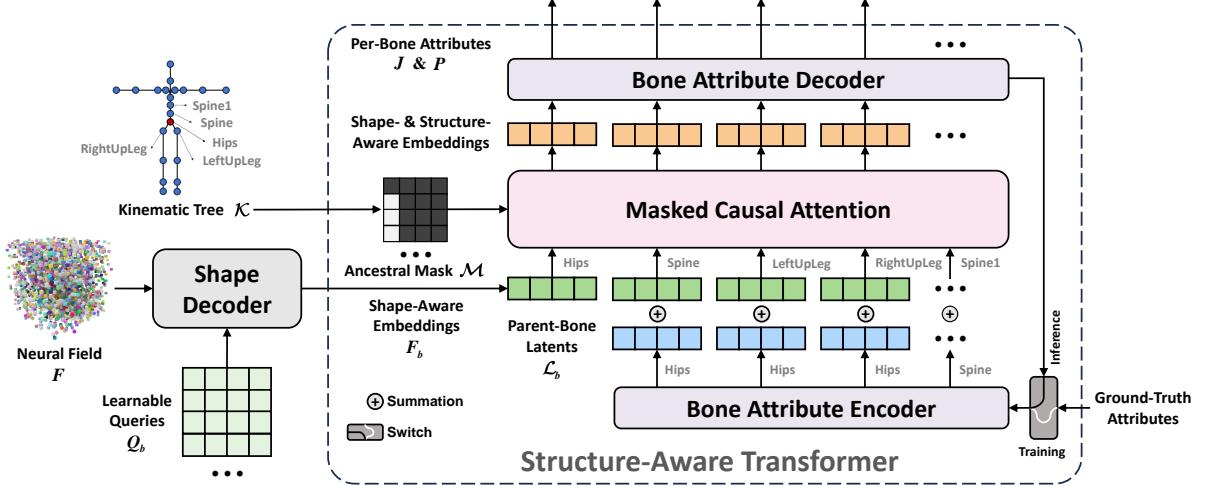


Figure 3. **Pipeline of the proposed structure-aware transformer.** The per-bone shape-aware embedding is first added with its parent bone’s latent, which is encoded from the autoregressive outputs (in inference) or the ground-truth values (in training). The summation is then fused with the ancestral bones’ features via the masked causal attention. Eventually, bone attributes are decoded from the output shape- and structure-aware embeddings. In inference, the whole process follows the paradigm of next-child-bone prediction.

B_t for humanoid characters is usually converted to finding the transformations from all poses to a predefined rest pose.

Based on this formulation, given a shape representation of a character composed of N particles, all the assets necessary for animating it (named as *animation assets*) are: **1)** K joints represented by bone head and tail positions, $\mathbf{J} \in \mathbb{R}^{K \times 6}$, which are connected based on any predefined skeleton topology; **2)** a 6-DoF transformation, $\mathbf{P} \in \mathbb{R}^{K \times 6}$, which transforms the input posed joints to a common rest pose; **3)** the blend weights $\mathbf{W} \in \mathbb{R}^{K \times N}$ between any joint-particle pair. The former two are also referred to as *bone attributes* since they are orthogonal to the input particles.

In the following subsections, we will elaborate on how these animation assets are obtained with our proposed unified framework, which is illustrated in Fig. 2. For the sake of clarity, we will first introduce the core part of this framework — a shape autoencoder that describes the input geometry as compact neural latent representations.

3.2. Particle-based Shape Autoencoder

Neural field encoding. The encoding part of our shape autoencoder is similar to the downsampled-point querying architecture of 3DShape2VecSet [58]. We start from a point cloud $\mathbf{X} \in \mathbb{R}^{N \times 3}$ sampled from the input character shape after re-centering and normalizing. In order to learn a compact representation for this shape, we want to aggregate the information in this possibly large point cloud into a smaller set (size M) of latent feature vectors denoted by $\mathbf{F} \in \mathbb{R}^{M \times C}$. To achieve this, \mathbf{X} is first downsampled to M points with farthest point sampling (FPS) [47]: $\tilde{\mathbf{X}} = \text{FPS}(\mathbf{X}) \in \mathbb{R}^{M \times 3}$. Then equipped with a positional encoding mapping denoted by $\text{PE} : \mathbb{R}^3 \rightarrow \mathbb{R}^C$, the latent features \mathbf{F} are obtained via a cross-attention between the

original and the downsampled points:

$$\mathbf{F} = \text{CrossAttn}(\text{PE}(\tilde{\mathbf{X}}), \text{PE}(\mathbf{X})) \in \mathbb{R}^{M \times C}. \quad (1)$$

In this way, the network can adaptively encode the spatial information as a neural field, instead of representing the actual spatial position explicitly. In practice, the shape encoder pretrained on ShapeNet [11] already has sufficient capacity for low-level geometry perception. Therefore, we freeze the parameters of the shape encoder during the training of our framework.

Geometry-aware attention. One of the benefits of our particle-based autoencoder is the support for various shape representations (as long as point sampling is possible on them). Nevertheless, for surface representations like mesh, much geometric information about the shape is lost during the sampling. This can lead to blend weight corruption problems that semantically distant but spatially close points may incorrectly share similar weight values. To address this issue, a non-intrusive way of injecting geometric awareness into the shape representation is proposed. Specifically, we extract the per-point normal values from the input mesh, as they carry rich geometry information. Then a lightweight fusing module is implemented using a simple attention layer as a branch way to the shape encoder, whose parameters are initialized so that the normal branch produces zero impact on the encoding results at the beginning. We randomly corrupt the normal values in training to make it more robust to low-quality meshes and also compatible with inputs like 3DGs in inference. Eventually, the attention mechanism adaptively decides what regions should care about the normals and whether the given normal values carry useful cues.

Spatially continuous decoding for blend weights. The neural field represented by latent features \mathbf{F} can be directly

queried with coordinates to obtain spatially continuous attributes, *i.e.*, the blend weights. Our shape decoder completes decoding in a flexible and learnable way by applying attention between the queries and the neural fields. Formally, given the N_q -point spatial queries (*e.g.*, all vertices of a mesh) denoted by $\mathbf{Q}_w \in \mathbb{R}^{N_q \times 3}$, the corresponding shape-aware embeddings \mathbf{F}_w is derived by

$$\mathbf{F}_w = \text{CrossAttn}(\text{PE}(\mathbf{Q}_w), \text{SelfAttn}(\mathbf{F})) \in \mathbb{R}^{N_q \times C}. \quad (2)$$

Then using an MLP-based lightweight decoding head Θ_w , we can finally get the blend weights $\mathbf{W} = \Theta_w(\mathbf{F}_w) \in \mathbb{R}^{K \times N_q}$ for all the query coordinates.

Learnable discrete decoding for bone attributes. When it comes to discrete per-bone attributes like joints’ positions and poses, the decoding process becomes more tricky. While following the same querying pattern [37, 38] of blend weights and predicting a joint regressor is possible, it’s inefficient for dense inputs and makes bone attributes sensitive to the spatial queries. Therefore, we adopt another way that involves learnable semantic queries. Specifically, for an attribute of the predefined K bones, we assign K learnable query embeddings denoted by $\mathbf{Q}_b \in \mathbb{R}^{K \times C}$. Then we use attention layers to integrate the shape features, similar to Eq. (2):

$$\mathbf{F}_b = \text{CrossAttn}(\mathbf{Q}_b, \text{SelfAttn}(\mathbf{F})) \in \mathbb{R}^{K \times C}. \quad (3)$$

Now the bone-wise embeddings \mathbf{F}_b contain both the global geometry information and bone semantic cues, *i.e.*, which region should a certain bone attend to. Then two individual bone attribute decoders are employed to regress the desired per-bone head-and-tail positions $\mathbf{J} \in \mathbb{R}^{K \times 6}$ and pose-to-rest transformation $\mathbf{P} \in \mathbb{R}^{K \times 6}$, respectively. We found in practice that the rigid transformation represented by dual quaternions presents better behaviors in optimization than a vanilla 6-DoF movement. Therefore, the pose decoder actually outputs 8D dual quaternions $\mathbf{P}_{dq} \in \mathbb{R}^{K \times 8}$ that can be trivially converted to a rigid transformation matrix.

3.3. Coarse-to-Fine Shape Representation

With the shape autoencoder, our model can produce promising blend weight predictions, but joint outputs, particularly for fine-grained regions like the hands, still face convergence challenges. We attribute these issues to the ambiguity of the input points and view the aforementioned pipeline as a coarse stage that provides rough but valuable localization information about the input character. To enhance the shape representation, we adopt two strategies: canonical transformation and hierarchical sampling. The canonical transformation resolves pose ambiguity by aligning the input shape to a common orientation, while the hierarchical sampling method ensures higher sampling density in key regions like the hands to improve accuracy without increasing computational cost. Further implementation details are provided in the supplementary material.

3.4. Structure-Aware Modeling of Bones

With a learnable query assigned to each bone, we can effectively predict the bone-wise attributes. Although those queries adaptively attend to different regions of the input shape, they lack awareness of the skeleton topology, leading to imperfections in predictions, especially for deep-level bones in a kinematic tree. Therefore, we propose to model the bones in a structure-aware way, including additional designs for both network architecture and loss functions.

Next-child-bone prediction via causal attention. Without structure awareness of the skeleton, the model tends to produce independent homogeneous predictions. For example, the pose of a finger bone largely depends on the poses of its ancestral bones like hand and arm. If the pose of the direct parent is known, predicting the finger’s pose becomes easier, as it only requires accounting for the subtle transformation relative to its parent. To this end, inspired by the decoder-only next-token prediction paradigm prevalently adopted in LLMs [10, 54], we develop a structure-aware transformer based on a novel *next-child-bone prediction* architecture, as depicted in Fig. 3.

Let us start from the bone-wise shape-aware embeddings, $\mathbf{F}_b = \{\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_K \in \mathbb{R}^C\}$, produced by the shape decoder (Sec. 3.2). Instead of directly feeding \mathbf{F}_b into a bone attribute decoder, we add several masked causal attention layers to enhance them. Each bone embedding, \mathbf{f}_i , $i \in [1, K]$, is regarded as an individual query token. Meanwhile, we use a bone attribute encoder to map the actual attribute values (joint positions or poses) into some same-shaped bone-wise latents $\mathcal{L}_b = \{\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_K \in \mathbb{R}^C\}$. Then given a kinematic tree denoted by \mathcal{K} , to let \mathbf{f}_i learn about the context information from its parent bone, we add the corresponding parent-bone latent \mathbf{l}_j , $j = \text{Anc}(i, \mathcal{K})(1)$ to it, where $\text{Anc}(i, \mathcal{K})$ returns a tuple of ancestral bones for the bone i recursively from the nearest up to the root node in \mathcal{K} , and here we only select the first element (*i.e.*, the direct parent bone).

We then feed the summation of \mathbf{f}_i and \mathbf{l}_j into a masked causal attention module that fuses that per-bone latent with those of other bones. We build an ancestral mask derived from the kinematic tree, ensuring that each bone learns from its ancestors and treats the states of its children as “future” information. The ancestral mask is defined by $\mathcal{M} \in \mathbb{R}^{K \times K}$ where $\mathcal{M}_{i,j} = 1$ iff $j \in \text{Anc}(i, \mathcal{K})$. Equipped with \mathcal{M} , the cross-attention operation outputs per-bone embeddings that have both shape and structure awareness. Eventually, they are sent into a bone attribute decoder to be translated back to the predicted attributes.

Similar to the causal attention in language models, our structure-aware transformer follows an autoregressive process at the inference stage. The ground-truth parent-bone attributes are replaced by the predictions produced by the decoder. It takes several progressive iterations to complete the

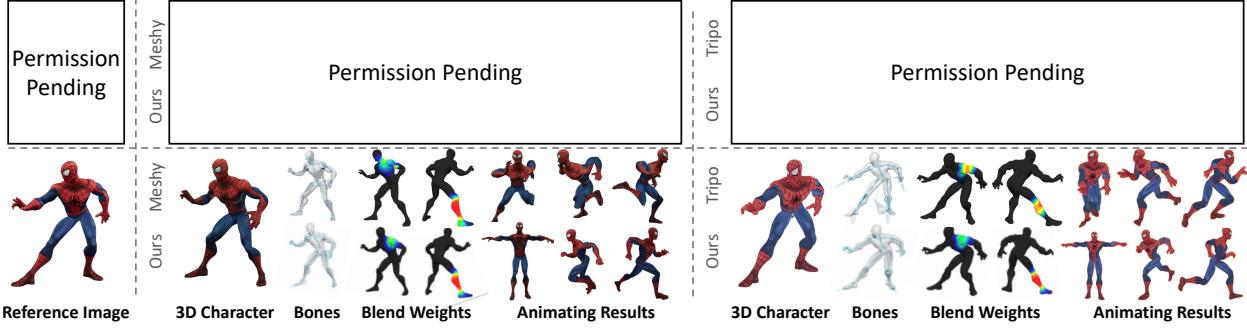


Figure 4. **Comparison with Meshy [36] and Tripo [53]**. We feed them the same image as reference and compare the performance based on their generated 3D models respectively. The blend weights of two joints, *i.e.*, Left Shoulder and Right Leg, are visualized. Given that these baselines can only apply preset motions and their rest-pose models cannot be exported, we apply a similar “running” sequence to all the methods for fair comparison. The T-pose models predicted by our method are also included as the front-view animating results.

decoding and one level in the kinematic tree is predicted in an iteration, hence the name “next-child-bone prediction”.

Body prior loss. Our predictions of bone attributes are directly supervised by the ground-truth values. Since a character can have multiple possible solutions in practice, we further constrain the optimizing process using prior losses, so that the predictions follow some essential body patterns. Specifically, three kinds of prior knowledge are leveraged as follows. 1) Bone connectivity: most bone heads have to be connected to the tails of their parent bones. 2) Bone symmetry: typically in the predefined rest pose, the left and right parts of the skeleton have to be symmetric about the spine. 3) Bone parallelism: in the rest pose, multiple bones belonging to the same limb have to share the same direction.

3.5. Training and Inference

The proposed framework is trained in an end-to-end data-driven manner, supervised by the L_1 losses with the ground-truth blend weights, bone positions, and pose-to-rest transformations, as well as the extra body prior losses (Sec. 3.4). As introduced in Sec. 3.3, a coarse-to-fine training strategy is adopted. In the coarse stage, we uniformly sample the input shape and predict only the bone positions, applying random 3D rotations for data augmentation to enhance generalization. Then in the fine stage, we transform the input character to canonical coordinates and apply the hierarchical sampling.

Once the training is finished, the model can take any particle-based character (in any pose with any global transformation) and infer its animation assets. With a fixed number of neural shape latents and a learnable discrete querying strategy, our framework achieves efficient inference, with speed less affected by input particle number. Generally, the whole feed-forwarding inference takes less than 1 second.

Although our framework requires a fixed definition of kinematic tree, and we typically choose a standard human skeleton to train it on, it can actually be extended to any predefined skeleton topology. If the framework has already been trained with the standard skeleton, adapting it to in-

clude extra bones is as simple as fine-tuning the final layer of the weight decoder and the extra learnable queries. In Sec. 4.2, we will present some animating results of characters with additional accessories like long ears or tails.

4. Experiments

4.1. Experimental Settings

Dataset. We utilize a collection of artist-designed 3D models from Mixamo [1], consisting of 95 high-quality characters and 2,453 diverse motion sequences. Each character is preprocessed to conform to a standard skeleton structure with $K = 52$ bones, and the motion sequences contain an average of 200 frames. We allocate 95% of the data for training and the remaining 5% for validation. During each training iteration, a character-motion pair is randomly selected, resulting in an effective dataset size equivalent to over 40 million frames.

Furthermore, our training framework accommodates additional skeleton topologies. We used Vroid Studio [45] to manually create 35 distinct anime characters, each featuring accessories like animal ears and tails. These characters were preprocessed to match the Mixamo skeleton definition, with the addition of extra bones for these accessories.

We select bipedal humanoid characters from the ‘ModelsResource-RigNetv1’ dataset [56] to construct a test set for quantitative comparison with existing automatic rigging algorithms. For qualitative evaluation, we gather a diverse collection of in-the-wild 3D characters from various sources, including the Objaverse dataset [18], generative tools [36, 53], and other artist-designed models.

Baselines. We choose two lines of existing methods that can produce animatable characters as our baselines. *Autorigging methods.* These methods share the same workflow as ours, *i.e.*, taking a 3D character as the input and outputting a rigged one with bones and blend weights. a) Meshy [36]: a commercial software that generates general 3D meshes from text/image prompts, and optionally rigs the meshes of bipedal humanoid inputs. Since it does not

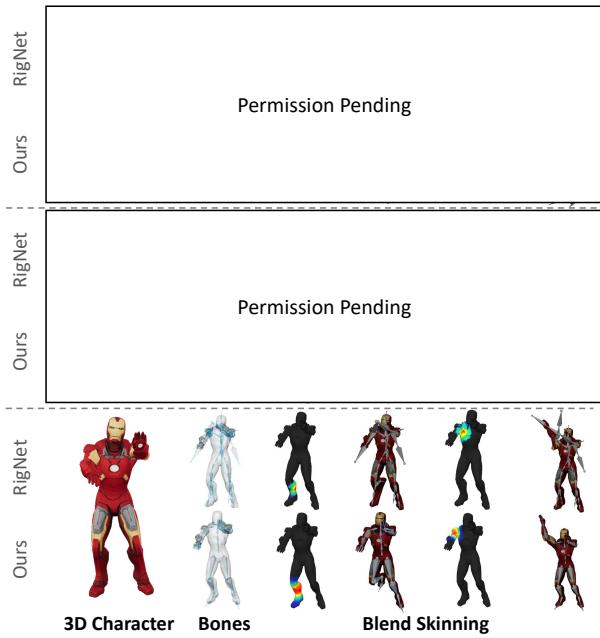


Figure 5. **Comparison with RigNet [57]**. We visualize the blend weights of selected joints and manually deform them to assess the impact of rigging quality on skinning results.

support direct 3D inputs, we feed an image of the test 3D character into it. b) Tripo [53]: another commercial software that basically has the same functions as Meshy, but differs in quality. c) RigNet [57]: a data-driven model that can rig any mesh-based geometry. We feed the same test characters to both RigNet and our model. *Template-based avatar generating methods*. By leveraging the SMPL template [37, 43] as a strong prior, this line of methods directly generates animatable humans. a) TADA [32]: mesh-based avatar generation. b) HumanGaussian [35] (HG): Gaussian-splatting-based avatar generation. The generated avatars of these two works are rigged with their corresponding template skeletons and blend weights.

Metrics. In the ablations conducted on the test split of Mixamo dataset, we report the percentage error between the predicted and ground-truth animation assets. As for the quantitative comparison with RigNet [57], we employ several metrics to evaluate the quality of skeleton predictions, *i.e.*, the IoU, Precision and Recall of bone matching, and the CD-J2J, CD-J2B, and CD-B2B (Chamfer distances between joints and bone line segments) [56]. For other in-the-wild test cases without ground truth, we provide extensive qualitative comparison by visualizing the skeletons, blend weights, and animating results. Due to the different skeleton topologies, the same motion sequence is usually not applicable to other baselines. For fairness, we manually set the bones from different methods to similar poses in 3D modeling software and render the animations.

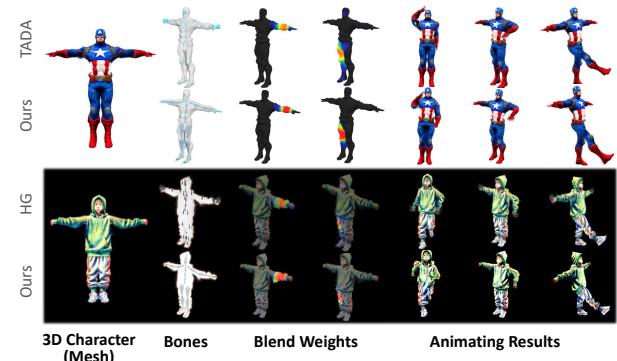


Figure 6. **Comparison with TADA [32] and HumanGaussian [35] (HG)**. We use the generated meshes from TADA and 3D Gaussians from HG for comparison. Note that the skeletons of these two baselines are identical to the shape-specific SMPL [37] templates (without bone tail), with their weights interpolated from the template meshes. Zoom in to better view the details.

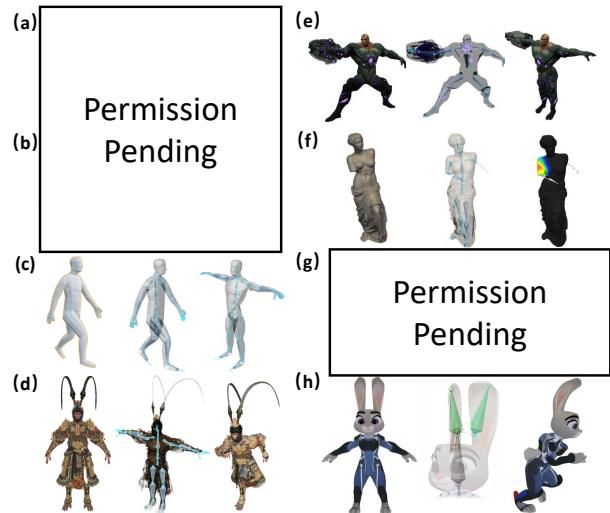


Figure 7. **Results of more cases to demonstrate the advantage of our method.** The detailed explanations can be found in the supplementary material.

4.2. Comparison Results

In Fig. 4, we compare the proposed method with two commercial software, Meshy [36] and Tripo [53]. Meshy tends to enlarge the influence regions of skinning weights, leading to incorrect deformations of clothes or muscles. It also fails to extract the rest pose of some complex input shapes and animates the characters in a weird way. Tripo does better in blend weights, but often produces inaccurate and unmatched bones. Plus, it also has difficulty determining the correct rest state for unusual input poses. Moreover, both of the two cannot predict finger bones, resulting in coarse hand poses in the running scenario. In contrast, our method produces complete and accurate bones as well as suitable blend weights, so the animations are also natural and fluid.

In Fig. 5, we evaluate the auto-rigging method

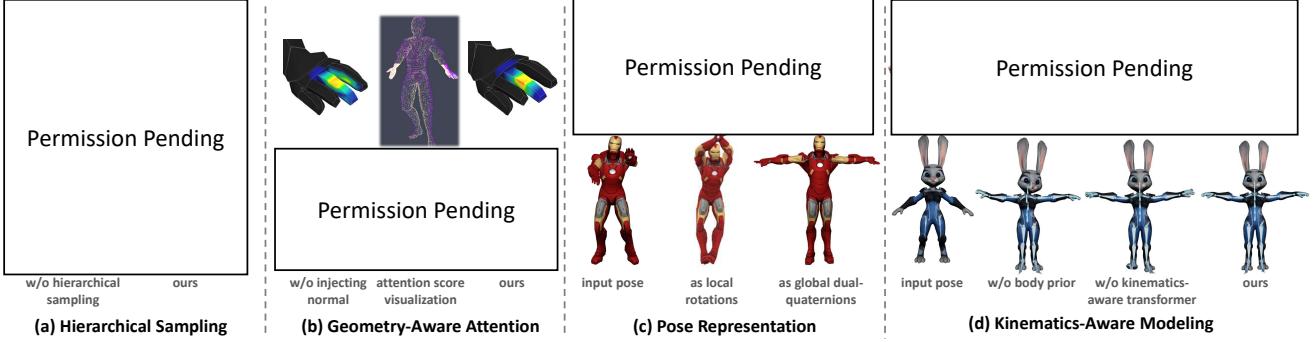


Figure 8. **Visualizations of some ablative experiments.** We show the effectiveness of the proposed modules and design choices by visualizing the predicted bones, blend weights, and pose transformations.

RigNet [57] with given test meshes. It can be observed that RigNet often generates unnecessary bone structures, *e.g.*, a cycle graph (the old man case) and a star topology (the bear case). These bone errors are also manifested in Tab. 2. Besides, its weights usually fail to maintain spatial smoothness, resulting in fragmented animating results.

In Fig. 6, we further compare two generative works that perform text-to-3D-avatar generation. TADA [32] produces textured meshes and HumanGaussian [35] chooses to use Gaussian splats. Our framework is directly compatible with both their representations. Since these two works generate shapes based on the well-defined SMPL [37] template, their joints’ locations and weights are good enough for animating. Nevertheless, our method can achieve comparable, if not superior, quality using a template-free pipeline, without any prior knowledge of the character’s shape or pose parameters. Furthermore, although these methods exhibit some deviations from the template mesh, they remain constrained by the preset body ratio of SMPL and are limited to shapes resembling realistic humans. In contrast, our method can be applied to a much wider range of character shapes.

We show more results produced by our method, focusing on the details and some tricky cases. Each sub-figure of Fig. 7 proves some advantages of our methods, which is interpreted one by one as follows. (a) Fine-grained control of fingers. (b) Capacity of abnormal shapes. (c) Complex input poses. (d) Efficiency for high polygon models. (e) Support for asymmetric inputs. (f) Adaptation to non-existing bones. (g),(h) Extension to extra bones.

4.3. Ablation Study

We conduct ablative experiments to validate the effectiveness of the proposed components and strategies, as shown in Tab. 3 and Fig. 8. It can be observed that the canonical transformation is vital for the proper convergence of our model, as it greatly reduces the distribution dispersion of inputs. The large improvement brought by hierarchical sampling also indicates the importance of the proposed coarse-to-fine shape representation. As illustrated in Fig. 8 (a), sampling more in the fine-grained hand regions leads to a more ac-

	IoU ↑	Precision ↑	Recall ↑	CD-J2J ↓	CD-J2B ↓	CD-B2B ↓
RigNet [57]	53.50%	47.27%	89.30%	6.63%	4.97%	2.88%
Ours	82.50%	81.07%	90.31%	4.49%	3.32%	1.58%

Table 2. **Quantitative comparison of skeleton prediction on the bipedal humanoid subset of the test dataset [56].**

	Weights Error ↓	Joints Error ↓	Poses Error ↓
w/o canonical transformation	6.27%	9.80%	41.8%
w/o hierarchical sampling	5.55%	2.42%	18.0%
w/o geometry-aware attention	5.16%	2.20%	14.2%
w/o structure-aware transformer	-	2.13%	14.9%
w/o body prior loss	-	2.13%	14.0%
w/o global pose representation	-	-	35.3%
Ours	4.70%	2.11%	13.6%

Table 3. **Ablation studies on the test split of the Mixamo dataset.** We report the percentage error of animation assets.

curate estimation of finger bones. Besides, the injection of additional normal information via geometry-aware attention boosts the performance, especially in weight prediction. We visualize the attention score of the input point clouds in Fig. 8 (b), where brighter color indicates more attention to normals rather than coordinates. The model adaptively learns to rely more on normals in regions like inner thigh and between fingers since coordinates become less discriminative there. Furthermore, accurately modeling pose transformations presents a complex challenge. When poses are predicted as local rotations rather than global dual quaternions, the model is limited to handling simpler inputs. Without a body prior loss, predicted limb and spine bones often exhibit slight but critical offsets, disrupting connectivity and symmetry. Additionally, in the absence of a structure-aware transformer, deeper-level bones, such as those in the fingers and feet, are susceptible to inaccuracies in pose-to-rest transformations. Their poses are heavily influenced by ancestral bones, where even minor inconsistencies can result in significant deformities.

5. Conclusion and Discussion

In this paper, we propose a novel framework for animation-ready 3D character production. To address non-trivial challenges and practical limitations of existing methods, we de-

velop several elaborate modules, including the coarse-to-fine shape representation, the particle-based shape autoencoder, and the structure-aware modeling of bones. Putting all these together, we provide an out-of-the-box and efficient solution to animating any 3D character. Comprehensive experiments demonstrate the superiority of our method and its considerable potential for future investigation.

Despite the merits, there is still room for improvement. Making the empirical hierarchical sampling more adaptive may be a meaningful future direction. It is also a promising avenue to extend the proposed framework to non-bipedal characters by proposing a more flexible way of modeling the topological structures of bones.

References

- [1] Adobe. Mixamo, 2024. <https://www.mixamo.com>. 2, 6, 1, 3, 5
- [2] David Baraff and Andrew Witkin. Large steps in cloth simulation. In *ACCGI*, page 43–54, New York, NY, USA, 1998. Association for Computing Machinery. 2
- [3] Ilya Baran and Jovan Popović. Automatic rigging and animation of 3D characters. *ACM TOG*, 26(3):72–es, 2007. 3
- [4] Jernej Barbic. *Real-time reduced large-deformation models and distributed contact for computer graphics and haptics*. PhD thesis, Carnegie Mellon University, 2007. 2
- [5] Jernej Barbič and Doug L James. Real-time subspace integration for st. venant-kirchhoff deformable models. *ACM TOG*, 24(3):982–990, 2005.
- [6] Jernej Barbič and Yili Zhao. Real-time large-deformation substructuring. *ACM TOG*, 30(4):1–8, 2011. 2
- [7] Adam W Bargteil, Chris Wojtan, Jessica K Hodgins, and Greg Turk. A finite element method for animating large viscoplastic flow. *ACM TOG*, 26(3):16–es, 2007. 2
- [8] Miklós Bergou, Max Wardetzky, Stephen Robinson, Basile Audoly, and Eitan Grinspun. Discrete elastic rods. In *ACM SIGGRAPH*, New York, NY, USA, 2008. Association for Computing Machinery. 2
- [9] Sofien Bouaziz, Sebastian Martin, Tiantian Liu, Ladislav Kavan, and Mark Pauly. Projective dynamics: fusing constraint projections for fast simulation. *ACM TOG*, 33(4), 2014. 2
- [10] Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 5
- [11] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. ShapeNet: An information-rich 3D model repository. *arXiv preprint arXiv:1512.03012*, 2015. 4
- [12] Yue Chang, Peter Yichen Chen, Zhecheng Wang, Maurizio M Chiaramonte, Kevin Carlberg, and Eitan Grinspun. LiCROM: Linear-subspace continuous reduced order modeling with neural fields. In *SIGGRAPH Asia*, pages 1–12, 2023. 2
- [13] Peter Yichen Chen, Jinxu Xiang, Dong Heon Cho, Yue Chang, GA Pershing, Henrique Teles Maia, Maurizio M Chiaramonte, Kevin Carlberg, and Eitan Grinspun. CROM: Continuous reduced-order modeling of PDEs using implicit neural representations. *arXiv preprint arXiv:2206.02607*, 2022. 2
- [14] Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, and Gang Yu. Executing your commands via motion diffusion in latent space. In *CVPR*, pages 18000–18010, 2023. 1
- [15] Stelian Coros, Sebastian Martin, Bernhard Thomaszewski, Christian Schumacher, Robert Sumner, and Markus Gross. Deformable objects alive! *ACM TOG*, 31(4):1–9, 2012. 2
- [16] Barbara Cutler, Julie Dorsey, Leonard McMillan, Matthias Müller, and Robert Jagnow. A procedural approach to authoring solid models. *ACM TOG*, 21(3):302–311, 2002. 2
- [17] Devikalyan Das, Christopher Wewer, Raza Yunus, Eddy Ilg, and Jan Eric Lenssen. Neural parametric gaussians for monocular non-rigid object reconstruction. In *CVPR*, pages 10715–10725, 2024. 1
- [18] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihns, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3D objects. In *CVPR*, pages 13142–13153, 2023. 6
- [19] Mathieu Desbrun and Marie-Paule Gascuel. Animating soft substances with implicit surfaces. In *ACCGI*, page 287–290, New York, NY, USA, 1995. Association for Computing Machinery. 2
- [20] Yixin Hu, Qingnan Zhou, Xifeng Gao, Alec Jacobson, Dennis Zorin, and Daniele Panozzo. Tetrahedral meshing in the wild. *ACM TOG*, 37(4):60, 2018. 2
- [21] Yukun Huang, Jianan Wang, Ailing Zeng, Zheng-Jun Zha, Lei Zhang, and Xihui Liu. DreamWaltz-G: Expressive 3D gaussian avatars from skeleton-guided 2D diffusion. *arXiv preprint arXiv:2409.17145*, 2024. 3
- [22] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. SC-GS: Sparse-controlled gaussian splatting for editable dynamic scenes. In *CVPR*, 2024. 3, 1
- [23] Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. MotionGPT: Human motion as a foreign language. *NeurIPS*, 36, 2024. 1
- [24] Chenfanfu Jiang, Craig Schroeder, Joseph Teran, Alexey Stomakhin, and Andrew Selle. The material point method for simulating continuum materials. In *ACM SIGGRAPH 2016 Courses*, New York, NY, USA, 2016. Association for Computing Machinery. 2
- [25] Tao Ju, Qian-Yi Zhou, Michiel Van De Panne, Daniel Cohen-Or, and Ulrich Neumann. Reusable skinning templates using cage-based deformations. *ACM TOG*, 27(5):1–10, 2008. 2
- [26] Jongmin Kim, Yeongho Seol, Taesoo Kwon, and Jehee Lee. Interactive manipulation of large-scale crowd animation. *ACM TOG*, 33(4):1–10, 2014. 2
- [27] Jeonghwan Kim, Hyeontae Son, Jinseok Bae, and Young Min Kim. Auto-rigging 3D Bipedal Characters in Arbitrary Poses. In *Eurographics - Short Papers*. The Eurographics Association, 2021. 3
- [28] Jiahui Lei, Yijia Weng, Adam Harley, Leonidas Guibas, and Kostas Daniilidis. MoSca: Dynamic gaussian fusion

- from casual videos via 4D motion scaffolds. *arXiv preprint arXiv:2405.17421*, 2024. 3, 1
- [29] J. P. Lewis, Matt Cordner, and Nickson Fong. Pose space deformation: a unified approach to shape interpolation and skeleton-driven deformation. In *ACCGI*, page 165–172, USA, 2000. ACM Press/Addison-Wesley Publishing Co. 2
- [30] Peizhuo Li, Kfir Aberman, Rana Hanocka, Libin Liu, Olga Sorkine-Hornung, and Baoquan Chen. Learning skeletal articulations with neural blend shapes. *ACM TOG*, 40(4):1–15, 2021. 2, 3, 1, 5
- [31] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM TOG*, 36(6):194:1–194:17, 2017. 1
- [32] Tingting Liao, Hongwei Yi, Yuliang Xiu, Jiaxiang Tang, Yangyi Huang, Justus Thies, and Michael J. Black. TADA! Text to animatable digital avatars. In *3DV*, pages 1508–1519, 2024. 2, 3, 7, 8, 4, 5
- [33] Anything World Limited. 3D animation and automated rigging | Anything World, 2024. <https://anything.world>. 2, 5, 6
- [34] Jing Lin, Ailing Zeng, Shunlin Lu, Yuanhao Cai, Ruimao Zhang, Haoqian Wang, and Lei Zhang. Motion-X: A large-scale 3D expressive whole-body human motion dataset. *NeurIPS*, 2023. 1
- [35] Xian Liu, Xiaohang Zhan, Jiaxiang Tang, Ying Shan, Gang Zeng, Dahua Lin, Xihui Liu, and Ziwei Liu. HumanGaussian: Text-Driven 3D Human Generation with Gaussian Splatting. In *CVPR*, pages 6646–6657, 2024. 2, 3, 7, 8, 4, 5
- [36] Meshy LLC. Meshy - convert text and images to 3D models, 2024. <https://www.meshy.ai>. 2, 6, 7, 4, 5
- [37] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM TOG*, 34(6):248:1–248:16, 2015. 3, 5, 7, 8, 1, 4
- [38] Jing Ma and Dongliang Zhang. TARig: Adaptive template-aware neural rigging for humanoid characters. *Computers & Graphics*, 114:158–167, 2023. 2, 3, 5, 6
- [39] Thalmann Magnenat, Richard Lapierre, and Daniel Thalmann. Joint-dependent local deformations for hand animation and object grasping. In *Proceedings of Graphics Interface*, pages 26–33. Canadian Inf. Process. Soc, 1988. 2
- [40] Vismay Modi, Nicholas Sharp, Or Perel, Shinjiro Sueda, and David I. W. Levin. Simplicits: Mesh-free, geometry-agnostic elastic simulation. *ACM TOG*, 43(4), 2024. 2
- [41] Joe J Monaghan. Smoothed particle hydrodynamics. *Reports on Progress in Physics*, 68(8):1703, 2005. 2
- [42] Matthias Müller, Richard Keiser, Andrew Nealen, Mark Pauly, Markus Gross, and Marc Alexa. Point based animation of elastic, plastic and melting objects. In *Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 141–151, 2004. 2
- [43] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3D hands, face, and body from a single image. In *CVPR*, pages 10975–10985, 2019. 3, 7, 1
- [44] Andreas Peer, Christoph Gissler, Stefan Band, and Matthias Teschner. An implicit sph formulation for incompressible linearly elastic solids. In *Computer Graphics Forum*, pages 135–148. Wiley Online Library, 2018. 2
- [45] pixiv Inc. VRoid Studio, 2024. <https://vroid.com/en/studio>. 6, 3
- [46] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. DreamFusion: Text-to-3D using 2D diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 3
- [47] Charles R. Qi, Li Yi, Hao Su, and Leonidas J. Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *NeurIPS*, 2017. 4
- [48] Thomas W Sederberg and Scott R Parry. Free-form deformation of solid geometric models. In *ACCGI*, pages 151–160, 1986. 2
- [49] Eftychios Sifakis and Jernej Barbic. Fem simulation of 3D deformable solids: a practitioner’s guide to theory, discretization and model reduction. In *ACM SIGGRAPH 2012 Courses*, New York, NY, USA, 2012. Association for Computing Machinery. 2
- [50] Alexey Stomakhin, Craig Schroeder, Lawrence Chai, Joseph Teran, and Andrew Selle. A material point method for snow simulation. *ACM TOG*, 32(4):1–10, 2013. 2
- [51] Daniel Strötter, Jean-Marc Thiery, Kai Hormann, Jiong Chen, Qingjun Chang, Sebastian Besler, Johannes Sebastian Mueller-Roemer, Tammy Boubekeur, André Stork, and Dieter W Fellner. A survey on cage-based deformation of 3D models. In *Computer Graphics Forum*, page e15060. Wiley Online Library, 2024. 2
- [52] Demetri Terzopoulos, John Platt, Alan Barr, and Kurt Fleischer. Elastically deformable models. In *ACCGIT*, pages 205–214, 1987. 2
- [53] Tripo. Tripo AI, 2024. <https://www.trip03d.ai>. 2, 6, 7, 4, 5
- [54] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 5
- [55] Qianqian Wang, Vickie Ye, Hang Gao, Jake Austin, Zhengqi Li, and Angjoo Kanazawa. Shape of motion: 4D reconstruction from a single video. *arXiv preprint arXiv:2407.13764*, 2024. 1
- [56] Zhan Xu, Yang Zhou, Evangelos Kalogerakis, and Karan Singh. Predicting animation skeletons for 3D articulated models via volumetric nets. In *3DV*, pages 298–307. IEEE, 2019. 6, 7, 8
- [57] Zhan Xu, Yang Zhou, Evangelos Kalogerakis, Chris Landreth, and Karan Singh. RigNet: Neural rigging for articulated characters. *ACM TOG*, 39(4):58:1–58:14, 2020. 2, 3, 7, 8, 5, 6
- [58] Biao Zhang, Jiapeng Tang, Matthias Nießner, and Peter Wonka. 3DShape2VecSet: A 3D shape representation for neural fields and generative diffusion models. *ACM TOG*, 42(4):92:1–92:16, 2023. 4, 3
- [59] Wojciech Zienonka, Timo Bolkart, Thabo Beeler, and Justus Thies. Gaussian eigen models for human heads. *arXiv preprint arXiv:2407.04545*, 2024. 1

Make-It-Animatable: An Efficient Framework for Authoring Animation-Ready 3D Characters

Supplementary Material

S1. Formulation of Low-Rank Dynamics

As discussed in Sec. 3.1, modeling the dynamics of an object typically requires a per-timestamp deformation of all the particles (*i.e.*, vertices for meshes, points for point clouds, and splats for 3D Gaussians) that make up its geometry. Suppose that we have an object composed of N particles and a desired dynamic sequence of T timestamps. The temporal deformations of all particles can be represented by a matrix $\mathbf{D} \in \mathbb{R}^{T \times N \times d}$, where d is the degrees of freedom (*e.g.*, $d = 6$ for rigid transformations). Generally, \mathbf{D} has a lot of redundancy when representing real-world dynamics. Therefore, we would like to seek a low-rank approximation of it. Mathematically, the low-rank decomposition of a matrix like \mathbf{D} has two different forms expressed by

$$\mathbf{D}^{T \times N \times d} \approx \mathbf{B}_t^{T \times K \times d} \mathbf{W}_s^{K \times N}, \quad (\text{S1})$$

$$\mathbf{D}^{T \times N \times d} \approx \mathbf{B}_s^{K \times N \times d} \mathbf{W}_t^{T \times K}, \quad (\text{S2})$$

where K is the desired rank. We call the matrix \mathbf{B} with dimension d as *basis*, and the other one \mathbf{W} as *weight*. Eqs. (S1) and (S2) both decouple the temporal and spatial dimensions by splitting them into basis and weight. We then name the ones with temporal dimension (T) as *temporal basis* \mathbf{B}_t and *temporal weight* \mathbf{W}_t . Correspondingly, the matrices with spatial dimension (N) are named as *spatial basis* \mathbf{B}_s or *spatial weight* \mathbf{W}_s .

In fact, both decomposition forms have been playing important roles in the applications of dynamic modeling. For example, Eq. (S1) is applied in the sparse control paradigms [22, 28, 55] and the linear blend skinning (LBS) [37, 43] algorithm, where \mathbf{B}_t is interpreted as the transformations/poses of control nodes/body joints, and \mathbf{W}_s is the blend weights of the nodes or joints. When applying Eq. (S2), \mathbf{B}_s and \mathbf{W}_t are interpreted as the blend shapes and their corresponding weights for the per-timestamp linear combination [17, 30, 31, 59].

In this work, we focus on the low-rank form of Eq. (S1) in modeling dynamic characters. The reason is twofold. First, the temporal weight \mathbf{W}_t in Eq. (S2) typically requires as much motion data as possible for one object to find a capacious low-rank space, which is prohibitively difficult in practice. In contrast, the spatial weight \mathbf{W}_s can be supervised much more easily by the well-defined blend weights in existing 3D models. Second, considering the availability of rich motion resources [1, 34] and powerful motion generation methods [14, 23], there is no need to model the temporal basis \mathbf{B}_t from scratch. What we have to do is

finding a transformation from the input character pose to a pre-defined rest pose, and then any desired animation is within easy reach.

Note that theoretically, joint/bone positions are only proxies or interfaces for the low-rank terms and are not indispensable to the dynamic modeling, as expressed by Eq. (S1). Consequently, some related works [55] choose to model such proxies in an implicit way. However, we still include the bones as one of the desired animation assets in our work for a self-contained and artist-friendly representation compatible with existing animating pipelines. Furthermore, the explicit existence of bones can bring much convenience when applying body priors to assist the optimization, as introduced in Sec. 3.4.

S2. Implementation Details

S2.1. Coarse-to-Fine Shape Representation

To address the limitations of the proposed framework and boost the performance, we introduce an additional design at the input side of the autoencoder, *i.e.*, the coarse-to-fine shape representation (Sec. 3.3). In addition to the overall framework presented in Fig. 2, we include two separate pipelines in Fig. S1 here for a clearer illustration of the coarse and fine training stages. In the following content, we will provide more details of this part, particularly regarding the motivations and implementation choices.

With the particle-based shape autoencoder, our method can already produce fairly good results of blend weights. However, the joint outputs are still unsatisfying in fine-grained regions like the hands. Meanwhile, the pose prediction can hardly converge. We attribute these issues to the ambiguity of the input points and treat the lite training process (taking uniformly sampled points as input and only predicting bone positions, as shown in Fig. S1 upper) as a coarse stage that provides rough but valuable localization information about the input character. To be specific, we exploit some of the coarse joint locations to gain a finer shape representation by applying two different strategies, *i.e.*, the canonical transformation and the hierarchical sampling.

Canonical transformation. Although we have normalized the input shapes to align their scales, they still differ in global transformations. Since the poses we want to predict are relative to the fixed origin, the same pose can have huge numerical differences under different coordinate systems. This will lead to a dramatically increased difficulty in

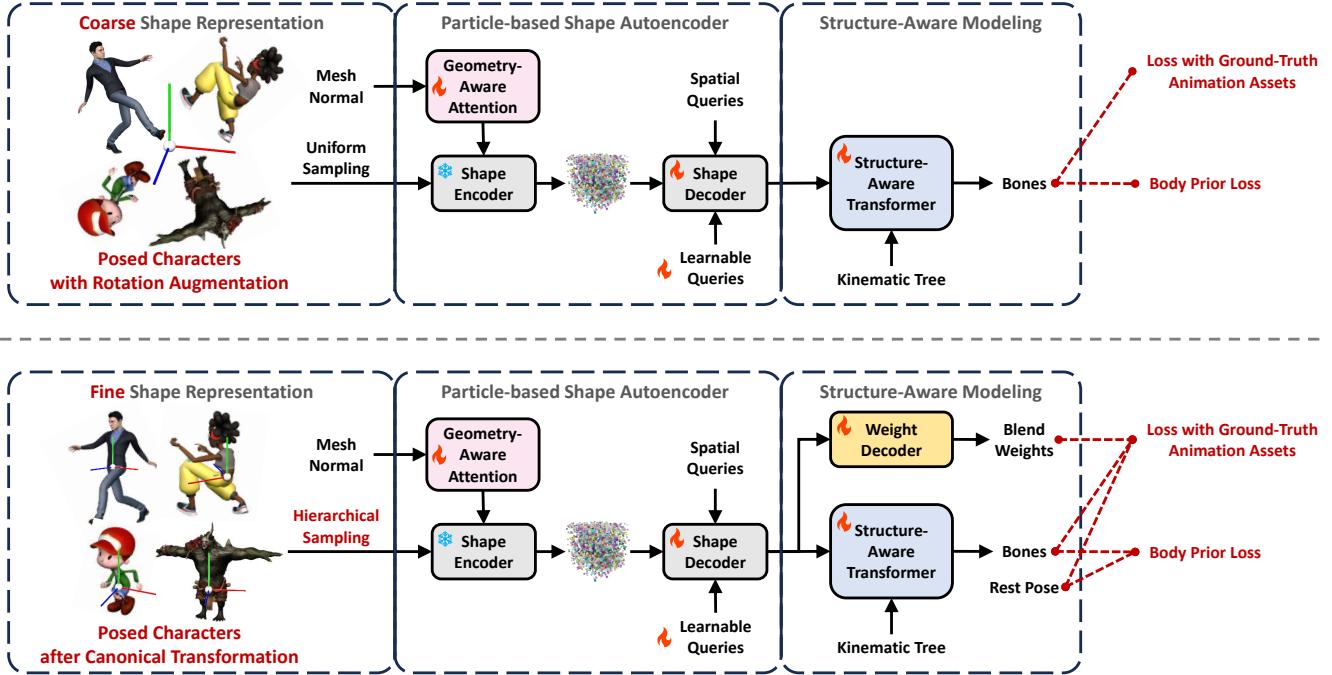


Figure S1. **The coarse (upper) and fine (lower) stages of training our framework.** In the coarse stage, the input shape is uniformly sampled and only the bone positions are predicted. We apply data augmentation to the inputs via random 3D rotations, so that the coarse model is generalizable to global transformations of in-the-wild cases with an acceptable accuracy. In the fine stage, we apply canonical transformation and hierarchical sampling to the shapes in advance based on the ground-truth bone positions. Then during inference, a 3D character is fed into the coarse framework to get its bone positions, which guide the establishment of coarse-to-fine shape representation later in the fine framework. Note that the body prior loss (Sec. 3.4) is directly applied to the bone positions. As for pose prediction, we take the ground-truth bones as a proxy and use the predicted pose to transform them, thereby indirectly affecting the pose optimization.

pose prediction. Therefore, we move and rotate the entire shape to a canonical position and orientation before sampling. In practice, we choose this transformation based on an empirically selected datum plane (referred to as the *hip plane*) determined by three joint positions, *i.e.*, the hip (root of the kinematic tree) and two upper thighs, which can be accurately located even in the coarse stage. The canonical transformation is applied so that: 1) the hip is located at the origin; 2) the normal of the hip plane is aligned with the z-axis; 3) the vector from the right to the left thigh is parallel to the x-axis. This process simplifies the input shape’s spatial distribution and eliminates the pose representation’s ambiguity. Furthermore, the chosen canonical transformation ensures a consistent upright and front-facing orientation of the input character, which is a common prerequisite of many auto-rigging methods [30, 38, 57]. By integrating the coarse localization, we now automate this preprocessing step, enabling our framework to effectively handle the inputs regardless of their initial spatial configuration (positions, rotations, and scales).

Hierarchical sampling. Some parts of the input character, *e.g.*, the hands, present fine-grained details within small

regions, which typically demand additional resources for accurate processing. However, the uniform sampling applied to the input shape, as well as the subsequent farthest point sampling (FPS) algorithm, usually results in sparsely distributed sample points on hands, which are far from enough to describe the geometry of fingers. Theoretically, increasing the number of sampling points N and the down-sampling number M can bring a larger representation capacity for the entire geometry including the hand regions, but that will significantly add to the computational overhead. For efficiency, we choose to keep the total sampling number unchanged and instead leverage the coarse prediction of hand joints. Specifically, we replace the uniform sampling with a hierarchical approach that ensures a designated proportion of sample points are distributed on both hands. Since the FPS algorithm in the shape encoder disrupts the non-uniformity of samples, we adapt it to a hierarchical algorithm as well.

S2.2. Networks and Hyperparameters

For the shape autoencoder, we use a point cloud of size $N = 32768$ as the input. All the points are sampled on the surface of the input mesh. In the fine training or the in-

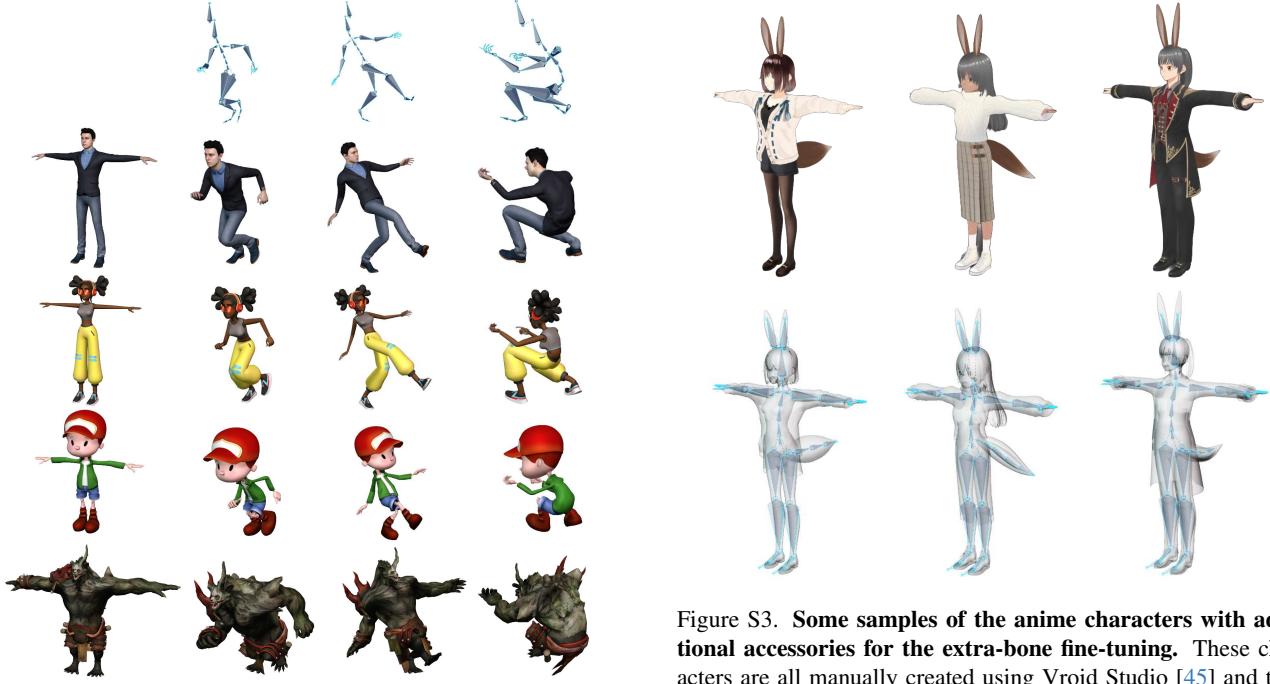


Figure S2. **Some samples from the collected Mixamo [1] dataset.** The dataset contains bipedal humanoids with different shapes, ranging from realistic humans to cartoon or fantasy creatures. Each character is preprocessed to be animatable by any of the motion sequences. The proposed framework is trained on this dataset.

ference stage equipped with the hierarchical sampling, 50% of the points are uniformly sampled and the other 50% are sampled near the hand joints. The hyperparameters of the shape autoencoder are consistent with the original setting of 3DShape2VecSet [58], which internally uses $M = 512$ latents (each of them is a 512-D vector) for the neural field. Our framework is trained on 8 NVIDIA A100 GPUs. The learning rate is linearly increased to $1e - 4$ within the first 1% iterations (warm-up), and then gradually decreased using the cosine decay schedule until reaching the minimum value of $1e - 5$.

S3. Data Details

S3.1. Mixamo Dataset

To obtain sufficient training data of 3D characters with high-quality geometry and animation assets, we collect the artist-designed 3D models from Mixamo [1] to form a dataset, which comprises texture meshes of 95 exquisite characters (each has an average of 15000 vertices), along with 2453 diverse motion sequences (each has an average of 200 frames). All the characters are preprocessed to share a standard skeleton structure with $K = 52$ bones (the leaf bones are removed). For characters with non-standard skeletal

Figure S3. **Some samples of the anime characters with additional accessories for the extra-bone fine-tuning.** These characters are all manually created using Vroid Studio [45] and then preprocessed to be compatible with the standard skeleton definition of Mixamo [1].

structures originally, the blend weights of those bones are transferred to their topologically nearest ancestral bones within the standard skeleton. If some standard bones are missing in a character (*e.g.*, armless person), we mask their corresponding values (weight channels and bone-wise attributes) in loss computing. We use 95% of the data for training and the remaining 5% for validation. During each training iteration, we randomly choose a character-motion pair to get the corresponding animation assets, resulting in an effective dataset size equivalent to over 40 million frames. In Fig. S2, we show some example characters and motions from the Mixamo dataset.

S3.2. Vroid Dataset

Our training framework can also be extended to different skeleton topologies. To demonstrate this capacity, we use Vroid Studio [45] to manually create 35 different anime characters (30 for training and 5 for validation) with additional accessories including two rabbit-like ears and a fox-like tail, as shown in Fig. S3. These characters differ in body shape, clothes, hair style, and the shape of ears and tails. They are all preprocessed to share the same skeleton definition of Mixamo [1] but with extra bones of the accessories. Therefore, they can also be animated with any of the Mixamo motion sequences during the extra-bone training (*i.e.*, fine-tuning the final layer of the weight decoder and the extra learnable queries, based on the standard-skeleton model pretrained on the Mixamo dataset). The experiments

Permission Pending

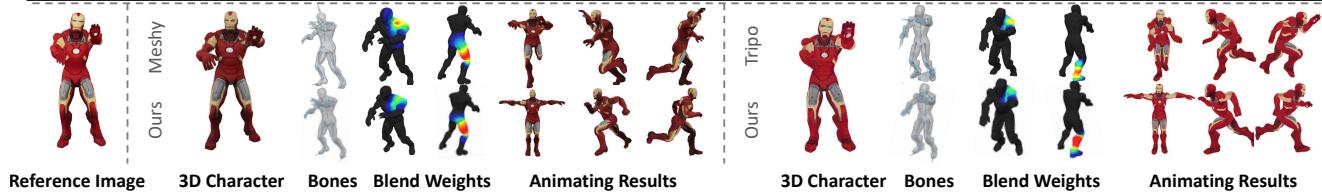


Figure S4. **Additional Comparison with generative 3D methods, i.e., Meshy [36] and Tripo [53].** We feed them the same image as reference and compare the performance based on their generated 3D models respectively. The blend weights of two joints, i.e., Left Shoulder and Right Leg, are visualized. The T-pose models predicted by our method are included as the front-view animating results. Zoom in to better view the details.

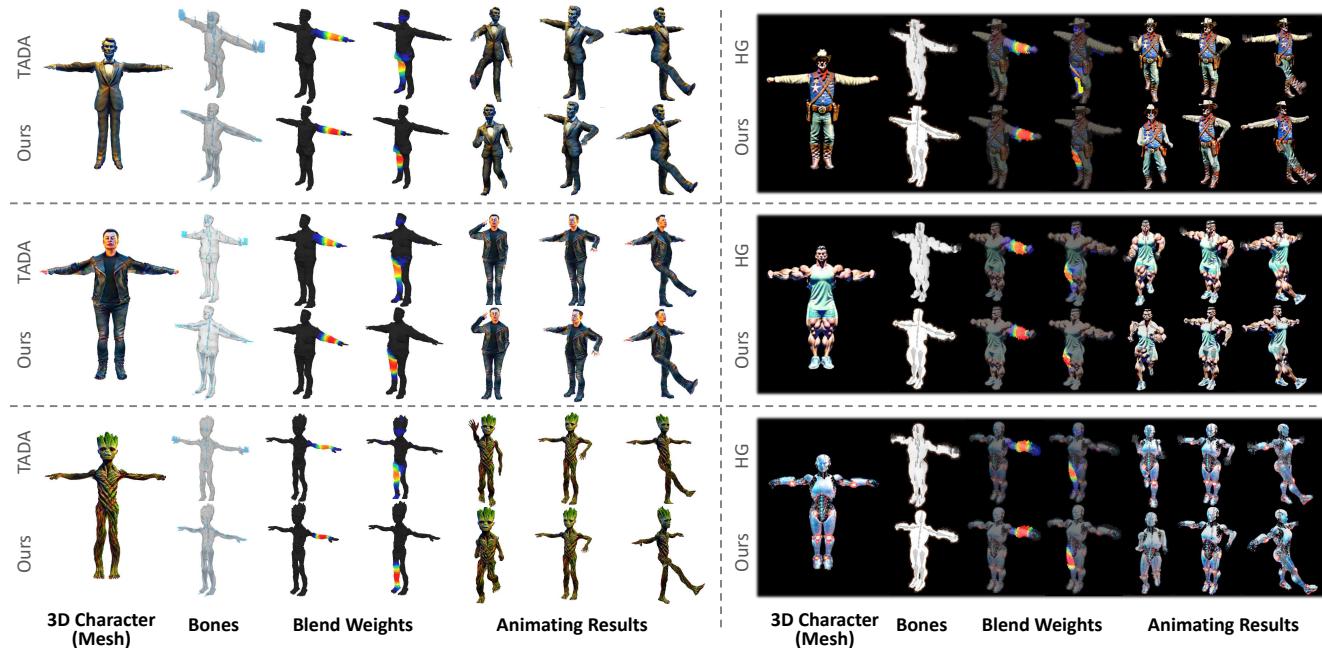


Figure S5. **Additional comparison with template-based avatar generation methods, i.e., TADA [32] and HumanGaussian [35] (HG).** We use the generated meshes from TADA and 3D Gaussians from HG for comparison. Note that the skeletons of these two baselines are identical to the shape-specific SMPL [37] templates (without bone tail), with their weights interpolated from the template meshes. Zoom in to better view the details.

show that with our framework, 30 training characters are sufficient to obtain a good model that produces promising predictions for extra bones.

S4. More Results

S4.1. Additional Comparison Results

More comparison cases. In addition to Figs. 4 and 6, we exhibit more cases of comparison with the baselines, includ-

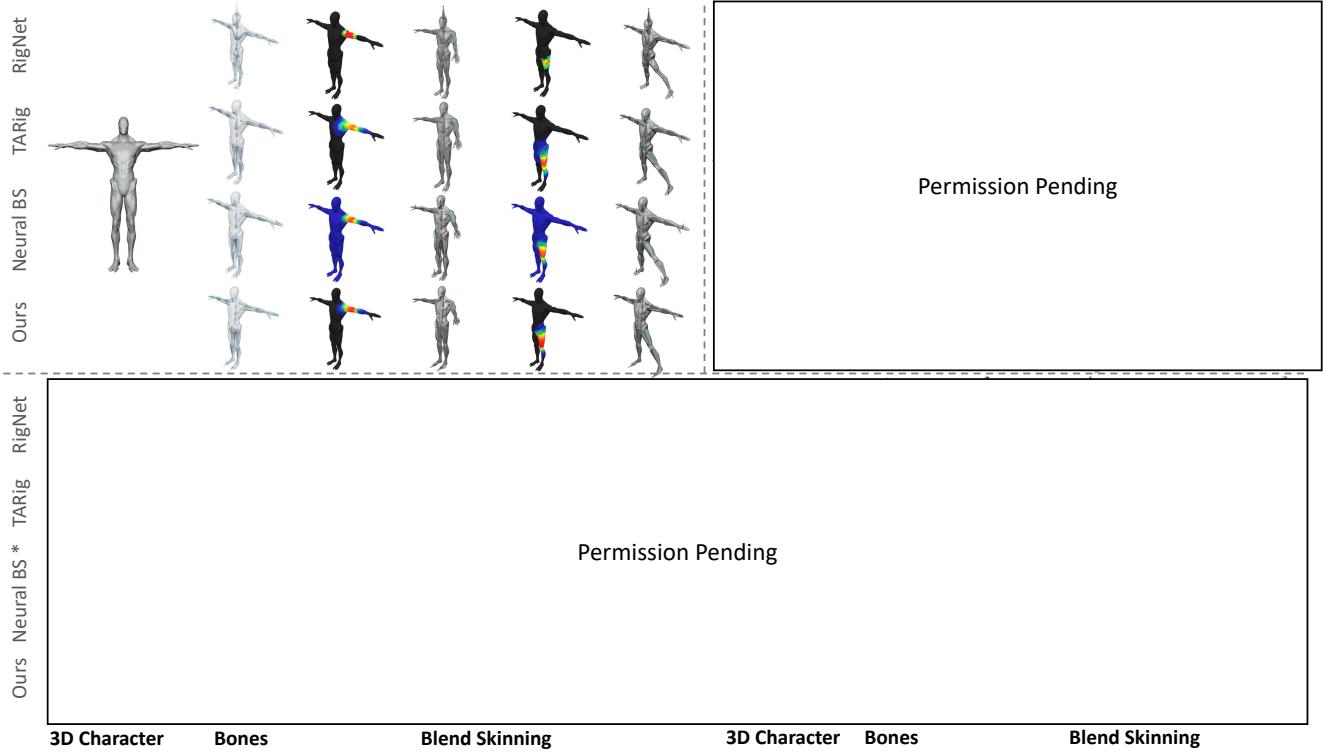


Figure S6. **Additional comparison with auto-rigging methods, i.e., RigNet [57], TARig [38], and Neural Blend Shapes [30] (Neural BS).** We visualize the blend weights of selected joints and manually deform them to assess the impact of rigging quality on skinning results. *: Neural Blend Shapes only support T-pose inputs, so for the non-rest cases (lower two), we feed it the T-pose meshes transformed by our pose-to-rest predictions.

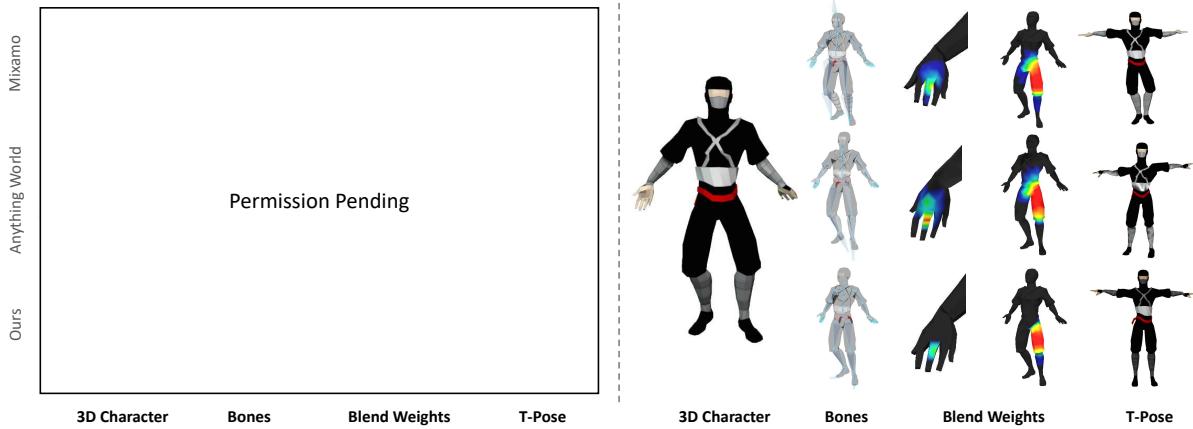


Figure S7. **Comparison with commercial auto-rigging software, i.e., Mixamo [1] and Anything World [33].** Note that these two tools can only deal with simple input poses (T- or A-pose is recommended) and often raise errors when faced with complex ones.

ing Meshy [36] and Tripo [53] (Fig. S4), TADA [32] and HumanGaussian [35] (Fig. S5). The results can prove the effectiveness, robustness, and generalizability of the proposed framework.

Comparison with more auto-rigging methods. For comparison with existing auto-rigging methods, we include more cases and compare with two more baselines here in Fig. S6 (in addition to Fig. 5), i.e., Neural Blend Shapes [30] and TARig [38]. Neural Blend Shapes only supports T-pose inputs and has quite limited generalizability to shapes dif-

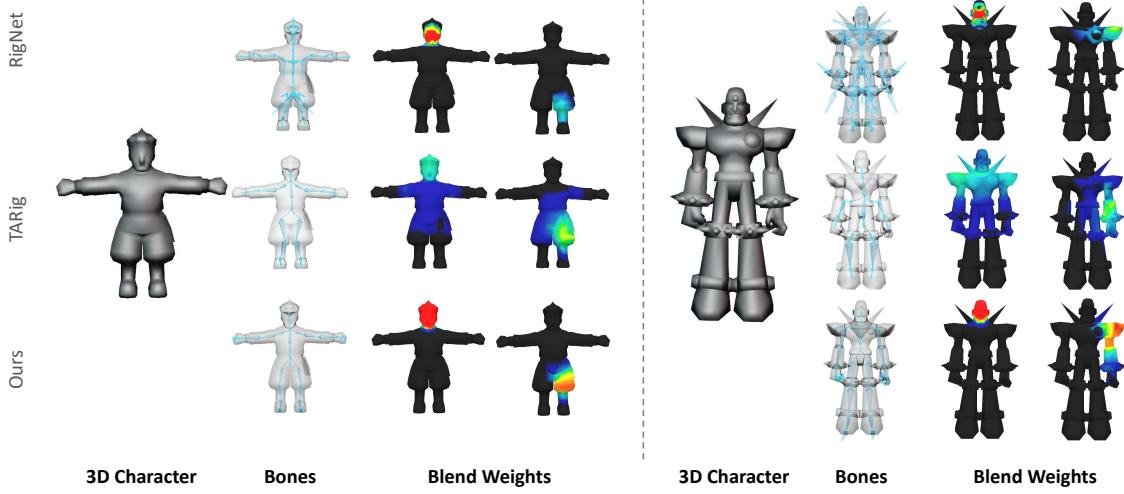


Figure S8. **Qualitative comparison with RigNet [57] and TARig [38] on cases from the test split of “ModelsResource-RigNetv1” dataset [56].** While both baselines are exactly trained on this dataset and ours are not, we still achieve the best performance.

ferent from the SMPL mesh. TARig’s skeleton predictions are better than Neural Blend Shapes and RigNet [57], but its blend weights still cannot meet the standard of practical application, producing spatial unsMOOTHNESS when deforming the meshes. Besides, all three baselines are unable to produce fine-grained hand bones, while our method handles the fingers well.

Furthermore, the commercial software, Mixamo [1] and Anything World [33], rely much on the symmetry or pose simplicity (*e.g.*, T-pose and A-pose) of the inputs and will raise errors when faced with complex ones. Therefore, we compare them separately on some additional cases here in Fig. S7. It can be observed that Anything World produces significant errors when extracting the bones of the left arm for the tiger (left case). Meanwhile, Mixamo fails to fail to distinguish the left and right sides of the ninja (right case) and produces a mirrored skeleton. Moreover, when faced with a non-rest input character like this ninja, the predicted pose-to-rest transformations of Mixamo and Anything World both suffer from unnatural deformations, while our results remain good.

Qualitative comparison on the ModelsResource dataset [56]. We also evaluate the proposed framework on the bipedal-humanoid subset of the “ModelsResource-RigNetv1” dataset [56]. Fig. S8 shows some cases for qualitative comparison with RigNet [57] and TARig [38]. Both of these two baselines are trained on the aforementioned dataset, but our model has never encountered a similar data distribution during training. Despite this, our method still achieves the best quality in rigging and skinning, demonstrating its strong generalizability.

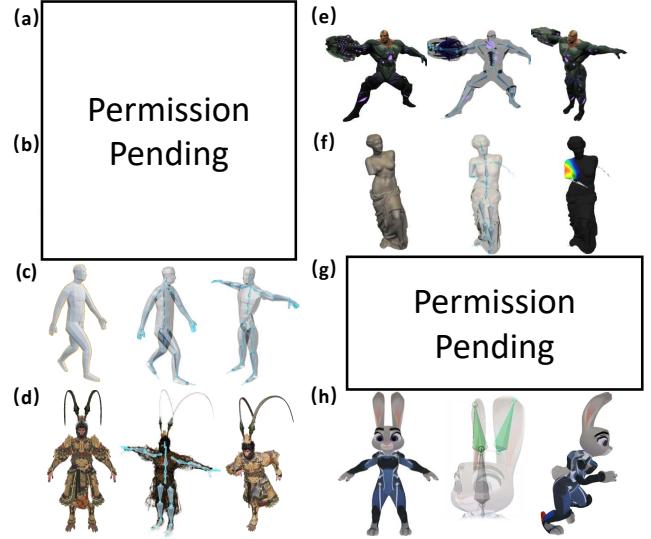


Figure S9. **Results of more cases to demonstrate the advantage of our method.** (a) Fine-grained control of fingers; (b) Capacity of abnormal shapes; (c) Complex input poses; (d) Efficiency for high polygon models; (e) Support of asymmetric inputs; (f) Adaptation to non-existing bones; (g) & (h): Extension to extra bones (*e.g.*, long ears and tails).

S4.2. Additional Visualizations

Detailed explanation of Fig. 7. We provide detailed explanation and analysis for each case in Fig. 7 (also included here in Fig. S9 for quick reference) as follows.

(a) Fine-grained control of fingers. Thanks to the coarse-to-fine shape representation (Sec. 3.3), our method shows remarkable accuracy at the hand regions, which is difficult for most existing approaches.

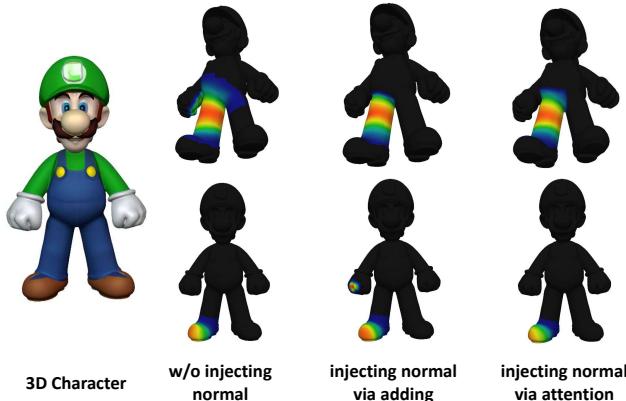


Figure S10. **Qualitative analysis of our geometry-aware attention module and its injecting method.** The proposed attention-based injection can benefit from normal information without any side effects.

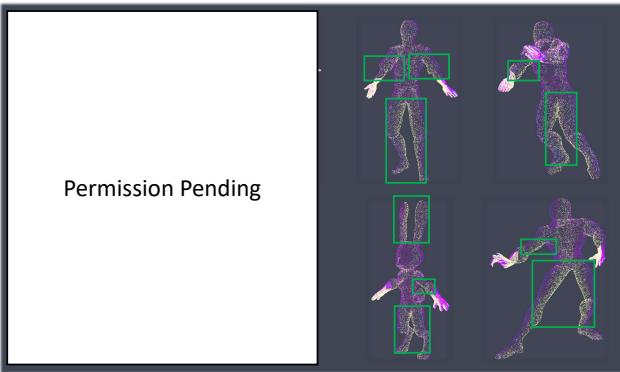


Figure S11. **Visualization of the attention score of our geometry-aware attention module.** These per-sampled-point values are extracted from the first attention head (out of 8 heads in total). The brighter color (yellower) indicates more attention to normals rather than coordinates. We also use green bounding boxes to label some clusters where high-attention-score points are densely distributed. It can be observed that the module adaptively learns to rely more on normals in regions like the inner thigh since coordinates become less discriminative there.

(b) Capacity of abnormal shapes. For those characters with an exaggerated body ratio (*e.g.*, extremely large head, short limbs, etc), our method can adaptively change the bone length to fit the shape. This is intractable for template-based human-generating works.

(c) Complex input poses. Poses far from the T/A-pose are fully supported. Our model can not only produce the well-fitted posed skeleton, but can also transform it into the T-pose for further animating applications.

(d) Efficiency for high polygon models. Benefiting from the particle-based shape autoencoder (Sec. 3.2), our framework is efficient and robust for different input resolutions, ranging from low-poly meshes like (c) to practical game-

level 3D models like (d). The Wukong model in (d) has over 1 million triangular faces, but our method can still make it animatable within 3 seconds.

(e) Support for asymmetric inputs. While many rigging methods assume symmetric inputs, our method can effectively deal with asymmetric ones. For example, the big gun of this cyborg in (e) is bound well to his arm bones.

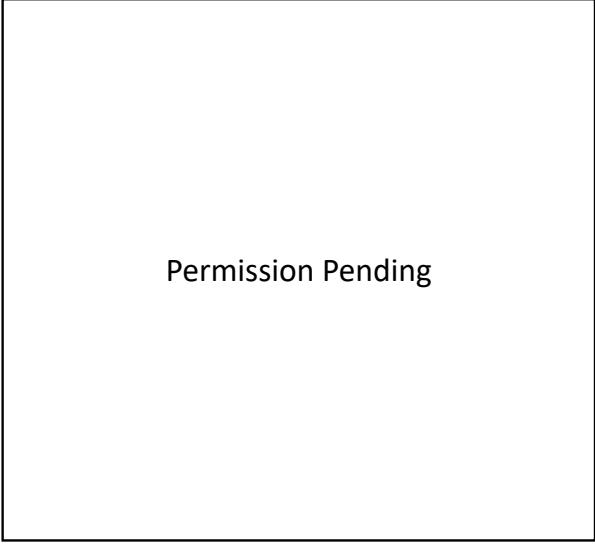
(f) Adaptation to non-existing bones. This armless statue is predicted to have extracorporeal arm bones. However, due to the spatial query mechanism of our framework, those dummy bones have no weights toward the actual vertices and can be simply removed without any side effects.

(g) & (h) Extension to extra bones. Once fine-tuned with some extra data (Sec. 3.5), our model is capable of producing positions and weights of non-standard bones. Here we show the long tail of a tiger in (g) and the long ears of a bunny in (h), which are all fully animatable.

Geometry-awareness: attention vs. add. Intuitively, simply concatenating or adding the normals to the coordinates when feeding points into the shape encoder can achieve the same goal as our geometry-aware attention. However, we found in practice that such a vanilla injecting strategy often leads to overfitting on the high-quality training mesh normals. The weight prediction depends so much on normals that some regions are influenced by far-away bones just because they have plausible normals, especially when faced with lower-quality meshes (*e.g.*, produced by generative models). We present a typical case in Fig. S10. While injecting normal information via adding resolves the weight corruption problem in the armpit region, it introduces a new issue of incorrect weights on hands. In contrast, the proposed attention mechanism benefits from normal information without any side effects.

Geometry-aware attention score. In addition to the two exemplar cases in Fig. 8, here we visualize the geometry-aware attention score of more characters in Fig. S11. The distribution of high-attention-score points shows patterns with statistical significance. Specifically, regions with possible spatial ambiguity, *e.g.*, inner thigh, armpit, between ears, *etc.*, are affected more by the normal values. This verifies the effectiveness of our geometry-aware attention, as it works exactly in the desired way to adaptively exploit the normal information.

Limitations of SMPL-based rigging. Template-based avatar generation methods [32, 35] benefit a lot from the well-defined SMPL mesh [37, 43], enabling the generation of animatable characters with no additional rigging cost. However, as discussed in Secs. 1 and 4.2, one of the unneglectable limitations of these methods is their inability to depart from realistic human shapes. Although



SMPL-X Fitting

Our Rigging and Skinning Results



Template-based Text-to-3D Generation Results

Figure S12. **Illustration of the limitations of SMPL-based rigging.** While SMPL provides a good template for skeletons and weights, it lacks the flexibility to handle exaggerated body shapes.

TADA [32] attempts to address this issue by predicting vertex deformations of the template mesh, it remains constrained by the preset body ratio of SMPL. Fig. S12 demonstrates some practical examples of cartoon characters with exaggerated body shapes that the SMPL model can hardly accommodate. As illustrated in the left part of Fig. S12, the large heads of these characters cannot be fitted by the template meshes, and the interpolated blend weights will definitely be incorrect in the head regions. In the right part of Fig. S12, we exhibit the results generated by two template-based methods using the same text prompt “cartoon brown bear with extremely large head and short legs”. Despite the specific prompting, both methods still produce body ratios resembling those of realistic humans, which limits their applicability in practical scenarios. In contrast, our method offers a promising solution by making bipedal characters of any shape ready for animation.