
AniDynRecon: Animatable 4D Dynamics Reconstruction from Sparse Point Clouds

Supplementary Document

Anonymous Author(s)

Affiliation
Address
email

1 A Notations

- 2 In this paper, we use notations summarized in Tab. S1. K denotes the number of joints, and N
3 represents the number of surface points. M refers to the number of components used for the Gaussian
4 Mixture Model (GMM). x indicates the position of a surface point, while p represents the position of
5 a joint. Δp is the delta position between a joint at the frame of rest pose and another frame. f stands
6 for a feature vector of rest pose encoded by PointNet++ Qi et al. (2017). \hat{w} is the output of the joint
7 weight MLP, and \tilde{w} represents the joint weight. w denotes the skinning weight. d is the offset vector
8 from a joint to a surface point, and $|d|$ is the distance between them. R represents the rotation matrix
9 of a joint, and r^t is the rotation of a joint represented as lie algebra. S stands for shape tokens. O
10 denotes the occupancy value. \bar{V} represents the vertices of a mesh, and \bar{F} refers to the faces of a mesh.

Notation	Type	Definition
K	scalar	number of joints
N	scalar	number of points
M	scalar	number of components in the GMM
x	vector	position of point
p	vector	position of joint
Δp	vector	delta position of joint
f	vector	feature vector
\hat{w}	scalar	output of joint weight MLP
\tilde{w}	scalar	joint weight
w	scalar	skinning weight
d	vector	offset from joint to surface point
$ d $	scalar	distance from joint to surface point
R	matrix	rotation matrix of joint
r^t	vector	rotation of a joint at time t
S	vector	shape token
O	scalar	occupancy value
\bar{V}	vector set	vertices of mesh
\bar{F}	vector set	faces of mesh

Table S1: Symbols and notations used in this paper.

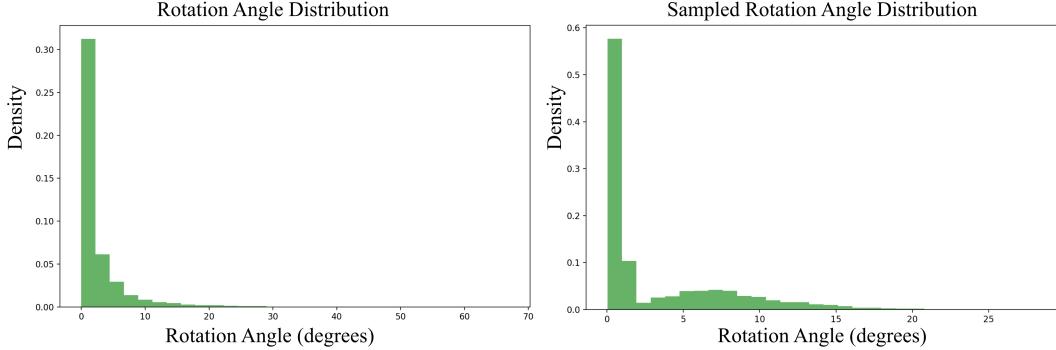


Figure S2: (**Left**). Rotation angle distribution in original training data. (**Right**). Sampled rotation angle distribution using GMM.

11 B Implementation Details

12 B.1 Final Training Objective

13 Our model is trained end-to-end, obtaining high quality when joint position and skinning weight field
14 are jointly trained without any supervision under ground-truth. The total training objective of our
15 network is:

$$\mathcal{L} = \lambda_{\text{locality}} \mathcal{L}_{\text{locality}} + \lambda_{\text{dynamics}} \mathcal{L}_{\text{dynamics}} + \lambda_{\text{joints_inside}} \mathcal{L}_{\text{joints_inside}} + \lambda_{\text{centrality}} \mathcal{L}_{\text{centrality}} + \lambda_{\text{shape}} \mathcal{L}_{\text{shape}} \quad (1)$$

16 We set the hyperparameters as follows: $\lambda_{\text{locality}} = 0.06$, $\lambda_{\text{dynamics}} = 0.001$, $\lambda_{\text{joints_inside}} =$
17 0.01 , $\lambda_{\text{centrality}} = 3.0$, $\lambda_{\text{shape}} = 1.0$. These values were determined based on the reconstructed
18 IoU metric in our early experiments on partial datasets. In the experiments, we only adjusted the
19 value of one weight at a time to determine these weights so that the network converges stably. Note
20 that the Chamfer distance and the correspondence error reported in our paper, as well as Lei and
21 Daniilidis (2022), are multiplied by 10.

22 B.2 Details of Skinning Weight Decoder

23 B.3 Network Settings

24 We design our skinning weight decoder \mathcal{W}
25 based on the observation that cross-attention be-
26 tween joints and vertices is similar to the blend-
27 ing weights in LBS. The architecture of the skin-
28 ning weight decoder is depicted in Fig. S1. Sur-
29 face points are first sampled via furthest point
30 sampling (FPS). Then, positional embeddings
31 (PE) of these points and subsampled points are
32 fed into a cross-attention block to produce shape
33 tokens \mathcal{S} . Subsequently, a self-attention block
34 encodes the PE of joints and shape tokens. Then,
35 another cross-attention block computes skinning
36 weights using the encoded features as key and
37 value, and the PE of surface points as query.

38 We leverage the 8-th frame of the input sequence
39 as rest pose to reconstruct the shape of object.
40 We follow the setting of Motion2Vecsets Cao
41 et al. (2024): the shape encoder extracts rest-
42 pose shape tokens from 2048 points randomly sam-
43 pled from the object’s surface and near-surface
regions.

44 In the implementation of *Hierarchical Skeletal Modeling* (HSM) of our network, we utilize a series
45 of set abstraction (SA) components to process and extract features from the selected 8-th frame point

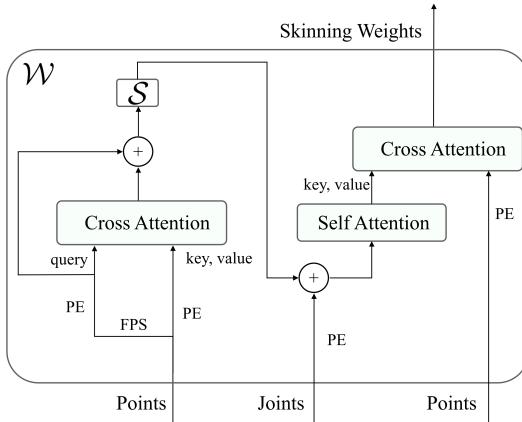


Figure S1: Architecture of our skinning weight decoder.

44 In the implementation of *Hierarchical Skeletal Modeling* (HSM) of our network, we utilize a series
45 of set abstraction (SA) components to process and extract features from the selected 8-th frame point

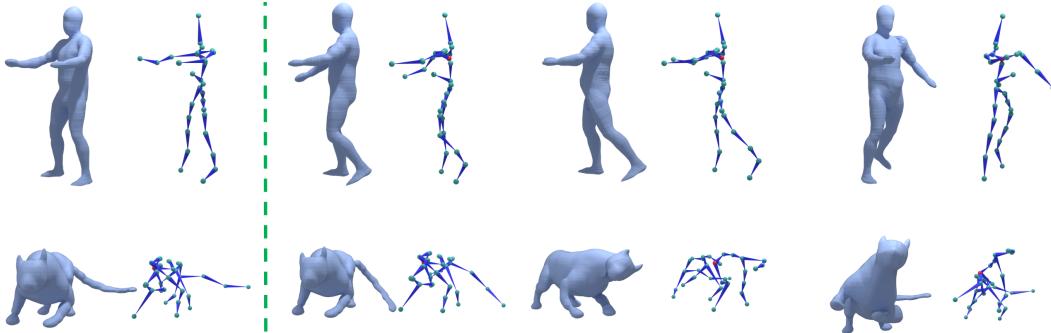


Figure S3: Editing on the D-FAUST Bogo et al. (2017) and DT4D Li et al. (2021) datasets. The left shows the rest poses with corresponding skeletons and the right illustrates the edited poses.

46 cloud data. For the first module, we specify radius ranges of [0.1, 0.2] and sample numbers of [16, 32].
 47 The multilayer perceptron (MLP) list is configured with two branches, each following the structure
 48 $16 \rightarrow 16 \rightarrow 32$. The subsequent modules gradually increase the radius ranges and sample numbers
 49 and use deeper MLPs. These modules are used sequentially for hierarchical feature extraction, and
 50 then the joint weight MLP will decode weights from these features.

51 C Experiment Details

52 **Dataset Details:** The DFAUST Bogo et al. (2017) dataset we used for human bodies contains 10
 53 subjects and 129 sequences. The DT4D Li et al. (2021) dataset for animals contains 38 identities
 54 with a total of 1227 animations.
 55 **Training Details:** We train our framework using the ADAM optimizer with a learning rate $\gamma = 10^{-4}$
 56 and batch size of 18, which occupies about 22 GB of GPU memory with 16 compute workers. The
 57 learning rate has a step decay with a decay rate of 0.3. The training time of each epoch is 80 minutes
 58 for human bodies and 100 minutes for animals. We train 100 epochs for both.

59 C.1 Data Augmentation

60 Our motion synthesis pipeline begins with fitting a Gaussian Mixture Model (GMM) to model
 61 the rotation Euler angles ρ of reconstructed joints based on the original training dataset. The
 62 rotation angle distribution of the original training data is shown in Fig. S2. The model complexity
 63 ($M = 2$ components) was determined through silhouette analysis to ensure optimal cluster separation.
 64 Following parameter estimation, we sample 10000 rotations and visualize their distribution in Fig. S2.
 65 The rotations ρ for data augmentation strategy are sampled from such a distribution and randomly
 66 applied to all reconstructed joints.

67 C.2 More Results

68 **More Results for Motion Editing.** The inside regularization of joints and the locality constraint
 69 of skinning weight play a pivotal role in guaranteeing accurate post-editing motion. Our proposed
 70 HSM module significantly contributes to achieving an optimal joint distribution, thereby greatly
 71 facilitating subsequent build of skeleton tree and motion editing tasks. When generating a novel pose,
 72 it merely requires rotating a single joint within the skeletal hierarchy. Subsequently, the rotation of
 73 the edited joint is propagated to its corresponding child nodes via the forward kinematics technique,
 74 as schematically illustrated in Fig. S3.

75 **More Results for 4D Reconstruction.** Fig. S4 shows the qualitative comparison on D-FAUST Bogo
 76 et al. (2017) human bodies. Fig. S5 shows the comparison on DeformingThings4D-Animals Li et al.
 77 (2021) dataset. More animations are shown in the **Video Demo**.

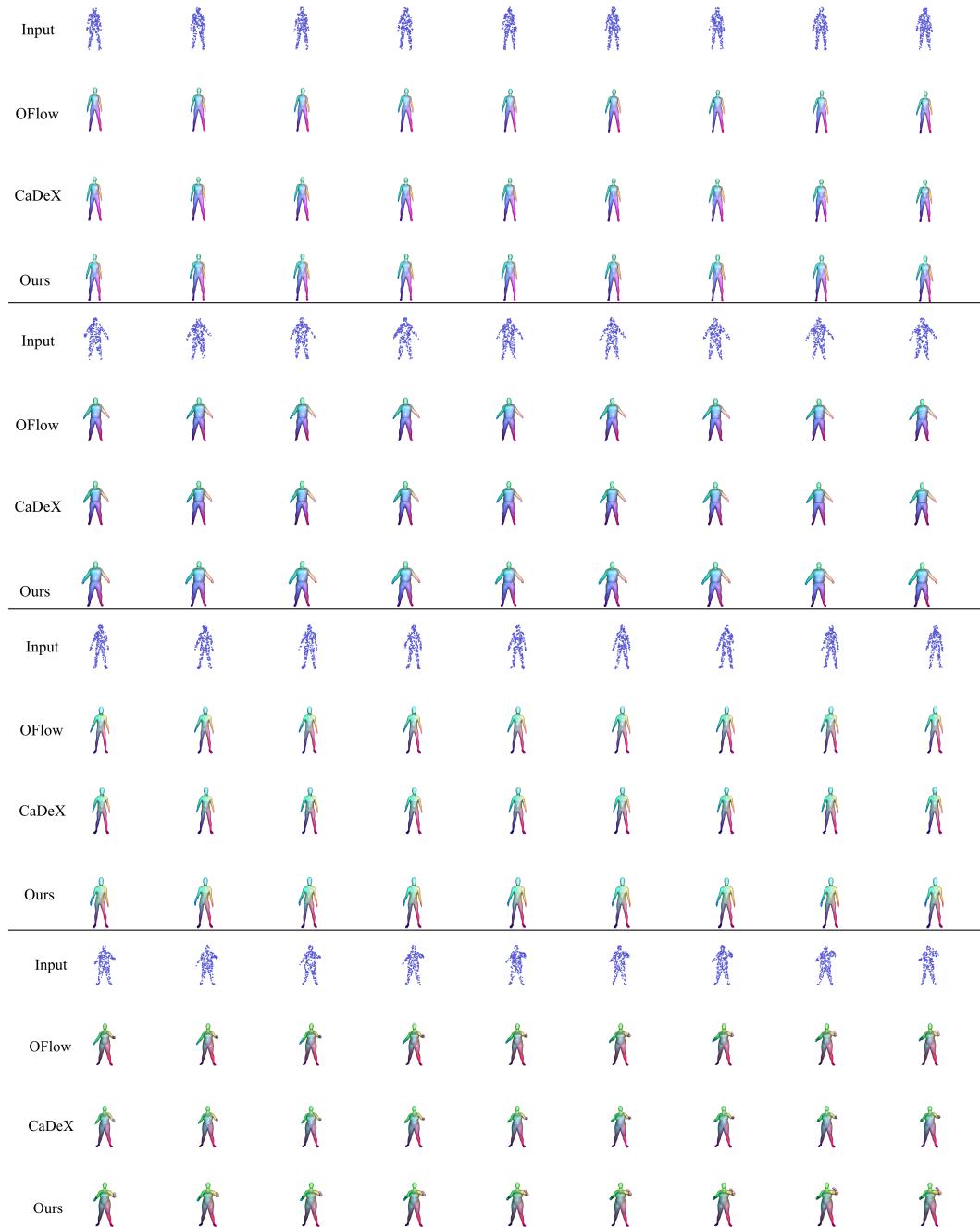
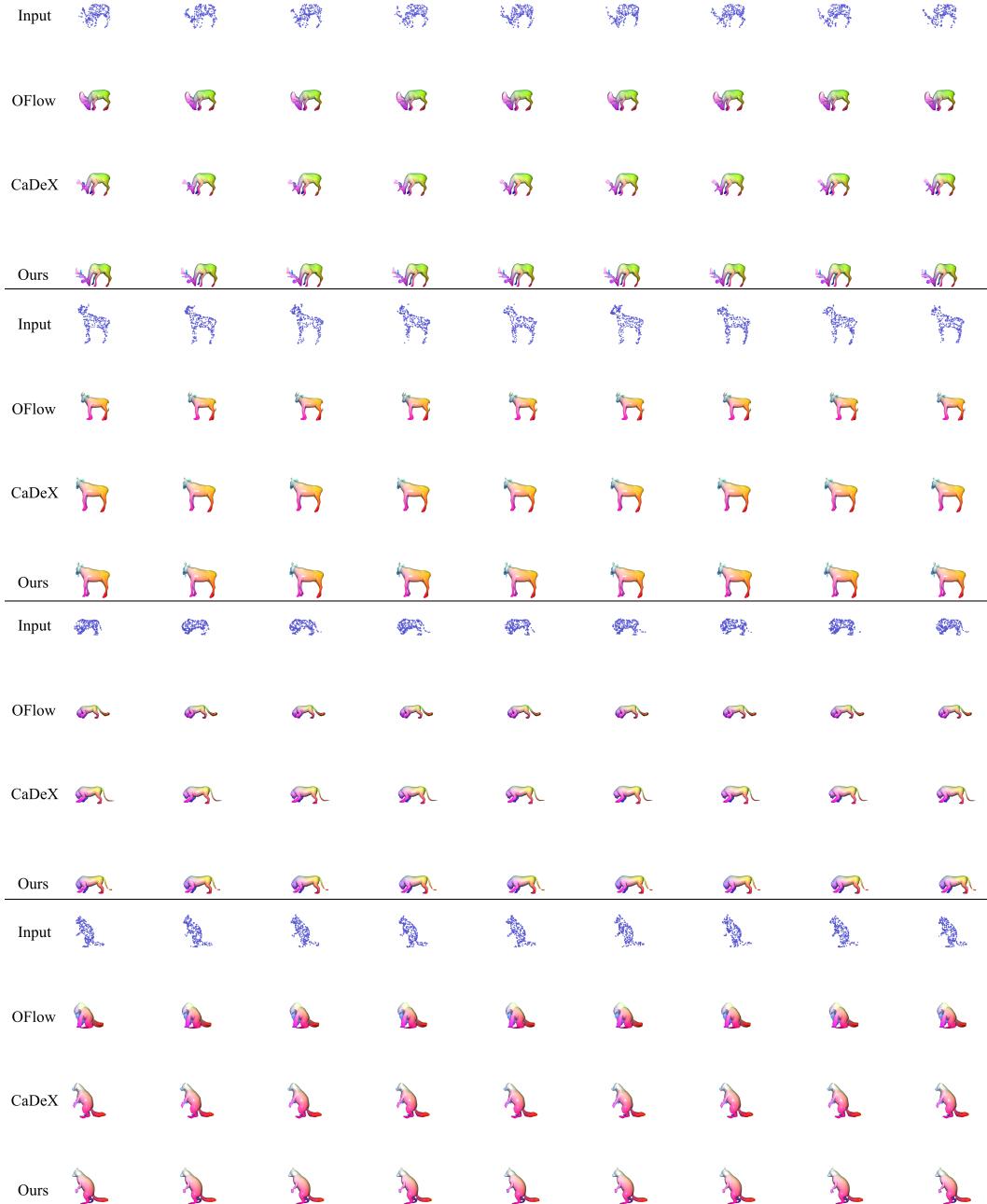


Figure S4: Modeling human bodies: each column corresponds to a time frame of the reconstruction.



78 **D Societal Benefits**

79 Our 4D reconstruction framework offers clear societal benefits: it can accelerate applications in
80 medical visualization (e.g. reconstructing patient-specific anatomies for surgical planning) and digital
81 heritage preservation (efficiently digitizing artifacts). However, there are potential downsides: the
82 same techniques could be repurposed to reconstruct people’s shapes from incomplete sensor data,
83 raising privacy and surveillance concerns.

84 **References**

- 85 Federica Bogo, Javier Romero, Gerard Pons-Moll, and Michael J Black. Dynamic faust: Registering human
86 bodies in motion. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
87 6233–6242, 2017.
- 88 Wei Cao, Chang Luo, Biao Zhang, Matthias Nießner, and Jiapeng Tang. Motion2vecsets: 4d latent vector set
89 diffusion for non-rigid shape reconstruction and tracking. In *Proceedings of the IEEE/CVF Conference on*
90 *Computer Vision and Pattern Recognition*, pages 20496–20506, 2024.
- 91 Jiahui Lei and Kostas Daniilidis. Cadex: Learning canonical deformation coordinate space for dynamic surface
92 representation via neural homeomorphism. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
93 and Pattern Recognition, pages 6624–6634, 2022.
- 94 Yang Li, Hikari Takehara, Takafumi Takiomi, Bo Zheng, and Matthias Nießner. 4dcomplete: Non-rigid motion
95 estimation beyond the observable surface. In *Proceedings of the IEEE/CVF International Conference on*
96 *Computer Vision*, pages 12706–12716, 2021.
- 97 Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning
98 on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017.