

Frame Averaging for Equivariant Shape Space Learning

A. Additional Details

Here we provide additional details for each experiment in section 5.

A.1. Global Euclidean Mesh \rightarrow mesh.

Architecture. We start by describing in full details the backbone architectures for both the encoder ϕ and decoder ψ , for all the methods in table 1. Then, we describe the relevant hyperparameters for our Frame Averaging version. Note that the backbone architectures are similar to the ones in [5], and are described here again for completeness. The network consists of the following layers:

$$\begin{aligned} \text{Conv}(n, d_{\text{in}}, d_{\text{out}}) &: \mathbb{R}^{n \times d_{\text{in}}} \rightarrow \mathbb{R}^{n \times d_{\text{out}}} \\ \text{Pool}(n, n') &: \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n' \times d} \\ \text{Linear}(d_{\text{in}}, d_{\text{out}}) &: \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}. \end{aligned}$$

The $\text{Conv}(n, d_{\text{in}}, d_{\text{out}})$ is the Chebyshev spectral graph convolutional operator [3] as implemented in [5, 12], where n denotes the number of vertices, and $d_{\text{in}}, d_{\text{out}}$ are the layer’s input and output features dimensions respectively. All the $\text{Conv}(n, d_{\text{in}}, d_{\text{out}})$ layers have 6 Chebyshev polynomials. The $\text{Pool}(n, n')$ layer is based on the mesh down/up sampling schemes defined in [12], where we used the same mesh samples as provided in [5] official implementation. The n, n' are the number of vertices in the input and the output, respectively. The linear layer is defined by $\mathbf{X}' = \mathbf{W}\mathbf{X} + \mathbf{b}$, where $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ and $\mathbf{b} \in \mathbb{R}^{d_{\text{out}}}$ are learnable parameters.

Moreover, the backbone architecture consists of the following blocks:

$$\begin{aligned} \text{EnBlock}(n, n', d_{\text{in}}, d_{\text{out}}) &: \mathbf{X} \in \mathbb{R}^{n \times d_{\text{in}}} \mapsto \mathbf{X}' \in \mathbb{R}^{n \times d_{\text{out}}} \\ \mathbf{X}' &= \text{Pool}(n, n', d_{\text{out}}) (\sigma (\text{Conv}(n, d_{\text{in}}, d_{\text{out}}))) \\ \text{DeBlock}(n, n', d_{\text{in}}, d_{\text{out}}) &: \mathbf{X} \in \mathbb{R}^{n \times d_{\text{in}}} \mapsto \mathbf{X}' \in \mathbb{R}^{n \times d_{\text{out}}} \\ \mathbf{X}' &= \sigma (\text{Conv}(n', d_{\text{in}}, d_{\text{out}}) (\text{Pool}(n, n', d_{\text{in}})(\mathbf{X}))) \end{aligned}$$

where σ is the ELU activation function. Then, the encoder

$\phi : \mathbb{R}^{6890 \times 3} \rightarrow \mathbb{R}^{72}$ consists of the following blocks:

$$\begin{aligned} &\text{EnBlock}(6890, 3445, 3, 32) \rightarrow \\ &\text{EnBlock}(3445, 1723, 32, 32) \rightarrow \\ &\text{EnBlock}(1723, 862, 32, 32) \rightarrow \\ &\text{EnBlock}(862, 431, 32, 64) \rightarrow \text{Linear}(27584, 72) \rightarrow . \end{aligned}$$

Similarly, the decoder $\psi : \mathbb{R}^{72} \rightarrow \mathbb{R}^{n \times 3}$ is consisted of the following blocks:

$$\begin{aligned} &\text{Linear}(72, 27584) \rightarrow \text{DeBlock}(431, 862, 64, 64) \rightarrow \\ &\text{DeBlock}(862, 1723, 64, 32) \rightarrow \\ &\text{DeBlock}(1723, 3445, 32, 32) \rightarrow \\ &\text{DeBlock}(3445, 6890, 32, 32) \rightarrow \\ &\text{Conv}(6890, 6890, 32, 3) \rightarrow . \end{aligned}$$

For the Frame Averaging versions of these backbone architectures, $\Phi = \langle \phi \rangle_{\mathcal{F}}$, and $\Psi = \langle \psi \rangle_{\mathcal{F}}$, we set $m = 0$ and $d = 24$. For the frames construction of Φ and Ψ , $\mathcal{F} : V \rightarrow 2^{E(3)} \setminus \emptyset$, we use the construction defined in section 3.3, with $\mathbf{w} = \mathbf{1}$. Note that for Φ the frames are calculated with respect to the input vertices $\mathbf{V} = \mathbf{X} \in \mathbb{R}^{6890 \times 3}$, and for Ψ the frames are calculated with respect to $\mathbf{V} = \mathbf{Z} \in \mathbb{R}^{24 \times 3}$.

Implementation details. For all methods in table 1, training was done using the ADAM [6] optimizer, with batch size 64. The number of epochs was set to 150, with 0.0001 fixed learning rate. Training was done on 2 Nvidia V-100 GPUs.

Evaluation. The quantitative measure reported in table 1, the average per-vertex distance error (MSE), is based on the official implementation of [5]. Given a collection of N test examples $\{\mathbf{X}_i\}_{i=1}^N$, and corresponding predictions $\{\mathbf{Y}_i\}_{i=1}^N$, the metric is defined by

$$\text{MSE} = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \|\mathbf{X}_{ij} - \mathbf{Y}_{ij}\|, \quad (1)$$

where M is the number of vertices.

A.2. Global Euclidean: Point cloud \rightarrow implicit

Architecture. In this experiment, the backbone encoder architecture is the same for the baseline VAE and our Frame Averaging version. The network is based on the design and implementation of OccupancyNetworks [9]. The encoder network consists of layers of the form

$$\text{FC}(d_{\text{in}}, d_{\text{out}}) : \mathbf{X} \in \mathbb{R}^{n \times d_{\text{in}}} \mapsto \nu(\mathbf{X}\mathbf{W} + \mathbf{1}\mathbf{b}^T)$$

$$\text{MaxPool} : \mathbf{X} \in \mathbb{R}^{n \times d} \mapsto [\max \mathbf{X}e_i]$$

$$\text{Linear}(d_{\text{in}}, d_{\text{out}}) : \mathbf{V} \in \mathbb{R}^{d_{\text{in}}} \mapsto \mathbf{W}'\mathbf{V} + \mathbf{b}'$$

where $\mathbf{W} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}}$, $\mathbf{W}' \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, $\mathbf{b} \in \mathbb{R}^{d_{\text{out}}}$, $\mathbf{b}' \in \mathbb{R}^{d_{\text{out}}}$ are the learnable parameters, $\mathbf{1} \in \mathbb{R}^n$ is the vector of all ones, $[\cdot]$ is the concatenation operator, e_i is the standard basis in \mathbb{R}^d , and ν is the ReLU activation. Then, we define the following block:

$$\text{EnBlock}(d_{\text{in}}, d_{\text{out}}) : \mathbf{X} \mapsto$$

$$[\text{FC}(d_{\text{in}}, d_{\text{out}}/2)(\mathbf{X}), \text{1MaxPool}(\text{FC}(d_{\text{in}}, d_{\text{out}}/2)(\mathbf{X}))],$$

and the encoder ϕ consists of the following blocks:

$$\text{FC}(3, 510) \rightarrow \text{EnBlock}(510, 510) \rightarrow$$

$$\text{EnBlock}(510, 510) \rightarrow \text{EnBlock}(510, 510) \rightarrow$$

$$\text{FC}(510, 255) \rightarrow \text{MaxPool} \xrightarrow{\times 2} \text{Linear}(255, 255),$$

where the final two linear layers outputs vectors $(\boldsymbol{\mu}, \boldsymbol{\eta})$, with $\boldsymbol{\mu} \in \mathbb{R}^{255}$ is the latent mean, and $\boldsymbol{\eta} \in \mathbb{R}^{255}$ is the latent log-standard-deviation.

The decoder, ψ , similarly to [1, 10], is a composition of 8 layers where the first layer is $\text{FC}(255 + 3, 512)$, the second and third layer are $\text{FC}(512, 512)$, the fourth layer is a concatenation of $\text{FC}(512, 255)$ to the input, the fifth to seventh layers are $\text{FC}(512, 512)$, and the final layer is $\text{Linear}(512, 1)$. For the Frame Averaging architecture of Φ and Ψ , we use the same as the above backbone encoder ϕ and decoder ψ . We set $m = 0$ and $d = 85$. For the frames construction of Φ and Ψ , $\mathcal{F} : V \rightarrow 2^{E(3)} \setminus \emptyset$, we use the construction defined in section 3.3, with $\mathbf{w} = \mathbf{1}$. Note that for Φ the frames are calculated with respect to the input point cloud $\mathbf{V} = \mathbf{X} \in \mathbb{R}^{6400 \times 3}$, and for Ψ the frames are calculated with respect to $\mathbf{V} = \mathbf{Z} \in \mathbb{R}^{85 \times 3}$.

Implementation details. For all methods in table 2, training was done using the ADAM [6] optimizer, with batch size 64. The number of epochs was set to 4000, with 0.0001 fixed learning rate. Training was done on 4 Nvidia QUADRO RTX 8000 GPUs.

Evaluation. We used the following Chamfer distance metrics to measure similarity between shapes:

$$d_C(\mathcal{X}_1, \mathcal{X}_2) = \frac{1}{2} (d_C^{\rightarrow}(\mathcal{X}_1, \mathcal{X}_2) + d_C^{\rightarrow}(\mathcal{X}_2, \mathcal{X}_1)) \quad (2)$$

where

$$d_C^{\rightarrow}(\mathcal{X}_1, \mathcal{X}_2) = \frac{1}{|\mathcal{X}_1|} \sum_{\mathbf{x}_1 \in \mathcal{X}_1} \min_{\mathbf{x}_2 \in \mathcal{X}_2} \|\mathbf{x}_1 - \mathbf{x}_2\|^2 \quad (3)$$

and the sets \mathcal{X}_i are point clouds. As the input test set are given as 3D point clouds, we use the test point clouds directly with this metric. However, as the reconstructions are represented as implicit surfaces, we first mesh the predictions zero levelset using the MarchingCubes [8] algorithm, and second, we generate a point cloud of 30000 points by uniformly sample from the predicted mesh.

A.3. Piecewise Euclidean: Mesh \rightarrow mesh

Architecture. Here we provide additional details regarding the architectures for the methods in table 3 applied on the DFaust splits. The AE architecture is the same as the one detailed in section A.1 for the global Euclidean mesh \rightarrow mesh. For ARAPReg [5], random split, we used the trained models provided with the official implementation. Note that ARAPReg is an auto-decoder, where in test time, latent codes are optimized to reconstruct test inputs. For the Frame Averaging version with the above encoder and decoder backbones, we describe next the design choices that have been made. We set the weights $\mathbf{W} \in \mathbb{R}^{6890 \times 24}$ according to the template skinning weights provided with SMPL [7]. For the random split, the latent space dimensions were set to $k = 24$, $m = 12$, $d = 20$. For both the unseen poses splits the latent space dimensions were set to $k = 24$, $m = 72$, $d = 24$. For the frames of Φ , $\mathcal{F}_i : V \rightarrow 2^{E(3)} \setminus \emptyset$, with $V = \mathbb{R}^{6890 \times 3}$ and $1 \leq i \leq 24$, we use the construction defined in section 3.3. We set the weights \mathbf{w} used for the weighted PCA of the i -th frame with $\mathbf{W}e_i$, where \mathbf{W} is the template skinning weights and e_i is the i -th element in the standard basis of \mathbb{R}^{24} . The frames of $\langle \psi \rangle_{\mathcal{F}}$, $\mathcal{F} : V \rightarrow 2^{E(3)} \setminus \emptyset$ (where $V = \mathbb{R}^{20 \times 3}$), are calculated with $\mathbf{w} = \mathbf{1}$. Note that exactly the same Frame Averaging architecture was also used in the experiment in table 4.

For the experiment with SMAL in table 3, we next provide the baseline AE architecture based on the one provided in ARAPReg. The encoder, ϕ , is consisted of the following blocks, where each block is defined as in section A.1:

$$\text{EnBlock}(3889, 1945, 3, 32) \rightarrow$$

$$\text{EnBlock}(1945, 973, 32, 32) \rightarrow$$

$$\text{EnBlock}(973, 973, 32, 32) \rightarrow$$

$$\text{EnBlock}(973, 973, 32, 64) \rightarrow \text{Linear}(62272, 96) \rightarrow .$$

The decoder, ψ , consists of the following blocks:

Linear(96, 62272) \rightarrow DeBlock(973, 973, 64, 64) \rightarrow
 DeBlock(973, 973, 64, 32) \rightarrow
 DeBlock(973, 1945, 32, 32) \rightarrow
 DeBlock(1945, 3889, 32, 32) \rightarrow
 Conv(3889, 3889, 32, 3) \rightarrow .

For our Frame Averaging architecture we set $k = 33$, $m = 12$, $d = 28$. The weights $\mathbf{W} \in \mathbb{R}^{3889 \times 33}$ are set according to the template skinning weights provided with the dataset. The frames for the encoder and the decoder are constructed in the same fashion as the one described above for the Frame Averaging DFaust architecture.

Implementation details. For our Frame Averaging method in tables 3 and 4, training was done using the ADAM [6] optimizer, with batch size 16. The number of epochs was set to 100, with 0.0001 fixed learning rate. Training was done on 4 Nvidia QUADRO RTX 8000 GPUs. For ARAPReg, the numbers reported for the random split and SMAL in table 3 are from the original paper [5]. For the unseen pose and unseen global pose splits in table 3, we retrained the ARAPReg and AE using the official implementation training details, where the only change is that the latent size was set to 144 from 72 (to be equal to ours). The numbers reported in table 4 for the baselines SNARF [2] and NASA [4] are from table 1 in [2]. The qualitative results for SNARF in figure 5 were produced using SNARF official implementation. Note that the meshes for the qualitative results in the original SNARF paper were extracted using a *Multiresolution IsoSurface Extraction* [9], whereas the meshing in the official SNARF implementation (used in this paper) is based on the standard MarchingCubes [8] algorithm, result in with more severe staircase artifacts.

Evaluation. The definition for the MSE metric reported in 3 is the same as defined above in equation 1. The numbers reported for our method in table 3, were produced with the evaluation code from the official implementation of SNARF. For completeness, we repeat the procedure for the calculation of the IoU metric here: Let \mathbf{X} be a ground truth shape and \mathbf{Y} is the corresponding prediction. Then, two samples of points in space, bbox (bounding box) and near surface, are constructed. The near surface sample is constructed by adding a Gaussian noise with $\sigma = 0.01$ to a sample of 200000 points from the shape \mathbf{X} . The bbox is constructed by uniformly sampling of 100000 points in a the bounding box of all shapes in the dataset. Then, the IoU is computed by

$$\text{IoU} = \frac{1}{|S|} \sum_{\mathbf{x} \in S} o(\mathbf{x}; \mathbf{X}) * o(\mathbf{x}; \mathbf{Y})$$

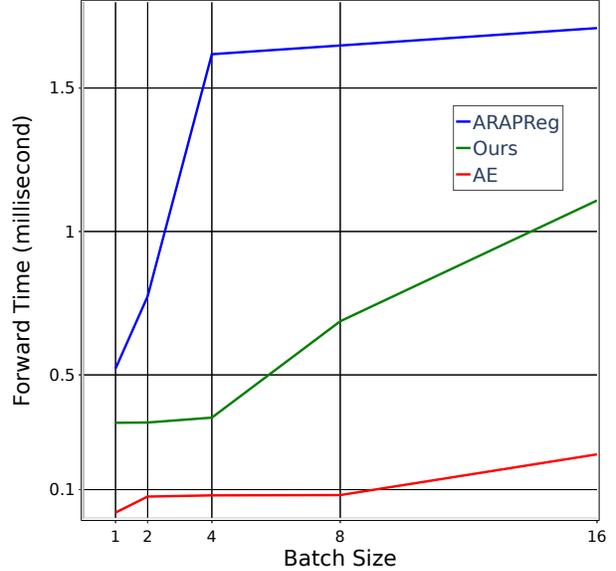


Figure A1. Timings for the methods in table 3 for various batch sizes.

where S is the prepared sample (bbox or near surface), and $o(\cdot; \mathbf{X}) \in \{0, 1\}$, $o(\cdot; \mathbf{Y}) \in \{0, 1\}$ are the occupancy functions of the shapes \mathbf{X} and \mathbf{Y} .

A.4. Interpolation in shape space

Here we provide additional details regarding the interpolation in shape space. Let $\mathbf{Z}^{(j)} = (\mathbf{q}^{(j)}, \mathbf{Q}^{(j)}) \in Z$, $j = 0, 1$, be two latent codes. Here we describe the case $k = 1$, as for $k > 1$ the following scheme can be applied in the same manner for each part. Note that $\mathbf{q}^{(j)}$ is the invariant features, whereas $\mathbf{Q}^{(j)}$ is the equivariant features. Let \mathbf{R} be the optimal rotation between $\hat{\mathbf{Q}}^{(0)}$ to $\hat{\mathbf{Q}}^{(1)}$, where $\hat{\mathbf{Q}}^{(j)}$ denotes the centered version of $\mathbf{Q}^{(j)}$. Then,

$$\mathbf{Q}_t = \text{slerp}(t, \mathbf{I}, \mathbf{R})(\hat{\mathbf{Q}}^{(0)} + t\mathbf{D}) + t \frac{1}{d} \mathbf{1}\mathbf{1}^T \mathbf{Q}^{(0)} + (1-t) \frac{1}{d} \mathbf{1}\mathbf{1}^T \mathbf{Q}^{(1)}$$

where $\mathbf{D} = \mathbf{R}^T \mathbf{Q}^{(1)} - \mathbf{Q}^{(0)}$, and $\text{slerp}(t, \mathbf{I}, \mathbf{R})$ denotes the spherical linear interpolation between \mathbf{I} to \mathbf{R} . For the invariant features, we perform regular linear interpolation $\mathbf{q}_t = t\mathbf{q}^{(0)} + (1-t)\mathbf{q}^{(1)}$. Then, the interpolant at time t is

$$\mathbf{Z}_t = (\mathbf{q}_t, \mathbf{Q}_t).$$

A.5. Timings

In figure A1 we report the training forward time for various batch sizes for the methods in table 3. Note that the increase in forward time for our method with respect to the baseline AE, stems from the computation of the backbone network for all the different parts and possible elements in

the frame. Note that for ARAPReg, the forward time in test reduces significantly as the computation of the ARAP regularizer is not needed. However, note that in the case of ARAReg with an AutoDecoder, prediction time increases significantly as it requires optimizing the latent codes.

B. Proofs

Proof of Proposition 1. The proof basically repeats the proof in [11] but for the weighted case. Let $g = (\mathbf{S}, \mathbf{u})$ be an element in $E(3)$. We need to show that $\mathcal{F}(\rho_1(g)\mathbf{V}) = g\mathcal{F}(\mathbf{V})$. The above equality for the translations part \mathbb{R}^3 in the group $E(3) = O(3) \times \mathbb{R}^3$ is trivial. For showing the above equality for $O(3)$, Let \mathbf{C}_V be the covariance matrix of \mathbf{V} . That is $\mathbf{C}_V = (\mathbf{V} - \frac{1\mathbf{w}^T}{1^T\mathbf{w}}\mathbf{V})^T \text{diag}(\mathbf{w})(\mathbf{V} - \frac{1\mathbf{w}^T}{1^T\mathbf{w}}\mathbf{V})$. A direct calculation shows that

$$\mathbf{C}_V = \mathbf{V}^T \left[\text{diag}(\mathbf{w}) \left(\mathbf{I} - \frac{1\mathbf{w}^T}{1^T\mathbf{w}} \right) \right] \mathbf{V}.$$

Then, the covariance matrix of $\rho_1(g)\mathbf{V}$ satisfies

$$\mathbf{C}_{\rho_1(g)\mathbf{V}} = \mathbf{S}^T \mathbf{V}^T \left[\text{diag}(\mathbf{w}) \left(\mathbf{I} - \frac{1\mathbf{w}^T}{1^T\mathbf{w}} \right) \right] \mathbf{V} \mathbf{S}.$$

Thus, \mathbf{r} is an eigen vector of \mathbf{C}_V if and only if $\mathbf{S}\mathbf{r}$ is an eigen vector of $\mathbf{C}_{\rho_1(g)\mathbf{V}}$. To finish the proof, notice that by definition $g\mathcal{F}(\mathbf{V}) = \{(\mathbf{S}\mathbf{R}, \mathbf{S}\mathbf{t} + \mathbf{u}) \mid \mathbf{R} = [\pm r_1, \pm r_2, \pm r_3]\}$. \square

Proof of Proposition 2. We need to show that

$$\Psi(\rho_Z(g)\mathbf{Z}) = \rho_Y(g)\Psi(\mathbf{Z})$$

which is equivalent to showing that for all $\mathbf{x} \in \mathbb{R}^3$

$$\hat{\Psi}(\rho_Z(g)\mathbf{Z}, \mathbf{x}) = \hat{\Psi}(\mathbf{Z}, \mathbf{R}^T(\mathbf{x} - \mathbf{t}))$$

which is again equivalent to showing that for all $\mathbf{x} \in \mathbb{R}^3$

$$\hat{\Psi}(\rho_Z(g)\mathbf{Z}, \mathbf{R}\mathbf{x} + \mathbf{t}) = \hat{\Psi}(\mathbf{Z}, \mathbf{x})$$

which by definition of ρ_V is finally equivalent to showing for all $\mathbf{x} \in \mathbb{R}^3$

$$\hat{\Psi}(\rho_P(g)(\mathbf{Z}, \mathbf{x})) = \rho_{\mathbb{R}}(g)\hat{\Psi}(\mathbf{Z}, \mathbf{x})$$

as required. \square

Proof of Theorem 1. If $g \in \mathcal{F}_j(\mathbf{X})$ then $g = (\mathbf{R}, \mathbf{t})$, where $\mathbf{t} = (\mathbf{1}^T \mathbf{w}_j)^{-1} \mathbf{X}^T \mathbf{w}_j$ as defined in equation 7. Therefore

$$\begin{aligned} \rho_X(g)^{-1} \mathbf{X}_j &= ((\mathbf{1} - \mathbf{w}_j) \mathbf{t}^T + \mathbf{w}_j \odot \mathbf{X} - \mathbf{1} \mathbf{t}^T) \mathbf{R} \\ &= \mathbf{w}_j \odot (\mathbf{X} - \mathbf{1} \mathbf{t}^T) \mathbf{R} \\ &= \mathbf{w}_j \odot \rho_X(g)^{-1} \mathbf{X} \end{aligned}$$

Furthermore, it can be directly checked that for hard weights $\mathcal{F}_j(\mathbf{X})$ is an equivariant frame that is defined only in terms of the rows of \mathbf{X} that belong to the j -th part. Now, $\mathbf{Z}_j = \langle \phi \rangle_{\mathcal{F}_j}(\mathbf{X}_j)$ where

$$\begin{aligned} \langle \phi \rangle_{\mathcal{F}_j}(\mathbf{X}_j) &= \frac{1}{|\mathcal{F}_j(\mathbf{X})|} \sum_{g \in \mathcal{F}_j(\mathbf{X})} \rho_{Z_j}(g) \phi(\rho_X(g)^{-1} \mathbf{X}_j) \\ &= \frac{1}{|\mathcal{F}_j(\mathbf{X})|} \sum_{g \in \mathcal{F}_j(\mathbf{X})} \rho_{Z_j}(g) \phi(\mathbf{w}_j \odot \rho_X(g)^{-1} \mathbf{X}). \end{aligned}$$

Therefore Frame Averaging now implies that $\langle \mathbf{X}_j \rangle_{\mathcal{F}}$ is equivariant to Euclidean motions of the j -th part. For the decoder, $\langle \mathbf{Z}_j \rangle_{\mathcal{F}}$ is equivariant to Euclidean motions of \mathbf{Z}_j by Frame Averaging, and $\mathbf{w}_j \odot \rho_Y(g_j) \mathbf{Y}$ transforms the decoded j -th part accordingly. \square

C. Additional results

We provide additional qualitative results for some of the experiments in section 5. Figure A2 shows additional qualitative results for the comparison with an implicit pose-conditioned method (SNARF). Figure A3 shows typical test reconstructions on the SMAL dataset. Note that the test set here does not include out of distribution poses, hence relatively easy. Note that our method is able to capture the pose better than the baselines (see the animal's tail for example), while achieving high fidelity (see base of the feet).

References

- [1] Matan Atzmon and Yaron Lipman. SALD: sign agnostic learning with derivatives. In *9th International Conference on Learning Representations, ICLR 2021*, 2021.
- [2] Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. *arXiv preprint arXiv:2104.03953*, 2021.
- [3] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29:3844–3852, 2016.
- [4] Boyang Deng, John P Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. Nasa neural articulated shape approximation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VII 16*, pages 612–628. Springer, 2020.
- [5] Qixing Huang, Xiangru Huang, Bo Sun, Zaiwei Zhang, Junfeng Jiang, and Chandrajit Bajaj. Arapreg: An as-rigid-as possible regularization loss for learning deformable shape generators. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5815–5825, 2021.
- [6] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [7] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015.
- [8] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *ACM siggraph computer graphics*, volume 21, pages 163–169. ACM, 1987.
- [9] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019.
- [10] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019.
- [11] Omri Puny, Matan Atzmon, Heli Ben-Hamu, Edward J Smith, Ishan Misra, Aditya Grover, and Yaron Lipman. Frame averaging for invariant and equivariant network design. *arXiv preprint arXiv:2110.03336*, 2021.
- [12] Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J Black. Generating 3d faces using convolutional mesh autoencoders. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 704–720, 2018.
- [13] Silvia Zuffi, Angjoo Kanazawa, David W Jacobs, and Michael J Black. 3d menagerie: Modeling the 3d shape and pose of animals. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6365–6373, 2017.

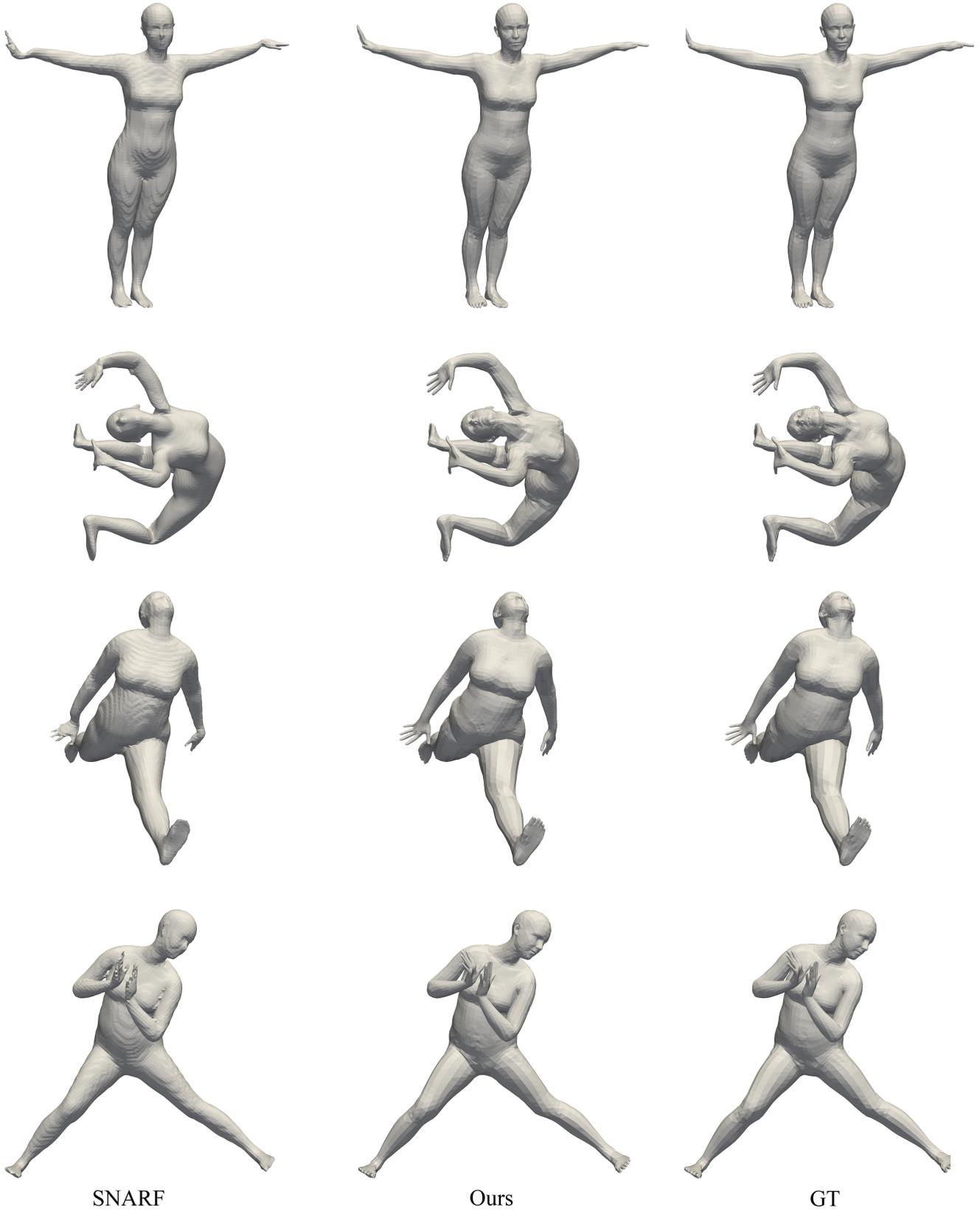
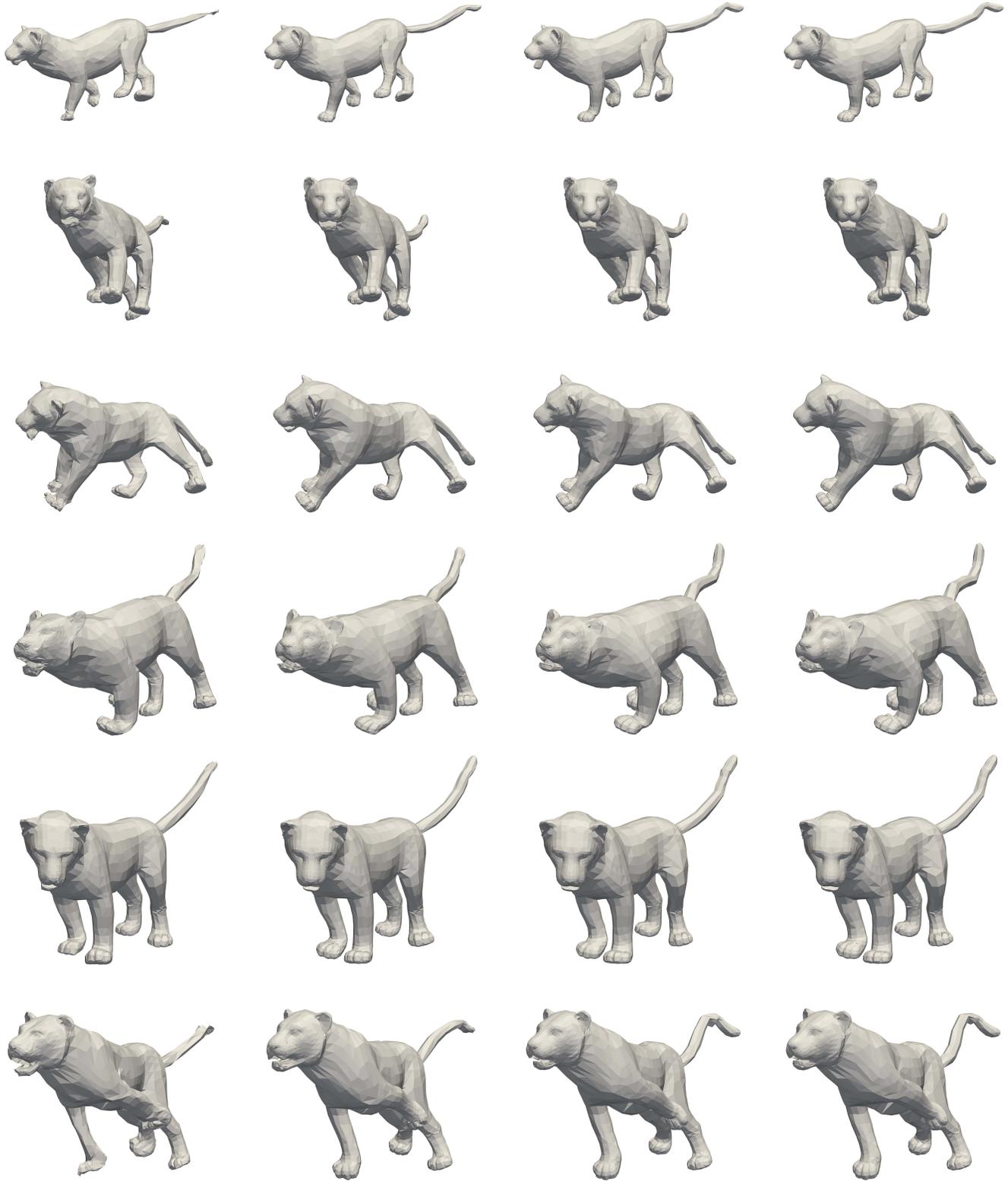


Figure A2. Comparison to SNARF on "out of distribution" test.



AE

ArapReg

Ours

GT

Figure A3. Piecewise Euclidean mesh \rightarrow mesh, qualitative results; SMAL [13] dataset.