

Overview

 \succ We propose to use meta-learning to effectively incorporate cloth-deformation priors for clothed humans, thus enabling fast fine-tuning (few minutes) for generating neural avatars given only a few monocular depth images of unseen clothed humans and their corresponding SMPL fittings as inputs.



Our Approach

> We follow the same pipeline of recently proposed SCANimate [1], which learns dynamic neural SDFs from dense full-body scans to represent subject/cloth-type specific avatars of clothed humans. Instead of learning subject/cloth-type specific models from scratch, we propose to meta-learn a prior model which can be fast fine-tuned to represent any subject/cloth-type specific neural avatars given only monocular depth frames.



MetaAvatar: Learning Animatable Clothed Human Models from Few Depth Images Shaofei Wang¹, Marko Mihajlovic¹, Qianli Ma^{1,2}, Andreas Geiger^{2,3}, Siyu Tang¹ ¹ETH Zürich, ²Max Planck Institute for Intelligent Systems, Tübingen, ³University of Tübingen



[1] Saito et al, SCANimate: Weakly Supervised Learning of Skinned Clothed Avatar Networks. CVPR, 2021.

Meta-learned Prior

- \succ Our key contribution is a meta-learned hypernetwork.
- \succ In practice, we decompose the training into two stages 1) meta-learn a static neural SDF 2) meta-learn a hypernetwork which predicts residuals to the parameters of the previously meta-learned static neural
 - SDF.

Training Stage 1: Learn Meta-SDF

Linear

SDF $f_{\phi^*}(\mathbf{x})$



 \succ Results show that our approach performs well with limited depth input data.

 $\rightarrow \longrightarrow \rightarrow$ Linear $\rightarrow \longrightarrow \rightarrow$ Linear

| | | 3D Input | | | 2.5D Input | - | $1 \rightarrow 1 \rightarrow 1$ | 100 | 50 | 20 | 10 | 5 | 1 | |
|---|------------------|----------|-------|-----------|-------------------|---|---|-------|-------|---------------|-------|-------|---------|--|
| | | NASA | LEAP | SCANimate | Ours | - Fine-t | une data (%) | 100 | 50 | 20 | 10 | 3 | <1 | |
| Subi 00122_00215 | | | | | | | Subj 00122, 00215 | | | | | | | |
| $\frac{500 \int 00122,00215}{E_{\rm T}}$ | | | | 0.5 | - <u>Ex.</u> | PS ↑ | 0.5 | 0.471 | 0.509 | 0.473 | 0.373 | 0.510 | | |
| EX. | P5 | 0.078 | 0.314 | 0.333 | 0.5 | _ | D_{m} | _ | 0.450 | 0 4 8 0 | 0.512 | 0 543 | 0 592 | |
| | $D_p \downarrow$ | 0.484 | 0.454 | 0.586 | 0.450 | Int | $D p \downarrow$ | | 0.120 | 0.100 | 0.252 | 0.201 | 0.372 | |
| Int. | $\hat{D_f}$ | 0.327 | 0.293 | 0.489 | 0.273 Int. | $D_f \downarrow$ | - | 0.275 | 0.510 | 0.555 | 0.391 | 0.430 | | |
| | $\frac{-}{NC}$ | 0 752 | 0.807 | 0 793 | 0 821 | | $NC\uparrow$ | - | 0.821 | 0.808 | 0.795 | 0.785 | 0.768 | |
| $\frac{110 \ \ 0.752 \ 0.007 \ 0.753 \ 0.021}{0.021}$ | | | | | | _ | Subj 00134, 03375 | | | | | | | |
| Subj 00134, 03375 | | | | | | - F v | DC ተ | 0.5 | 0 176 | $\frac{1}{0}$ | 0.463 | 0/30 | 0 3 8 7 | |
| Ex. | PS ↑ | 0.182 | 0.224 | 0.481 | 0.5 | <u> </u> | | 0.5 | | | | | 0.307 | |
| Int. | $D_{\rm m}$ | 0 595 | 0.483 | 0.629 | 0 518 | - | $D_p\downarrow$ | - | 0.518 | 0.545 | 0.5/6 | 0.603 | 0.619 | |
| | $D p \downarrow$ | 0.575 | 0.100 | 0.542 | 0.267 | Int. | $D_f \downarrow$ | - | 0.367 | 0.400 | 0.438 | 0.471 | 0.489 | |
| | $D_f \downarrow$ | 0.409 | 0.340 | 0.342 | 0.307 | | $NC\uparrow$ | _ | 0.773 | 0.762 | 0.753 | 0.745 | 0.737 | |
| | $NC\uparrow$ | 0.693 | 0.780 | 0.755 | 0.773 | | Average per model training/fine tuning time (hours) | | | | | | | |
| Averge per-model training/fine-tuning time (hours) | | | | | - | Average per-model training/me-tuning time (nours) | | | | | | | | |
| | 0 | >10 | | <u> </u> | 1.60 | - | | 1.60 | 0.8 | 0.32 | 0.16 | 0.08 | 0.02 | |
| | | /10 | - | ~10 | 1.00 | _ | | | | | | | | |

Comparison to baslines











Few-shot learning