
AniDynRecon: Animatable 4D Dynamics Reconstruction from Sparse Point Clouds

Anonymous Author(s)

Affiliation

Address

email

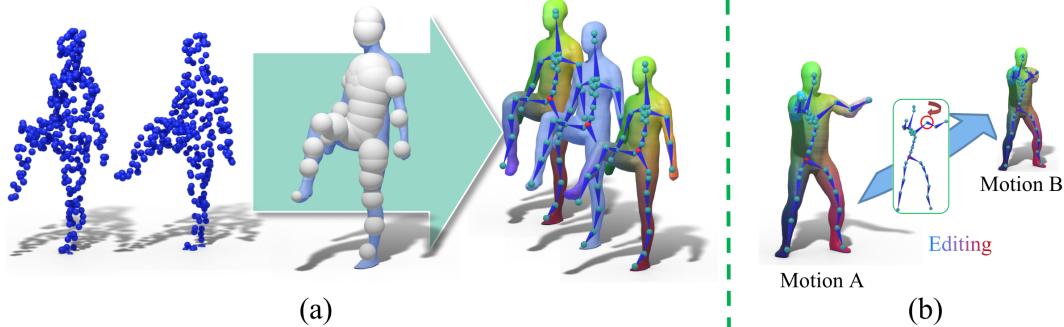


Figure 1: AniDynRecon holistically reconstructs fully animatable skeletal structures—including hierarchical topology, skinning weights, and surface geometry—directly from sparse and dynamic point clouds. (a) Reconstructed skeletal and geometric representations. (b) Motion editing results based on reconstructed skeleton structure.

Abstract

1 Reconstructing dynamic objects from sparse point cloud sequences without corre-
2 spondence is a challenging task. Existing methods usually rely on object templates
3 as strong geometric priors to ensure the accuracy and temporal consistency of
4 geometric reconstruction. In this paper, we explore this highly ill-posed problem
5 of simultaneously reconstructing object skeletons, skinning weights, motion fields,
6 and surfaces from sparse point cloud sequences in an unsupervised manner. Specif-
7 ically, we propose a *Hierarchical Skeletal Modeling (HSM)* module, which adopts
8 a compact, interpretable skeletal abstraction inspired by the Medial Axis Trans-
9 form. The skeleton representation facilitates jointly learn plausible joints, skinning
10 weights in an unsupervised manner for capturing dynamics while reconstructing
11 shapes. In addition, we introduce a *Combinatorial Motion Sampling (CMS)* strat-
12 egy to adaptively generate new motion sequence for data augmentation utilizing
13 the learned hierarchical skeletons. This scheme enhances the motion diversity of
14 objects and improves the generalization of the model, especially for unseen out-
15 of-domain scenarios. Extensive experiments on multiple benchmarks demonstrate
16 that AniDynRecon outperforms existing methods in terms of accuracy, robustness,
17 and structural interpretability, while enabling flexible downstream applications
18 such as motion synthesis and editing. Our project page is <https://anidynrecon-anonymous.github.io/>.

20 **1 Introduction**

21 Dynamic object reconstruction is a fascinating task in computer vision and graphics, given its
22 widespread applications across animation, AR, VR, robotics, and more. However, accurately capturing
23 and modeling the dynamic objects from sparse point cloud data remains a significant challenge.
24 First, it demands simultaneous modeling geometry and motion with strong temporal consistency,
25 which is a challenging ill-posed problem. Despite significant efforts for shape modeling Loper et al.
26 (2023); Li et al. (2017); Osman et al. (2020); Romero et al. (2022); Zuffi et al. (2017), these methods
27 typically rely on pre-defined templates or strong assumptions about the object’s structure. Neural
28 implicit functions Palafox et al. (2021); Ma et al. (2020); Saito et al. (2021) further enhanced the
29 power of their expression. However, reconstruction based on template categories significantly limits
30 their generalizability and adaptability. Additionally, previous approaches Yao et al. (2025b) are
31 frequently suffering from issues such as inaccurate joint localization, suboptimal skinning weight
32 locality, and reduced robustness in the presence of noisy or incomplete data; See Sec. 4.1. Second,
33 while recent advancements in deep learning have shown encouraging results Lei and Daniilidis
34 (2022); Niemeyer et al. (2019); Cao et al. (2024); Tang et al. (2021b), these data-driven approaches
35 are inherently limited in their generalization capability by the scope of the training data. In particular,
36 the scarcity of sparse point cloud data further exacerbates the lack of generalization and robustness in
37 such methods; See Sec. 5.3.

38 In this paper, we present AniDynRecon, a novel unsupervised end-to-end framework that reconstructs
39 animatable skeletal structures—capturing skeletal hierarchy, skinning, and geometry—directly from
40 dynamic point clouds. We introduce *Hierarchical Skeletal Modeling (HSM)* for dynamics recon-
41 struction from the sparse point cloud inputs. Unlike prior work Yao et al. (2025b) that overlooks
42 skeletal representations, we adopt a template-free skeletal abstraction inspired by the Medial Axis
43 Transform Blum (1967) as a compact and expressive low-dimensional prior. This design offers two
44 main advantages: (1) it effectively captures geometric structures of 3D shapes for shape recovery, as
45 also evidenced by Lin et al. (2021); Dou et al. (2022); Tang et al. (2019) for static shape modeling,
46 which helps obtain plausible joints without template, and (2) it can be naturally integrated with our
47 unsupervised skinning and geometry learning modules. It enables animatable reconstructions for
48 flexible downstream applications such as motion editing, as shown in Fig. 1. This helps generate
49 ample motion data for training a robust model. While recent work Yao et al. (2025b) models 4D
50 dynamics using skeletons but lacks correct joint positions for editing. In contrast, our method follows
51 a standard animation pipeline, learning controllable joints to support motion sampling.

52 To overcome the challenge caused by limited 4D dynamic data, we propose a *Combinatorial Motion-*
53 *Based Sampling (CMS)* strategy that augments temporal dynamics through intra-dataset motion
54 recombination. CMS synthesizes diverse training samples by modeling the distribution of local
55 motion segments and recombining them across different temporal windows, leveraging inherent
56 motion variations. We show the effectiveness of CMS in significantly improving the robustness and
57 accuracy of the learned model, especially for unseen scenarios.

58 We conduct extensive experiments to demonstrate that AniDynRecon outperforms state-of-the-
59 art (SOTA) methods in accuracy, robustness, and interpretability. In summary, our contributions are
60 as follows:

- 61 • We propose AniDynRecon, an end-to-end framework that holistically reconstructs fully
62 animatable skeletal structures—including hierarchy, skinning, and geometry—directly from
63 dynamic point clouds.
- 64 • We introduce *Hierarchical Skeletal Modeling (HSM)*, a compact and expressive representa-
65 tion inspired by the Medial Axis Transform, which unsupervisedly captures tree-structured
66 joint hierarchies and local transformations.
- 67 • We develop a *Combinatorial Motion-Based Sampling (CMS)* strategy that leverages intra-
68 dataset temporal variations to augment training data, thereby improving robustness under
69 sparse or limited supervision.
- 70 • We demonstrate that AniDynRecon achieves the SOTA performance on multiple dynamic
71 point cloud benchmarks in terms of accuracy, robustness, and structural interpretability.

72 **2 Related Work**

73 **3D Representation** 3D representations can be broadly categorized into explicit and implicit types.
74 Explicit representations directly encode geometry, such as meshes Groueix et al. (2018); Liao et al.
75 (2018); Pan et al. (2019); Tang et al. (2019, 2021a), point clouds Achlioptas et al. (2018); Fan
76 et al. (2017); Zhao et al. (2021), octrees Riegler et al. (2017); Häne et al. (2017); Tatarchenko et al.
77 (2017); Wang et al. (2018) and voxels Choy et al. (2016); Gadelha et al. (2017); Han et al. (2017);
78 Riegler et al. (2017); Stutz and Geiger (2018). Parametric models, such as SMPL Loper et al. (2023),
79 STAR Osman et al. (2020), FLAME Li et al. (2017), MANO Romero et al. (2022) and SMAL Zuffi
80 et al. (2017), have been widely used to represent specific shape categories (e.g., human bodies, faces,
81 hands and animals) through explicit geometric representations. These models leverage fixed templates
82 and parameters to efficiently adjust pose and shape. Recently, implicit representations model 3D
83 shapes using continuous functions, such as in Signed Distance Field (SDF) Chabra et al. (2020);
84 Chen et al. (2023); Chen and Zhang (2019); Park et al. (2019) or occupancy networks Chibane
85 et al. (2020); Mescheder et al. (2019). Implicit methods are more efficient and robust in handling
86 complex shapes, as they provide smooth, continuous surfaces and can adapt to varying levels of detail.
87 These representations are also memory-efficient and excel in tasks like shape interpolation and 3D
88 reconstruction, particularly with noisy or sparse data. Recent advancements, such as DeepSDF Park
89 et al. (2019) and GridPull Chen et al. (2023), demonstrate the advantages of neural implicit models
90 in recovering high-quality 3D shapes from incomplete data, outperforming explicit methods in terms
91 of flexibility and robustness.

92 **Dynamic Reconstruction** Recent advancements have extended traditional 3D representations into
93 the fourth dimension (4D), enabling more effective modeling of dynamic object behaviors over
94 time Bozic et al. (2021); Fan et al. (2017); Jiang et al. (2021); Lei and Daniilidis (2022); Niemeyer
95 et al. (2019); Palafox et al. (2021); Tang et al. (2021b, 2022); Zou et al. (2023). Prominent works
96 include OFLOW Niemeyer et al. (2019), LPDC Tang et al. (2021b) and CaDeX capture both
97 shape and deformation of high quality. Recently, Motion2Vecsets Cao et al. (2024) employs a
98 set of latent codes to represent the deformation of points. However, it is highly dependent on the
99 correspondence between points across frames. Drawing inspiration from DynoSurf Yao et al. (2025b),
100 which represents the deformation field using control points, our method refines the distribution of
101 their control points and highly expands the generalization ability.

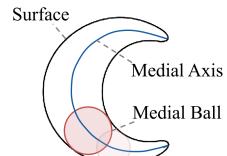
102 **Automatic Skinning and Rigging** Linear blend skinning (LBS) Le and Deng (2012) has been exten-
103 sively employed for real-time animation of articulated characters. However, the manual adjustment
104 of joints, skeleton structures, and blending weights by artists proves to be a time-consuming en-
105 deavor Bang and Lee (2018). In this regard, traditional method *Pinocchio* Baran and Popović (2007)
106 introduced a template-based approach that enables automatic fitting of a template skeleton to a target
107 mesh, thereby generating an animation-ready rig. Other geometric-based skinning methods exploit
108 the theory of heat diffusion Wareham and Lasenby (2008), bounded biharmonic energy Jacobson et al.
109 (2011), physics-inspired approaches Kavan and Sorkine (2012) and geodesic voxel binding Dionne
110 and de Lasa (2013). Data-driven methods can capture information implied in ground-truth rigs edited
111 by experienced animators, such as RigNet Xu et al. (2020). However, it suffers from poor robustness,
112 which is not end-to-end trainable. HumanRig Chu et al. (2024) and Make-it-Animatable Guo et al.
113 (2024) rely on template skeletons, which also lack the ability to rig unseen data categories.

114 All these approaches are unsuitable for the task of 4D reconstruction, which relies solely on point
115 cloud, and is challenging to be trained for strong generalization. Uniquely, when coupled with the
116 theory of MAT we are able to train a robust model for reconstructing animation-ready meshes from
117 unseen point cloud sequences, generating diverse training samples.

118 **3 Preliminaries**

119 **3.1 Medial Axis Transform**

120 Medial axis transform (MAT) is a well-known method representing a closed,
121 oriented, and bounded two-manifold surface S as media-axis \mathcal{M} and its radius
122 function \mathcal{R} Blum (1967). The media-axis \mathcal{M} is defined as the set of points in
123 the interior with more than one closest point on the boundary surface, where
124 the closest distance is defined as the corresponding radii \mathcal{R} . Each medial axis



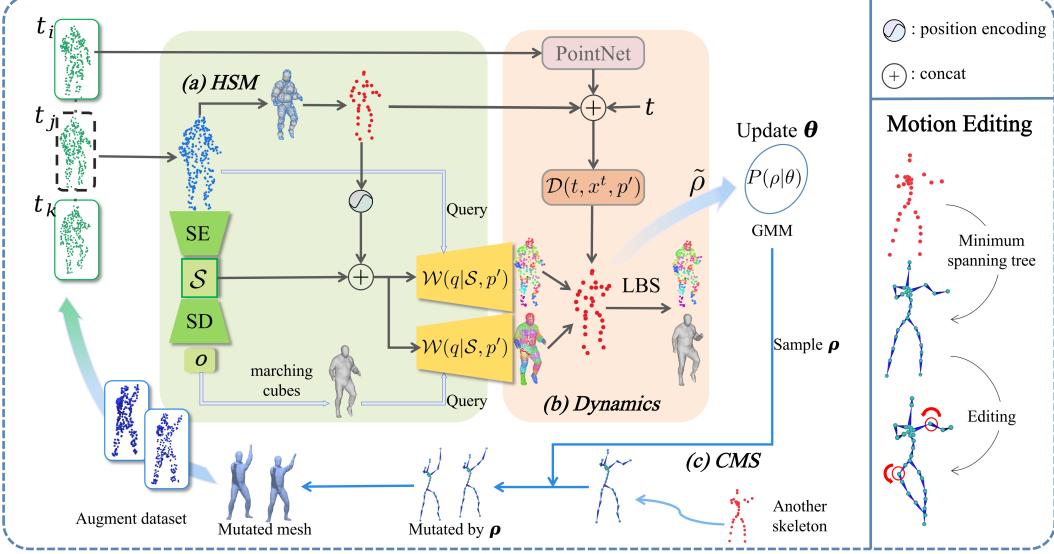


Figure 2: Overview of our method AniDynRecon. **(a) Hierarchical Skeletal Modeling (HSM)** (Sec. 4.1). Given a point cloud sequence, we initially select one frame to serve as the rest pose and encode it as shape tokens \mathcal{S} . Subsequently, we learn skeletal points along with their medial balls in an unsupervised manner. The occupancy field O of the rest pose is decoded from the shape tokens, and the skinning weight field can be decoded by combining the shape token with the control points. The final mesh of the rest pose is extracted from the occupancy field using marching cubes. **(b) Dynamics.** (Sec. 4.2). With skinning weights and the deformed skeletal points output from MLP \mathcal{D} , the mesh is deformed via LBS. The deformation results is trained to match the point cloud sequence. **(c) Combinatorial Motion-Based Sampling (CMS)** (Sec. 4.3). With the animatable reconstruction result, we model the distribution of joint rotation $\tilde{\rho}$ by fitting an GMM, then we sample new ρ s to obtain new motion sequences to drive the dynamics for building an augment dataset.

125 point, combined with its radius, defines a medial ball. Inspired by this technique, we use the center of
 126 medial balls as control points because they have a good quality of centrality and evenly distribution
 127 across mesh.

128 3.2 Linear Blend Skinning

129 Linear Blend Skinning (LBS) Le and Deng (2012) is a widely used technique used in computer
 130 graphics and animation to simulate the deformation of a 3D mesh model, particularly when the model
 131 is subjected to skeletal animation. Given skinning weight w_{ij} , the j -th joint's initial position p'_j of
 132 rest pose, delta coordinate Δp_j^t and rotation matrix R_j^t of joint j at frame t , the deformed position of
 133 point i can be computed by LBS: $x_i^t = \sum_{j=1}^K w_{ij}(R_j^t d'_{ij} + p'_j + \Delta p_j^t)$, where p'_j means position of
 134 initial joints of rest pose and $d'_{ij} = x'_i - p'_j$ is the offset from joint to surface point.

135 4 Method

136 Our pipeline consists of three parts: Sec. 4.1 introduces unsupervised joint extraction via hierarchical
 137 structural modeling; Sec. 4.2 tackles joint optimization of shape, joints, skinning, and motion; Sec. 4.3
 138 presents a synthetic-data-driven augmentation strategy to improve generalization. The overall pipeline
 139 is shown in Fig. 2.

140 4.1 Hierarchical Skeletal Modeling

141 **Initial state of joints.** We adopt the skeletal point prediction framework from Point2Skeleton Lin
 142 et al. (2021), which establishes skeletal points as local centroids of surface point clusters. To ensure
 143 high-quality joint extraction, we initialize the training process using a Chamfer Distance loss between
 144 predicted joints and the input point cloud in rest pose. This geometric constraint effectively distributes
 145 control points across different anatomical regions during the initial training phase.

146 Following joint position stabilization, we progressively introduce radius parameters r based on the
 147 concept of medial axis transform. Each sphere, centered at a joint position, is optimized to maximize
 148 contact with the rest pose surface and the corresponding radius is optimized as large as possible
 149 through the skeletal sphere loss $\mathcal{L}_{\text{centrality}}$ in Point2Skeleton Lin et al. (2021) (see the results of joints
 150 and their radii in Fig. 3):

$$\mathcal{L}_{\text{centrality}} = \left[\sum_{x \in \{x_i\}} \left(\min_{p' \in \{p'_j\}} \|p' - x\|_2 - r(p_x^{\min}) + \sum_p (\min_x \|p - x\|_2 - r(p')) \right) \right] - \sum_{p' \in \{p'_j\}} r(p') \quad (1)$$

151 where p_x^{\min} means the closest joint to the input point x . Notably, joint positions are computed as convex combinations of
 152 surface points, with the corresponding combination weights (hereafter
 153 termed joint weights) being normalized through softmax activation:

$$\tilde{w}_{ij} = \frac{\exp(\hat{w}_{ij})}{\sum_i \exp(\hat{w}_{ij})} \quad (2)$$

155 where \tilde{w}_{ij} means joint weight between the j -th bone and the i -th
 156 vertex and \hat{w}_{ij} is the output of joint weight MLP: $\hat{w}_{ij} = \text{MLP}(f)$
 157 (f is an embedding of rest pose encoded by PointNet++ Qi et al.
 158 (2017b)). Same way as Point2Skeleton Lin et al. (2021), we use
 159 joint weights to reconstruct the position of the joints:

$$p'_j = \sum_{i=1}^N \tilde{w}_{ij} x'_i \quad (3)$$

160 which naturally gives the joint weights locality. (Joints tend to move towards surface points with
 161 higher weights). At the same time, we calculate the chamfer distance between the rest pose joints
 162 p' and surface points x' to ensure that the joints are evenly dispersed, which aids in the subsequent
 163 LBS. More specifically, dispersed joints prevent input point cloud regions from being covered by
 164 multiple joints with overlapping skinning weights. Another advantage of this approach is that our
 165 joints possess neutrality thanks to medial axis transform, meaning all joints are guaranteed to be
 166 distributed within the inscribed sphere of the point cloud.

167 Then in order to satisfy the definition of standard LBS, the joint weights should be normalized to
 168 temporary skinning weights $w'_{ij} = \tilde{w}_{ij} / \sum_j \tilde{w}_{ij}$ for further training a skinning weight field.

169 4.2 Jointly Train Shape, Skinning Weights and Motion

170 **Occupancy field.** We first extract a shape token $\mathcal{S} = \text{SE}(x') = \text{CrossAttn}(PE(x'), PE(x^d))$ (x^d is
 171 subsampled point cloud from x' using furthest point sampling and PE is the positional embedding)
 172 from the point cloud of the rest pose, which can subsequently be decoded into occupancy values
 173 $O = \text{SD}(q|\mathcal{S}) = \text{CrossAttn}(\text{SelfAttn}(\mathcal{S}), PE(q))$ at locations of query points q using the same
 174 shape encoder and decoder as those employed in Motion2Vecsets Cao et al. (2024). The occupancy
 175 values O , which is necessary for further mesh reconstruction using marching cubes Lorensen and
 176 Cline (1998), is supervised by ground truth values \hat{O} using a BCE loss $\mathcal{L}_{\text{shape}} = \text{BCE}(O, \hat{O})$. Thanks
 177 to the occupancy field, we can use our joints as query points for shape decoder to check if they are
 178 inside the mesh volume:

$$\mathcal{L}_{\text{joints_inside}} = \sum_{j=1}^K \text{BCE}(\text{SD}(p'_j|\mathcal{S}), 1) \quad (4)$$

179 However, the supervision of $\mathcal{L}_{\text{joints_inside}}$ have a bad influence on the training process of shape encoder
 180 and decoder because it uses shape token \mathcal{S} as condition. To avoid this, We freeze the network
 181 parameters and gradients of shape encoder and decoder before the calculation of $\mathcal{L}_{\text{joints_inside}}$. This
 182 strategy ensures that the reconstruction of rest pose is correctly supervised and the joint positions are
 183 constrained at the same time.

184 **Skinning weight field.** To model the skinning weight field, we designed a simple skinning weight
 185 decoder, which consists of a cross-attention block and a fully connected (FC) layer. Our skinning

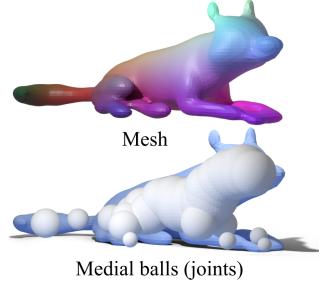


Figure 3: The medial balls approximately align with the mesh volume.

186 weight decoder takes surface points as query and joints along with shape token as key and value,
 187 outputting skinning weights w_{ij} at the query position q :

$$\mathcal{W}(q|\mathcal{S}, p') = FC(\text{CrossAttn}(PE(q), \text{Concat}(\mathcal{S}, PE(p')))) \quad (5)$$

188 When the query q is one point of input point cloud $q = x'_i$, the corresponding skinning weights will
 189 be $w_i = \{w_{i0}, w_{i1}, \dots, w_{iK}\}$. With these skinning weights, we can supervise deformation directly
 190 on input point cloud using Eq. (7). Besides, We train our skinning weight decoder \mathcal{W} based on
 191 the supervision $\|w_{ij} - w'_{ij}\|$ of joint weights in the early several epochs for faster convergence.
 192 Once the skinning weight field was trained and mesh of rest pose containing vertices \bar{V} and triangle
 193 faces \bar{F} was reconstructed from the occupancy field, we can use the vertices \bar{V} to query skinning
 194 weights \bar{W} on mesh surface for deformation. Given the reconstructed mesh $\{\bar{V}, \bar{F}\}$, skinning weights
 195 \bar{W} correspond to its vertices and joints p at any time, it is easy to get the deformed mesh $\{\hat{V}, \hat{F}\}$
 196 whose vertices are driven by LBS $\hat{V} = \text{LBS}(\hat{V}, \bar{W}, p)$. Additionally, we define a locality loss to
 197 further concentrate the skinning weights around their respective joints, clearing the tiny weights
 198 corresponding to points far from the joints. The locality loss is a simple *ReLU* function:

$$\mathcal{L}_{\text{locality}} = \sum_{i=1}^N \sum_{j=1}^K \text{ReLU}((d'_{ij} - D)w_{ij}) \quad (6)$$

199 where D is a hyper-parameter denoting the maximum distance beyond which skinning weights are
 200 unconstrained. The skinning weight field is visualized in Fig. 4, 7, and an example of pose editing
 201 using our skinning weights is shown in Fig. 7.

202 **Deformation of joints.** As the deformation field of sequences is represented as deformation of
 203 joints through LBS in our framework, we use an MLP module \mathcal{D} to capture dynamics contained
 204 in input point cloud sequence. Initially, we extract input point cloud features using PointNet Qi
 205 et al. (2017a) (we opt against PointNet++ Qi et al. (2017b) due to its excessive GPU memory
 206 consumption when training on point cloud sequences). Simultaneously, we apply positional em-
 207 bedding to encode time t and initial joint positions p' . These latent codes are then fed into a
 208 deformation MLP, yielding the joints' delta rotation and delta position $\{\Delta p^t, r^t\} = \mathcal{D}(t, x^t, p') =$
 209 $\text{MLP}(\text{Concat}(PE(t), PE(p'), \text{PointNet}(x^t)))$ at time t . Here the rotation matrices of joints are
 210 represented as lie algebra. Finally, we utilize the chamfer distance between the point collections
 211 $\tilde{x}^t \in \{\tilde{x}_i^t\}$ of the input point cloud and output collection $x^t \in \{x_i^t\}$ of LBS in each frame as a loss
 212 function $\mathcal{L}_{\text{dynamic}}$:

$$\mathcal{L}_{\text{dynamic}} = \sum_{t=1}^T \left(\sum_{x^t \in \{x_i^t\}} \min_{\tilde{x}^t \in \{\tilde{x}_i^t\}} \|x^t - \tilde{x}^t\|_2 + \sum_{\tilde{x}^t \in \{\tilde{x}_i^t\}} \min_{x^t \in \{x_i^t\}} \|\tilde{x}^t - x^t\|_2 \right) \quad (7)$$

213 Thanks to LBS and supervision of chamfer distance, our pipeline reconstructs dynamic mesh sequence
 214 from noisy point cloud sequence with no correspondence.

215 4.3 Combinatorial Motion-Based Sampling

216 To obtain more diverse motions for training a robust model, we sample new motions based on the
 217 original training dataset. Given joints in a novel pose, the mesh extracted from the occupancy field
 218 using marching cubes can be deformed to the new pose through LBS. AniDynRecon allows us to
 219 manually craft novel motion sequences that conform to the distribution of the training data.

220 However, manual editing of unstructured control points presents significant challenges. For instance,
 221 animating a simple arm movement requires individually manipulating all control points within the
 222 arm region. To address this limitation, we construct a skeleton tree using Prim's algorithm based on
 223 our joint structure, following the approach established in RigGS Yao et al. (2025a). This enables a
 224 more intuitive animation workflow: by randomly selecting a joint in the skeleton tree and applying
 225 a random rotation, child nodes are automatically transformed to new positions through forward
 226 kinematics. Instead of disturbing the synthetic data to a random distribution, we fit a Gaussian
 227 Mixture Model (GMM) for rotation euler angles ρ of joints using the original training data:

$$P(\rho|\theta) = \sum_{m=1}^M \pi_m \mathcal{N}(\rho|\mu_m, \Sigma_m) \quad (8)$$

Input	Method	Seen Individual				Unseen Individual		
		F-1 ↑	IoU ↑	CD ↓	Corr ↓	IoU ↑	CD ↓	Corr ↓
DT4D-A Li et al. (2021)	OFlow Niemeyer et al. (2019)	-	70.6%	0.104	0.204	57.3%	0.175	0.285
	CaDex Lei and Daniilidis (2022)	-	80.3%	0.061	0.133	64.7%	0.127	0.239
	DynoSurf Yao et al. (2025b)	0.401	-	0.174	0.463	-	-	-
	Ours	0.858	81.7%	0.055	0.166	74.8%	0.073	0.216
D-FAUST Bogo et al. (2017)	OFlow Niemeyer et al. (2019)	-	79.9%	0.073	0.122	69.6%	0.095	0.149
	CaDex Lei and Daniilidis (2022)	-	85.5%	0.056	0.100	75.4%	0.074	0.126
	DynoSurf Yao et al. (2025b)	0.402	-	0.118	0.175	-	-	-
	Ours	0.854	83.0%	0.060	0.173	76.5%	0.070	0.197

Table 1: Performance comparison on DT4D-A Li et al. (2021) and D-FAUST Bogo et al. (2017) datasets. IoU (%) indicates Intersection over Union (higher is better), CD indicates l_1 -Chamfer Distance (lower is better), and Corr indicates correspondence error (lower is better), with cell colors to indicate the **best** and **second best**.

where M is the number of components in the GMM, $\mathcal{N}(\rho|\mu_m, \Sigma_m)$ is the probability density function of the m -th Gaussian distribution with mean μ_m and covariance matrix Σ_m , $\theta = \{\pi, \mu, \Sigma\}$ represents the parameters for GMM. The training objective of GMM is to maximize the log-likelihood function of the observed data $\ln P(\tilde{\rho}|\theta) = \sum \ln P(\rho|\theta)$ where $\tilde{\rho} = \{\rho\}$ is the observed dataset. Once P is fitted, we sample K euler angles from GMM as the rotation of joints for synthetic pose.

5 Experiments

Datasets: Our experiments encompass two distinct categories of dynamic objects: human bodies (Sec. 5.1) and animals (Sec. 5.2). To ensure a fair comparison, we adopt the same experimental setup, datasets, and data split as CaDeX Lei and Daniilidis (2022). Firstly, human body dataset D-FAUST Bogo et al. (2017), containing 10 subjects and 109 sequences, is split into training (70%), validation (10%), and test (20%) subsets. Secondly, DeformingThings4D-Animals Li et al. (2021), including 38 indentities with a total of 1227 animations, is divided into training (75%), validation (7.5%), and test (17.5%) subsets.

Baselines: We compare our method with several state-of-the-art approaches, including OFlow Niemeyer et al. (2019), CaDeX Lei and Daniilidis (2022), and DynoSurf Yao et al. (2025b).

Evaluation metrics: The Intersection over Union (IoU) metric assesses the overlap between the predicted and ground truth meshes. F-score with thresholds of 1% (F-1%). The Chamfer distance computes the average nearest-neighbor distance between two point sets. The l_2 -distance error quantifies the Euclidean distance between corresponding points on the predicted and ground truth meshes.

Runtime: The enhanced capabilities of our pipeline result in a more challenging training process compared with those aforementioned methods. To ensure the reliability and consistency of our experimental results, all the experiments are conducted on a single 4090 GPU for 10 days. We present more experiment details in the supplementary material.

5.1 Human Bodies

The input sequence of human bodies comprises 17 frames, with each frame of the input point cloud containing 300 points, evenly distributed over time to supervise deformation. The testing set encompasses two difficulty levels: unseen motion and unseen individuals Niemeyer et al. (2019). We show our results in Fig. 5. As demonstrated in Tab. 1, our method achieves advanced performance in testing on unseen individual datasets, which is more challenging than unseen motion sets. Subsequently, to

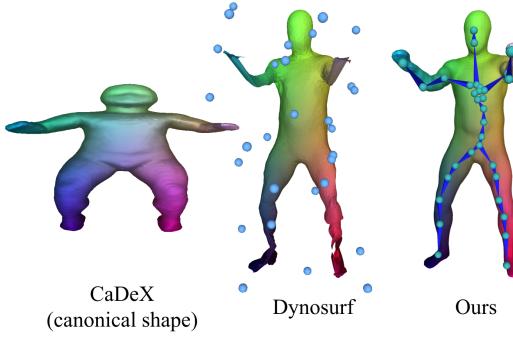


Figure 4: Visual comparison of our joints distribution with CaDeX Lei and Daniilidis (2022) and Dynosurf Yao et al. (2025b).

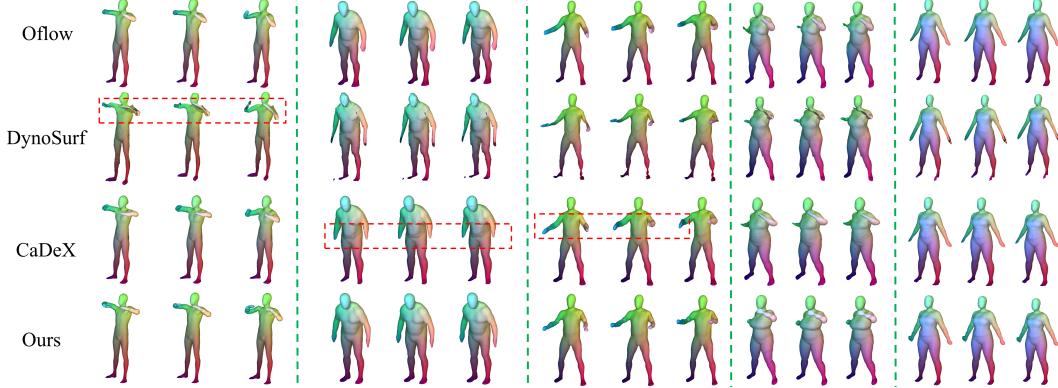


Figure 5: Comparisons on the D-FAUST Bogo et al. (2017) dataset.

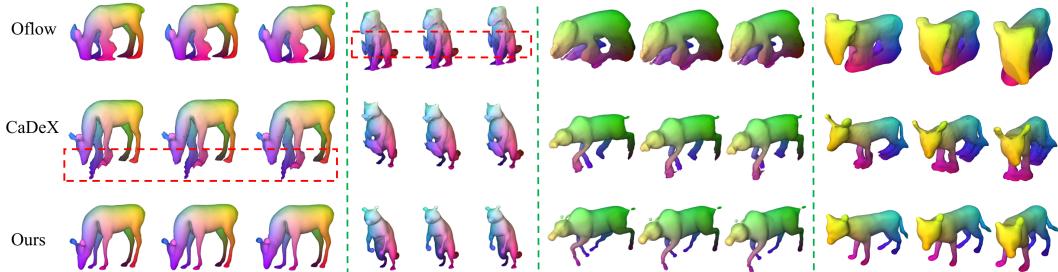


Figure 6: Comparisons on the DT4D Li et al. (2021) dataset.

267 make a fair comparison with DynoSurf Yao et al. (2025b) which is only trained and visualized on
 268 unseen motion datasets, we adhered to the same comparison procedure as in DynoSurf Yao et al.
 269 (2025b). The results presented in Tab. 1 highlight the advanced performance of our work against
 270 theirs. In addition to the occupancy field of shape and deformation field reconstructed by methods
 271 such as CaDeX Lei and Daniilidis (2022), our pipeline stands out by extracting control points and a
 272 skinning weight field of significantly higher quality, as shown in Fig. 4. Unlike messy canonical shapes
 273 in CaDeX Lei and Daniilidis (2022) and messy control points in DynoSurf Yao et al. (2025b), our
 274 high-quality elements can be effectively utilized in the processes of generation and motion editing.
 275 As for the training strategy, we first train our joints with a regularization of chamfer distance loss
 276 between input point cloud and control points in the early 48 epochs. After that, the positions of joints
 277 will be optimized using the centrality loss. We refer the reader to the supplementary material for
 278 more details. It is obviously that in Fig. 4, our joints distribute evenly inside mesh but the result of
 279 DynoSurf Yao et al. (2025b) show that most of joints scatter outside the mesh irregularly.

280 5.2 Animals

281 Similar to human bodies, we utilize
 282 a 17-frame point cloud sequence as
 283 input. However, due to the greater
 284 morphological diversity across animal
 285 categories and the increased difficulty
 286 in capturing deformations compared
 287 to humans, which have a more uni-
 288 form skeletal structure, we employ
 289 512 points for deformation supervision,
 290 consistent with CaDeX Lei and Daniilidis (2022). The
 291 results is shown in Fig. 6. We present more details in the supplementary material.

Input	Method	Unseen Individual		
		IoU ↑	CD ↓	Corr ↓
D-FAUST Bogo et al. (2017)	Ours(w/o joints_inside)	76.1%	0.072	0.162
	Ours(w/o boosting)	76.5%	0.071	0.209
	Full	76.5%	0.070	0.197

Table 2: Quantitative ablation studies on the D-FAUST Bogo et al. (2017) dataset.

291 5.3 Ablation Study

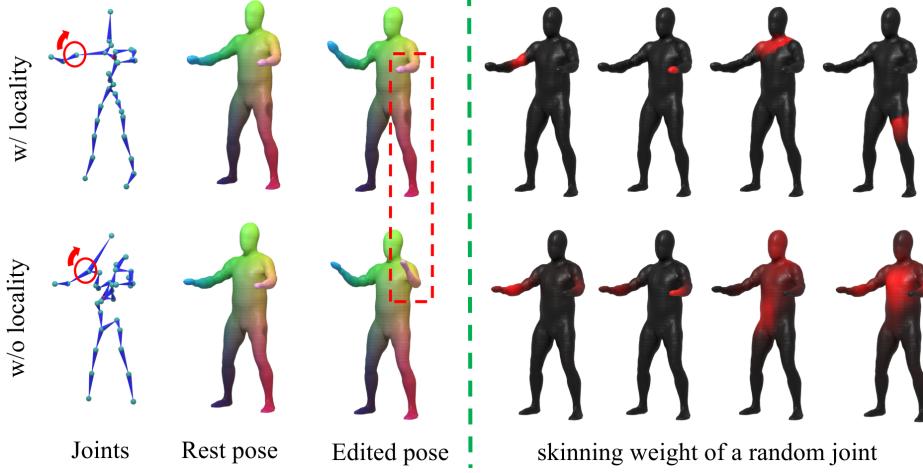


Figure 7: The locality loss $\mathcal{L}_{\text{locality}}$ ensures a good distribution of joints and skinning weights, consequently leading to correct motion editing results.

292 To investigate the impact of different components in our
 293 framework, we conduct ablation studies focusing on three
 294 key aspects: skinning weight locality regularization, joint
 295 position inside constraints, and data augmentation strategy.

296 For skinning weight locality regularization, in Fig. 7 we
 297 demonstrate that without locality loss, the distribution of
 298 joints will be messy and editing control points of one
 299 human arm will cause another arm to move, highlighting
 300 the importance of locality regularization for achieving
 301 plausible deformations. Fig. 7 also illustrates the skinning
 302 weight distribution for a random joint without locality loss
 303 supervision. For joint position inside constraints, we set
 304 the weight of joint inside loss $\mathcal{L}_{\text{joints_inside}}$ to zero and find
 305 that some joints are generated outside the mesh, which is fatal
 306 for motion editing, as shown in Fig. 8. As indicated in Tab. 2,
 307 this negatively impacts deformation capture. Finally, we remove
 308 our data augmentation strategy and find that the correspondence
 metric drops, particularly for the task of
 unseen individuals, also shown in Tab. 2.

309 6 Conclusion

310 We present a framework named AniDynRecon for reconstructing fully animatable 4D dynamic
 311 objects directly from sparse and noisy point cloud observations. By introducing Hierarchical Skeletal
 312 Modeling (HSM), we provide an interpretable and compact representation that captures both geo-
 313 metric and structural information in an unsupervised manner. The integration of Combinatorial
 314 Motion-Based Sampling (CMS) further enhances the model’s ability to generalize by increasing the
 315 diversity of training dynamics, especially in scenarios with limited data. Together, these components
 316 enable high-fidelity reconstruction and realistic animation, supporting downstream applications such
 317 as motion editing and synthesis. Extensive experiments demonstrate that our method consistently out-
 318 performs existing approaches in accuracy, robustness, and structural coherence. However, our method
 319 still encounters limitations in cases where real world deformations are unsuitable to be represented
 320 using LBS, such as muscle-induced elastic dynamics. This issue becomes more pronounced in animal
 321 categories, which have more uncertain biologic structures compared to human bodies. Addressing
 322 such cases is beyond the scope of this paper. Future improvements could explore how to expand the
 323 expressive capabilities of LBS using a warpping function.

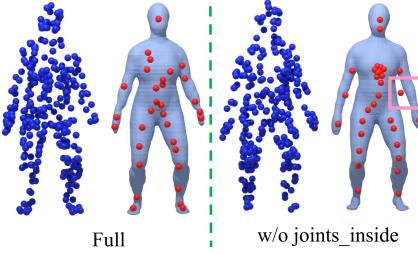


Figure 8: $\mathcal{L}_{\text{joints_inside}}$ avoids generating joints outside of the mesh.

324 **References**

- 325 Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and
326 generative models for 3d point clouds. In *International conference on machine learning*, pages 40–49. PMLR,
327 2018.
- 328 Seungbae Bang and Sung-Hee Lee. Spline interface for intuitive skinning weight editing. *ACM Transactions on*
329 *Graphics (TOG)*, 37(5):1–14, 2018.
- 330 Ilya Baran and Jovan Popović. Automatic rigging and animation of 3d characters. *ACM Transactions on graphics*
331 (*TOG*), 26(3):72–es, 2007.
- 332 Harry Blum. A transformation for extracting new descriptions of shape. *Models for the perception of speech*
333 *and visual form*, pages 362–380, 1967.
- 334 Federica Bogo, Javier Romero, Gerard Pons-Moll, and Michael J Black. Dynamic faust: Registering human
335 bodies in motion. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
336 6233–6242, 2017.
- 337 Aljaz Bozic, Pablo Palafox, Michael Zollhofer, Justus Thies, Angela Dai, and Matthias Nießner. Neural
338 deformation graphs for globally-consistent non-rigid reconstruction. In *Proceedings of the IEEE/CVF*
339 *Conference on Computer Vision and Pattern Recognition*, pages 1450–1459, 2021.
- 340 Wei Cao, Chang Luo, Biao Zhang, Matthias Nießner, and Jiapeng Tang. Motion2vecsets: 4d latent vector set
341 diffusion for non-rigid shape reconstruction and tracking. In *Proceedings of the IEEE/CVF Conference on*
342 *Computer Vision and Pattern Recognition*, pages 20496–20506, 2024.
- 343 Rohan Chabra, Jan E Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard
344 Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *Computer*
345 *Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX*
346 *16*, pages 608–625. Springer, 2020.
- 347 Chao Chen, Yu-Shen Liu, and Zhizhong Han. Gridpull: Towards scalability in learning implicit representations
348 from 3d point clouds. In *Proceedings of the ieee/cvf international conference on computer vision*, pages
349 18322–18334, 2023.
- 350 Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the*
351 *IEEE/CVF conference on computer vision and pattern recognition*, pages 5939–5948, 2019.
- 352 Julian Chibane, Thiemo Alldieck, and Gerard Pons-Moll. Implicit functions in feature space for 3d shape
353 reconstruction and completion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
354 *recognition*, pages 6970–6981, 2020.
- 355 Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified
356 approach for single and multi-view 3d object reconstruction. In *Computer Vision–ECCV 2016: 14th European*
357 *Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VIII 14*, pages 628–644.
358 Springer, 2016.
- 359 Zedong Chu, Feng Xiong, Meiduo Liu, Jinzhi Zhang, Mingqi Shao, Zhaoxu Sun, Di Wang, and Mu Xu.
360 Humanrig: Learning automatic rigging for humanoid character in a large scale dataset. *arXiv preprint*
361 *arXiv:2412.02317*, 2024.
- 362 Olivier Dionne and Martin de Lasa. Geodesic voxel binding for production character meshes. In *Proceedings of*
363 *the 12th ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 173–180, 2013.
- 364 Zhiyang Dou, Cheng Lin, Rui Xu, Lei Yang, Shiqing Xin, Taku Komura, and Wenping Wang. Coverage
365 axis: Inner point selection for 3d shape skeletonization. In <Top Cited Articles in CGF 2022-2023>;
366 *EUROGRAPHICS 2022; Computer Graphics Forum 2022*, pages 419–432, 2022.
- 367 Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction
368 from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
369 pages 605–613, 2017.
- 370 Matheus Gadelha, Subhransu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. In
371 *2017 international conference on 3d vision (3DV)*, pages 402–411. IEEE, 2017.
- 372 Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché
373 approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and*
374 *pattern recognition*, pages 216–224, 2018.

- 375 Zhiyang Guo, Jinxu Xiang, Kai Ma, Wengang Zhou, Houqiang Li, and Ran Zhang. Make-it-animatable: An
 376 efficient framework for authoring animation-ready 3d characters. *arXiv preprint arXiv:2411.18197*, 2024.
- 377 Xiaoguang Han, Zhen Li, Haibin Huang, Evangelos Kalogerakis, and Yizhou Yu. High-resolution shape
 378 completion using deep neural networks for global structure and local geometry inference. In *Proceedings of*
 379 *the IEEE international conference on computer vision*, pages 85–93, 2017.
- 380 Christian Häne, Shubham Tulsiani, and Jitendra Malik. Hierarchical surface prediction for 3d object reconstruc-
 381 tion. In *2017 International Conference on 3D Vision (3DV)*, pages 412–420. IEEE, 2017.
- 382 Alec Jacobson, Ilya Baran, Jovan Popovic, and Olga Sorkine. Bounded biharmonic weights for real-time
 383 deformation. *ACM Trans. Graph.*, 30(4):78, 2011.
- 384 Boyan Jiang, Yinda Zhang, Xingkui Wei, Xiangyang Xue, and Yanwei Fu. Learning compositional representation
 385 for 4d captures with neural ode. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
 386 *Recognition*, pages 5340–5350, 2021.
- 387 Ladislav Kavan and Olga Sorkine. Elasticity-inspired deformers for character articulation. *ACM Transactions*
 388 *on Graphics (TOG)*, 31(6):1–8, 2012.
- 389 Binh Huy Le and Zhigang Deng. Smooth skinning decomposition with rigid bones. *ACM Transactions on*
 390 *Graphics (TOG)*, 31(6):1–10, 2012.
- 391 Jiahui Lei and Kostas Daniilidis. Cadex: Learning canonical deformation coordinate space for dynamic surface
 392 representation via neural homeomorphism. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
 393 *and Pattern Recognition*, pages 6624–6634, 2022.
- 394 Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and
 395 expression from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017.
- 396 Yang Li, Hikari Takehara, Takafumi Taketomi, Bo Zheng, and Matthias Nießner. 4dcomplete: Non-rigid motion
 397 estimation beyond the observable surface. In *Proceedings of the IEEE/CVF International Conference on*
 398 *Computer Vision*, pages 12706–12716, 2021.
- 399 Yiyi Liao, Simon Donne, and Andreas Geiger. Deep marching cubes: Learning explicit surface representations.
 400 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2916–2925, 2018.
- 401 Cheng Lin, Changjian Li, Yuan Liu, Nenglun Chen, Yi-King Choi, and Wenping Wang. Point2skeleton: Learning
 402 skeletal representations from point clouds. In *Proceedings of the IEEE/CVF conference on computer vision*
 403 *and pattern recognition*, pages 4277–4286, 2021.
- 404 Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned
 405 multi-person linear model. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 851–866.
 406 2023.
- 407 William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm.
 408 In *Seminal graphics: pioneering efforts that shaped the field*, pages 347–353. 1998.
- 409 Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades, Gerard Pons-Moll, Siyu Tang, and Michael J Black.
 410 Learning to dress 3d people in generative clothing. In *Proceedings of the IEEE/CVF Conference on Computer*
 411 *Vision and Pattern Recognition*, pages 6469–6478, 2020.
- 412 Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy
 413 networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on*
 414 *computer vision and pattern recognition*, pages 4460–4470, 2019.
- 415 Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Occupancy flow: 4d reconstruction
 416 by learning particle dynamics. In *Proceedings of the IEEE/CVF international conference on computer vision*,
 417 pages 5379–5389, 2019.
- 418 Ahmed AA Osman, Timo Bolkart, and Michael J Black. Star: Sparse trained articulated human body regressor. In
 419 *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings,*
 420 *Part VI 16*, pages 598–613. Springer, 2020.
- 421 Pablo Palafox, Aljaž Božič, Justus Thies, Matthias Nießner, and Angela Dai. Npms: Neural parametric models
 422 for 3d deformable shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
 423 pages 12695–12705, 2021.

- 424 Junyi Pan, Xiaoguang Han, Weikai Chen, Jiapeng Tang, and Kui Jia. Deep mesh reconstruction from single rgb
 425 images via topology modification networks. In *Proceedings of the IEEE/CVF International Conference on*
 426 *Computer Vision*, pages 9964–9973, 2019.
- 427 Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning
 428 continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference*
 429 *on computer vision and pattern recognition*, pages 165–174, 2019.
- 430 Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d
 431 classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern*
 432 *recognition*, pages 652–660, 2017a.
- 433 Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning
 434 on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017b.
- 435 Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high
 436 resolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
 437 3577–3586, 2017.
- 438 Javier Romero, Dimitrios Tzionas, and Michael J Black. Embodied hands: Modeling and capturing hands and
 439 bodies together. *arXiv preprint arXiv:2201.02610*, 2022.
- 440 Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J Black. Scanimate: Weakly supervised learning of
 441 skinned clothed avatar networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
 442 *Pattern Recognition*, pages 2886–2897, 2021.
- 443 David Stutz and Andreas Geiger. Learning 3d shape completion from laser scan data with weak supervision. In
 444 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1955–1964, 2018.
- 445 Jiapeng Tang, Xiaoguang Han, Junyi Pan, Kui Jia, and Xin Tong. A skeleton-bridged deep learning approach for
 446 generating meshes of complex topologies from single rgb images. In *Proceedings of the ieee/cvf conference*
 447 *on computer vision and pattern recognition*, pages 4541–4550, 2019.
- 448 Jiapeng Tang, Xiaoguang Han, Mingkui Tan, Xin Tong, and Kui Jia. Skeletonnet: A topology-preserving
 449 solution for learning mesh reconstruction of object surfaces from rgb images. *IEEE transactions on pattern*
 450 *analysis and machine intelligence*, 44(10):6454–6471, 2021a.
- 451 Jiapeng Tang, Dan Xu, Kui Jia, and Lei Zhang. Learning parallel dense correspondence from spatio-temporal
 452 descriptors for efficient and robust 4d reconstruction. In *Proceedings of the IEEE/CVF Conference on*
 453 *Computer Vision and Pattern Recognition*, pages 6022–6031, 2021b.
- 454 Jiapeng Tang, Lev Markhasin, Bi Wang, Justus Thies, and Matthias Nießner. Neural shape deformation priors.
 455 *Advances in Neural Information Processing Systems*, 35:17117–17132, 2022.
- 456 Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional
 457 architectures for high-resolution 3d outputs. In *Proceedings of the IEEE international conference on computer*
 458 *vision*, pages 2088–2096, 2017.
- 459 Peng-Shuai Wang, Chun-Yu Sun, Yang Liu, and Xin Tong. Adaptive o-cnn: A patch-based deep representation
 460 of 3d shapes. *ACM Transactions on Graphics (TOG)*, 37(6):1–11, 2018.
- 461 Rich Wareham and Joan Lasenby. Bone glow: An improved method for the assignment of weights for mesh
 462 deformation. In *Articulated Motion and Deformable Objects: 5th International Conference, AMDO 2008,*
 463 *Port d'Andratx, Mallorca, Spain, July 9–11, 2008. Proceedings 5*, pages 63–71. Springer, 2008.
- 464 Zhan Xu, Yang Zhou, Evangelos Kalogerakis, Chris Landreth, and Karan Singh. Rignet: Neural rigging for
 465 articulated characters. *arXiv preprint arXiv:2005.00559*, 2020.
- 466 Yuxin Yao, Zhi Deng, and Junhui Hou. Riggs: Rigging of 3d gaussians for modeling articulated objects in
 467 videos. *arXiv preprint arXiv:2503.16822*, 2025a.
- 468 Yuxin Yao, Siyu Ren, Junhui Hou, Zhi Deng, Juyong Zhang, and Wenping Wang. Dynosurf: Neural deformation-
 469 based temporally consistent dynamic surface reconstruction. In *European Conference on Computer Vision*,
 470 pages 271–288. Springer, 2025b.
- 471 Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of*
 472 *the IEEE/CVF international conference on computer vision*, pages 16259–16268, 2021.

- 473 Kaifeng Zou, Sylvain Faisan, Boyang Yu, Sébastien Valette, and Hyewon Seo. 4d facial expression diffusion
474 model. *ACM Transactions on Multimedia Computing, Communications and Applications*, 2023.
- 475 Silvia Zuffi, Angjoo Kanazawa, David W Jacobs, and Michael J Black. 3d menagerie: Modeling the 3d shape
476 and pose of animals. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
477 pages 6365–6373, 2017.

478 **NeurIPS Paper Checklist**

479 **1. Claims**

480 Question: Do the main claims made in the abstract and introduction accurately reflect the
481 paper's contributions and scope?

482 Answer: [Yes]

483 Justification: The abstract and introduction clearly outline the contributions, including the
484 AniDynRecon framework, HSM, CMS strategy, and superior performance on benchmarks.
485 These claims are substantiated in Sec. 4 (Method), 5 (Experiments), and Fig. 2.

486 Guidelines:

- 487 • The answer NA means that the abstract and introduction do not include the claims
488 made in the paper.
- 489 • The abstract and/or introduction should clearly state the claims made, including the
490 contributions made in the paper and important assumptions and limitations. A No or
491 NA answer to this question will not be perceived well by the reviewers.
- 492 • The claims made should match theoretical and experimental results, and reflect how
493 much the results can be expected to generalize to other settings.
- 494 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
495 are not attained by the paper.

496 **2. Limitations**

497 Question: Does the paper discuss the limitations of the work performed by the authors?

498 Answer: [Yes]

499 Justification: The paper explicitly acknowledges limitations in its Conclusion (Sec. 6),
500 noting that the spanning tree algorithm can produce incorrect skeleton structures in some
501 cases and that addressing these failures is left to future work.

502 Guidelines:

- 503 • The answer NA means that the paper has no limitation while the answer No means that
504 the paper has limitations, but those are not discussed in the paper.
- 505 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 506 • The paper should point out any strong assumptions and how robust the results are to
507 violations of these assumptions (e.g., independence assumptions, noiseless settings,
508 model well-specification, asymptotic approximations only holding locally). The authors
509 should reflect on how these assumptions might be violated in practice and what the
510 implications would be.
- 511 • The authors should reflect on the scope of the claims made, e.g., if the approach was
512 only tested on a few datasets or with a few runs. In general, empirical results often
513 depend on implicit assumptions, which should be articulated.
- 514 • The authors should reflect on the factors that influence the performance of the approach.
515 For example, a facial recognition algorithm may perform poorly when image resolution
516 is low or images are taken in low lighting. Or a speech-to-text system might not be
517 used reliably to provide closed captions for online lectures because it fails to handle
518 technical jargon.
- 519 • The authors should discuss the computational efficiency of the proposed algorithms
520 and how they scale with dataset size.
- 521 • If applicable, the authors should discuss possible limitations of their approach to
522 address problems of privacy and fairness.
- 523 • While the authors might fear that complete honesty about limitations might be used by
524 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
525 limitations that aren't acknowledged in the paper. The authors should use their best
526 judgment and recognize that individual actions in favor of transparency play an impor-
527 tant role in developing norms that preserve the integrity of the community. Reviewers
528 will be specifically instructed to not penalize honesty concerning limitations.

529 **3. Theory assumptions and proofs**

530 Question: For each theoretical result, does the paper provide the full set of assumptions and
531 a complete (and correct) proof?

532 Answer: [NA]

533 Justification: The paper focuses on methodological contributions (AniDynRecon framework,
534 HSM, CMS) and empirical validation. It does not introduce theoretical results requiring for-
535 mal proofs, but instead describes algorithmic components, loss functions, and experimental
536 evaluations.

537 Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

548 4. Experimental result reproducibility

549 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
550 perimental results of the paper to the extent that it affects the main claims and/or conclusions
551 of the paper (regardless of whether the code and data are provided or not)?

552 Answer: [Yes]

553 Justification: The paper details datasets (D-FAUST, DT4D-A), data splits, baselines, evalua-
554 tion metrics (IoU, CD, F-score), and implementation specifics (PointNet/PointNet++ usage,
555 GPU setup). Key components like loss functions (Sec. 4), CMS strategy (Sec. 4.3), and
556 training phases (Sec. 5.1, 5.2) are sufficiently described. However, hyperparameters (e.g.,
557 learning rates) and full optimization details are partially covered, relying on future code
558 release for full reproducibility.

559 Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.

582 (c) If the contribution is a new model (e.g., a large language model), then there should
583 either be a way to access this model for reproducing the results or a way to reproduce
584 the model (e.g., with an open-source dataset or instructions for how to construct
585 the dataset).

586 (d) We recognize that reproducibility may be tricky in some cases, in which case
587 authors are welcome to describe the particular way they provide for reproducibility.
588 In the case of closed-source models, it may be that access to the model is limited in
589 some way (e.g., to registered users), but it should be possible for other researchers
590 to have some path to reproducing or verifying the results.

591 5. Open access to data and code

592 Question: Does the paper provide open access to the data and code, with sufficient instruc-
593 tions to faithfully reproduce the main experimental results, as described in supplemental
594 material?

595 Answer: [No]

596 Justification: Our code will be made publicly available upon acceptance.

597 Guidelines:

- 598 • The answer NA means that paper does not include experiments requiring code.
- 599 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 600 • While we encourage the release of code and data, we understand that this might not be
601 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
602 including code, unless this is central to the contribution (e.g., for a new open-source
603 benchmark).
- 604 • The instructions should contain the exact command and environment needed to run to
605 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 606 • The authors should provide instructions on data access and preparation, including how
607 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 608 • The authors should provide scripts to reproduce all experimental results for the new
609 proposed method and baselines. If only a subset of experiments are reproducible, they
610 should state which ones are omitted from the script and why.
- 611 • At submission time, to preserve anonymity, the authors should release anonymized
612 versions (if applicable).
- 613 • Providing as much information as possible in supplemental material (appended to the
614 paper) is recommended, but including URLs to data and code is permitted.

617 6. Experimental setting/details

618 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
619 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
620 results?

621 Answer: [Yes]

622 Justification: The paper gives full experimental details, including data splits and GPU
623 configuration in Sec. 5 and Supplementary, and specifies all key training settings, learning
624 rate, batch size, number of epochs, plus hyperparameter values and their selection procedure
625 in Supplementary.

626 Guidelines:

- 627 • The answer NA means that the paper does not include experiments.
- 628 • The experimental setting should be presented in the core of the paper to a level of detail
629 that is necessary to appreciate the results and make sense of them.
- 630 • The full details can be provided either with the code, in appendix, or as supplemental
631 material.

632 7. Experiment statistical significance

633 Question: Does the paper report error bars suitably and correctly defined or other appropriate
634 information about the statistical significance of the experiments?

635 Answer: [No]

636 Justification: All reported metrics in the tables and figures are presented as single-point
637 values without any accompanying error bars, confidence intervals, or description of statistical
638 significance tests.

639 Guidelines:

- 640 • The answer NA means that the paper does not include experiments.
- 641 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
642 dence intervals, or statistical significance tests, at least for the experiments that support
643 the main claims of the paper.
- 644 • The factors of variability that the error bars are capturing should be clearly stated (for
645 example, train/test split, initialization, random drawing of some parameter, or overall
646 run with given experimental conditions).
- 647 • The method for calculating the error bars should be explained (closed form formula,
648 call to a library function, bootstrap, etc.)
- 649 • The assumptions made should be given (e.g., Normally distributed errors).
- 650 • It should be clear whether the error bar is the standard deviation or the standard error
651 of the mean.
- 652 • It is OK to report 1-sigma error bars, but one should state it. The authors should
653 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
654 of Normality of errors is not verified.
- 655 • For asymmetric distributions, the authors should be careful not to show in tables or
656 figures symmetric error bars that would yield results that are out of range (e.g. negative
657 error rates).
- 658 • If error bars are reported in tables or plots, The authors should explain in the text how
659 they were calculated and reference the corresponding figures or tables in the text.

660 8. Experiments compute resources

661 Question: For each experiment, does the paper provide sufficient information on the com-
662 puter resources (type of compute workers, memory, time of execution) needed to reproduce
663 the experiments?

664 Answer: [Yes]

665 Justification: The supplementary material (Implementation Details) reports that training
666 uses 16 compute workers, occupies about 22GB of GPU memory, and requires 80 minutes
667 per epoch for human bodies and 100 minutes per epoch for animals over 100 epochs.
668 Additionally, All runs were conducted on a single NVIDIA RTX4090 GPU for 10 days.

669 Guidelines:

- 670 • The answer NA means that the paper does not include experiments.
- 671 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
672 or cloud provider, including relevant memory and storage.
- 673 • The paper should provide the amount of compute required for each of the individual
674 experimental runs as well as estimate the total compute.
- 675 • The paper should disclose whether the full research project required more compute
676 than the experiments reported in the paper (e.g., preliminary or failed experiments that
677 didn't make it into the paper).

678 9. Code of ethics

679 Question: Does the research conducted in the paper conform, in every respect, with the
680 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

681 Answer: [Yes]

682 Justification: This work is purely foundational, involves no human subjects or sensitive
683 data, poses no dual-use or safety risks, and the authors have affirmed in the NeurIPS Paper
684 Checklist that it conforms fully to the NeurIPS Code of Ethics.

685 Guidelines:

- 686 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.

- 687 • If the authors answer No, they should explain the special circumstances that require a
688 deviation from the Code of Ethics.
689 • The authors should make sure to preserve anonymity (e.g., if there is a special consider-
690 ation due to laws or regulations in their jurisdiction).

691 **10. Broader impacts**

692 Question: Does the paper discuss both potential positive societal impacts and negative
693 societal impacts of the work performed?

694 Answer: [Yes]

695 Justification: The paper have discussed potential societal impacts in supplementary material.

696 Guidelines:

- 697 • The answer NA means that there is no societal impact of the work performed.
698 • If the authors answer NA or No, they should explain why their work has no societal
699 impact or why the paper does not address societal impact.
700 • Examples of negative societal impacts include potential malicious or unintended uses
701 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
702 (e.g., deployment of technologies that could make decisions that unfairly impact specific
703 groups), privacy considerations, and security considerations.
704 • The conference expects that many papers will be foundational research and not tied
705 to particular applications, let alone deployments. However, if there is a direct path to
706 any negative applications, the authors should point it out. For example, it is legitimate
707 to point out that an improvement in the quality of generative models could be used to
708 generate deepfakes for disinformation. On the other hand, it is not needed to point out
709 that a generic algorithm for optimizing neural networks could enable people to train
710 models that generate Deepfakes faster.
711 • The authors should consider possible harms that could arise when the technology is
712 being used as intended and functioning correctly, harms that could arise when the
713 technology is being used as intended but gives incorrect results, and harms following
714 from (intentional or unintentional) misuse of the technology.
715 • If there are negative societal impacts, the authors could also discuss possible mitigation
716 strategies (e.g., gated release of models, providing defenses in addition to attacks,
717 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
718 feedback over time, improving the efficiency and accessibility of ML).

719 **11. Safeguards**

720 Question: Does the paper describe safeguards that have been put in place for responsible
721 release of data or models that have a high risk for misuse (e.g., pretrained language models,
722 image generators, or scraped datasets)?

723 Answer: [NA]

724 Justification: The paper does not introduce or release any high-risk pretrained models,
725 generative systems, or scraped datasets.

726 Guidelines:

- 727 • The answer NA means that the paper poses no such risks.
728 • Released models that have a high risk for misuse or dual-use should be released with
729 necessary safeguards to allow for controlled use of the model, for example by requiring
730 that users adhere to usage guidelines or restrictions to access the model or implementing
731 safety filters.
732 • Datasets that have been scraped from the Internet could pose safety risks. The authors
733 should describe how they avoided releasing unsafe images.
734 • We recognize that providing effective safeguards is challenging, and many papers do
735 not require this, but we encourage authors to take this into account and make a best
736 faith effort.

737 **12. Licenses for existing assets**

738 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
739 the paper, properly credited and are the license and terms of use explicitly mentioned and
740 properly respected?

741 Answer: [Yes]

742 Justification: The paper credits all external assets used in the experiments, including datasets
743 like D-FAUST and DT4D-A; citations are provided in the references section, and no violations
744 of license terms in the text or supplementary material.

745 Guidelines:

- 746 • The answer NA means that the paper does not use existing assets.
- 747 • The authors should cite the original paper that produced the code package or dataset.
- 748 • The authors should state which version of the asset is used and, if possible, include a
749 URL.
- 750 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 751 • For scraped data from a particular source (e.g., website), the copyright and terms of
752 service of that source should be provided.
- 753 • If assets are released, the license, copyright information, and terms of use in the
754 package should be provided. For popular datasets, paperswithcode.com/datasets
755 has curated licenses for some datasets. Their licensing guide can help determine the
756 license of a dataset.
- 757 • For existing datasets that are re-packaged, both the original license and the license of
758 the derived asset (if it has changed) should be provided.
- 759 • If this information is not available online, the authors are encouraged to reach out to
760 the asset's creators.

761 **13. New assets**

762 Question: Are new assets introduced in the paper well documented and is the documentation
763 provided alongside the assets?

764 Answer: [NA]

765 Justification: The paper does not introduce or release any new datasets, code libraries,
766 or pretrained model weights as standalone assets; accordingly, there is no accompanying
767 documentation provided.

768 Guidelines:

- 769 • The answer NA means that the paper does not release new assets.
- 770 • Researchers should communicate the details of the dataset/code/model as part of their
771 submissions via structured templates. This includes details about training, license,
772 limitations, etc.
- 773 • The paper should discuss whether and how consent was obtained from people whose
774 asset is used.
- 775 • At submission time, remember to anonymize your assets (if applicable). You can either
776 create an anonymized URL or include an anonymized zip file.

777 **14. Crowdsourcing and research with human subjects**

778 Question: For crowdsourcing experiments and research with human subjects, does the paper
779 include the full text of instructions given to participants and screenshots, if applicable, as
780 well as details about compensation (if any)?

781 Answer: [NA]

782 Justification: The submission contains no crowdsourcing studies or other human-subject
783 experiments. There are no participant instructions, screenshots, or compensation details
784 because no such research was conducted.

785 Guidelines:

- 786 • The answer NA means that the paper does not involve crowdsourcing nor research with
787 human subjects.
- 788 • Including this information in the supplemental material is fine, but if the main contribu-
789 tion of the paper involves human subjects, then as much detail as possible should be
790 included in the main paper.
- 791 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
792 or other labor should be paid at least the minimum wage in the country of the data
793 collector.

794 **15. Institutional review board (IRB) approvals or equivalent for research with human
795 subjects**

796 Question: Does the paper describe potential risks incurred by study participants, whether
797 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
798 approvals (or an equivalent approval/review based on the requirements of your country or
799 institution) were obtained?

800 Answer: [NA]

801 Justification: The paper does not involve any experiments with human subjects or crowd-
802 sourced participants, and thus there are no associated risks, disclosures, or IRB approvals
803 required or mentioned in the submission.

804 Guidelines:

- 805 • The answer NA means that the paper does not involve crowdsourcing nor research with
806 human subjects.
- 807 • Depending on the country in which research is conducted, IRB approval (or equivalent)
808 may be required for any human subjects research. If you obtained IRB approval, you
809 should clearly state this in the paper.
- 810 • We recognize that the procedures for this may vary significantly between institutions
811 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
812 guidelines for their institution.
- 813 • For initial submissions, do not include any information that would break anonymity (if
814 applicable), such as the institution conducting the review.

815 **16. Declaration of LLM usage**

816 Question: Does the paper describe the usage of LLMs if it is an important, original, or
817 non-standard component of the core methods in this research? Note that if the LLM is used
818 only for writing, editing, or formatting purposes and does not impact the core methodology,
819 scientific rigorosity, or originality of the research, declaration is not required.

820 Answer: [NA]

821 Justification: According to the NeurIPS 2025 LLM policy and the Paper Checklist guide-
822 lines, LLM is used solely for writing, editing, or formatting and does not affect the core
823 methodology, declaration is not required.

824 Guidelines:

- 825 • The answer NA means that the core method development in this research does not
826 involve LLMs as any important, original, or non-standard components.
- 827 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
828 for what should or should not be described.