# HiLo: Detailed and Robust 3D Clothed Human Reconstruction with High-and Low-Frequency Information of Parametric Models

Yifan Yang<sup>1,2\*</sup> Dong Liu<sup>1\*</sup> Shuhai Zhang<sup>1,2</sup> Zeshuai Deng<sup>1</sup> Zixiong Huang<sup>1</sup> Mingkui Tan<sup>1,2,3†</sup> <sup>1</sup>South China University of Technology <sup>2</sup>Pazhou Lab

<sup>3</sup>Key Laboratory of Big Data and Intelligent Robot, Ministry of Education

{seyoungyif, sesmildong, mszhangshuhai, sedengzeshuai, sesmilhzx}@mail.scut.edu.cn mingkuitan@scut.edu.cn

## Abstract

Reconstructing 3D clothed human involves creating a detailed geometry of individuals in clothing, with applications ranging from virtual try-on, movies, to games. To enable practical and widespread applications, recent advances propose to generate a clothed human from an RGB image. However, they struggle to reconstruct detailed and robust avatars simultaneously. We empirically find that the high-frequency (HF) and low-frequency (LF) information from a parametric model has the potential to enhance geometry details and improve robustness to noise, respectively. Based on this, we propose HiLo, namely clothed human reconstruction with high- and low-frequency information, which contains two components. 1) To recover detailed geometry using HF information, we propose a progressive HF Signed Distance Function to enhance the detailed 3D geometry of a clothed human. We analyze that our progressive learning manner alleviates large gradients that hinder model convergence. 2) To achieve robust reconstruction against inaccurate estimation of the parametric model by using LF information, we propose a spatial interaction implicit function. This function effectively exploits the complementary spatial information from a low-resolution voxel grid of the parametric model. Experimental results demonstrate that HiLo outperforms the state-of-the-art methods by 10.43% and 9.54% in terms of Chamfer distance on the Thuman2.0 and CAPE datasets, respectively. Additionally, HiLo demonstrates robustness to noise from the parametric model, challenging poses, and various clothing styles.<sup>1</sup>

# 1. Introduction

The creation of 3D realistic digital human plays a pivotal role in the realm of mixed reality [24, 36, 57], remote pre-



Figure 1. Visualization comparisons on in-the-wild images, our HiLo achieves more accurate and detailed reconstruction on challenging poses and diverse clothes.

sentation [5, 57], film [16, 49], and gaming [51]. Traditional methods often require expensive and specialized equipment combined with complex artistic efforts to customize the avatars[17, 46, 62], which limits the ability of individuals to create personalized avatars easily. To address the limitation, recent approaches [1, 2, 27, 30, 34, 44, 45, 53–55, 64] capture a 3D avatar from an RGB image of a clothed human, thus eliminating the need for costly scanning equip-

Authorized licensed use limited to: SOUTH CHINA UNIVERSITY OF TECHNOLOGY. Downloaded on April 25,2025 at 09:42:26 UTC from IEEE Xplore. Restrictions apply.

<sup>&</sup>lt;sup>†</sup>Corresponding author, \*Equal contribution.

<sup>&</sup>lt;sup>1</sup>Code link: https://github.com/YifYang993/HiLo.git

ment and making it easier for a broader range of users to create personalized avatars.

Despite the convenience of recent advances, the input image usually lacks details about delicate human body parts and diverse clothes from multiple angles. Moreover, the limited viewpoints and lack of accurate depth information make the reconstruction vulnerable to noise, *e.g.*, inaccurate shape and pose of the estimated parametric body model [15, 58, 59]. Therefore, a detailed and robust 3D human reconstruction from an RGB image is challenging.

In spite of the impressive results of the previous methods, they have not fully addressed the problem of detailed and robust reconstruction simultaneously. Specifically, PIFu [44] produces overly smoothed or non-human body shapes on the unseen side of the human from the input image. ECON [54] requires Poisson Surface Reconstruction and replacement of body parts, introducing an extent of computational overhead (c.f. Sec. 4.4). Additionally, there is a risk of body part misalignment when the mid-term data is inaccurate. Considering that clothes need to conform to the surface of naked bodies, the geometry of the parametric model provides effective semantic regularization for reconstructing clothed humans. PaMIR [64], ICON [53], and D-IF [55] use parametric human bodies [28, 41] to regularize the reconstruction. However, the performance of these methods degrades significantly when facing noise on the parameters from the estimated naked bodies (c.f. Sec. 4.4).

To achieve robust 3D clothed human reconstruction with *detailed* geometry, we aim to explore how to further use the regularization from the parametric model to facilitate this goal. Our exploration is based on two common observations. First, high-frequency (HF) information enhances details [33, 43]. Considering that Signed Distance Function (SDF) [37] describes the geometry of a parametric model by representing a distance to the object surface boundary, we investigate the effectiveness of SDF in improving the geometry details of clothed humans. Second, low-frequency (LF) information is relatively robust to noise [10, 11, 26, 61]. Since inaccurate parametric model estimation within an error range has an insignificant impact on the corresponding low-resolution voxel grid [48], we seek to use the voxel grid to mitigate the noise of the estimated body. As shown in Fig. 2, qualitative results demonstrate that SDF boosts the reconstruction details while the voxel grid improves robustness to noise. However, how to effectively combine highand low-frequency information to generate details and mitigate noise simultaneously is still an open question.

In this paper, we propose a high- and low-frequency paradigm *HiLo*, which stands for *clothed human reconstruction with high- and low-frequency information*. To achieve HF detail, we further enhance SDF with HF function [43]. Intuitively, by amplifying the variation of adjacent points that share similar SDFs, we allow for better



Figure 2. We empirically demonstrate the effectiveness of the high-frequency (HF) regularization from naked bodies in enhancing geometry details in Toy Experiment <sup>(6)</sup>. We also verify the effectiveness of the low-frequency (LF) regularization in improving robustness to noise in Toy Experiment <sup>(6)</sup>.

delineation and capturing of fine details in the 3D human. Moreover, to alleviate the convergence difficulty caused by the large gradients amplified by HF function (c.f. Sec. 3), we introduce a progressive HF SDF that learns detailed 3D geometry in a coarse-to-fine manner. To achieve robustness, we seek to capture the LF complementary information of the low-resolution voxel grid from the naked human body. To this end, we design a spatial interaction implicit function, which promotes the interaction of global and local information across different voxels via an attention mechanism.

We qualitatively and quantitatively evaluate our HiLo on in-the-wild images and benchmark datasets. The experimental results verify the superiority of HiLo over previous approaches in three key aspects: 1) 3D geometry details (see Fig. 1). 2) Robust reconstruction. 3) Convergence speed. We summarize our contributions in three folds:

- To enhance the geometry details and improve robustness against noise during the clothed human reconstruction process, we propose to explore the high-frequency (HF) information and low-frequency (LF) information from a parametric body model simultaneously.
- To facilitate detailed reconstruction, we introduce a progressive HF function to enhance the signed distance function (SDF) of a parametric model, providing regularization during the reconstruction process. This function learns an HF SDF in a progressive manner to alleviate the convergence difficulty associated with HF informa-

tion. Experimental results show that HiLo reconstructs a more detailed clothed human.

 To ensure robust reconstruction, we employ LF information of low-resolution voxel grids from the parametric model to regularize the reconstruction. We propose a spatial interaction implicit function that reasons complementary information between different voxels. Experimental results show that HiLo is robust to various levels of noise.

# 2. Preliminaries

# 2.1. Signed Distance Function

Signed distance function (SDF) [37] is a continuous function that takes a given spatial point p with spatial coordinate  $x \in \mathbb{R}^n$  and outputs the distance  $s \in \mathbb{R}$  of the point to the closest point on the surface  $\partial\Omega$  of an object  $\Omega$ :

$$\mathcal{F}_{s}(\mathbf{p}) = s, \ s = \begin{cases} d(x,\partial\Omega) & \text{if } x \notin \Omega, \\ -d(x,\partial\Omega) & \text{if } x \in \Omega, \end{cases}$$
(1)

where  $d(x, \partial \Omega) = \inf_{y \in \partial \Omega} (||x - y||_2)$ . The sign of the distance implies whether the point is inside (negative) or outside (positive) of the surface  $\partial \Omega$ . s = 0 denotes the point **p** locates on the  $\partial \Omega$ .

# 2.2. Voxel Grid and Mesh Voxelization

**Voxel Grids**  $\Omega_v \in \mathbb{R}^{d \times h \times w}$  is a representative data structure for describing a 3D object  $\Omega$ . Specifically,  $\Omega_v$  is a three-dimensional matrix of 3D space with depth d, height h and width w.  $\Omega_v$  is composed of equally distributed and equally sized cube-shaped volumetric elements called voxels. The term voxel is the 3D counterpart to a 2D pixel. The resolution of a voxel grid is determined by the size of the voxels and the dimensions of the grid. Lower resolution implies larger voxels, resulting in a coarser representation of the space. Thus, we use the low-resolution voxel grid to represent the low-frequency information of a 3D object.

Mesh voxelization  $\mathcal{V}$  is a computational technique that plays a crucial role in converting irregular continuous 3D geometric models [4, 28, 41] such as triangular mesh and point clouds into regular and discrete voxel grids. In this process, a 3D mesh  $\Omega_m$ , which is a collection of connected triangles, is converted into a grid of voxels  $\Omega_v \in \mathbb{R}^3$ . In this paper, we use  $\mathcal{V}$  to transfer SMPL-X mesh  $\mathcal{M}$  to a lowresolution voxelized mesh  $\mathcal{M}_v \in \mathbb{R}^{32 \times 32 \times 32}$ , and then use a 3D CNN to encode  $\mathcal{M}_v$  for  $\mathcal{M}_v^{3D} \in$  for more flexibility. Based on our observation from Tab. 3, the low-frequency natural of  $\mathcal{M}_v^{3D}$  aids multiple existing methods in mitigating various noise levels in SMPL-X shape and pose.

# 3. Clothed Human Reconstruction with Highand Low-frequency Information

We aim to robustly infer detailed 3D clothed avatars from RGB images  $\mathcal{I}$ . Recent advances [8, 9, 19, 21] tend to use parametric naked body  $\mathcal{M}$  such as SMPL [28] or SMPL-X [41] estimated from  $\mathcal{I}$  to provide semantic regularization on clothed human avatars. We empirically verify that high-frequency (HF) and low-frequency (LF) information from  $\mathcal{M}$  are able to refine geometry and improve robustness to the reconstruction of the clothed human (Sec. 4.1). Based on this, we propose clothed human reconstruction with high- and low-frequency information, namely HiLo, which balances the HF and LF information to achieve detailed and robust reconstruction simultaneously. As shown in Fig. 3, HiLo contains two key components: (1) To refine the geometry of clothed human with HF information, we propose to use progressive high-frequency function to enhance the signed distance function (SDF) of  $\mathcal{M}$  (c.f. Sec. 2), and alleviate the convergence difficulty introduced by large gradients in a coarse-to-fine enhancement manner. (2) To achieve robust reconstruction using low-frequency information, we explore complementary information from lowresolution voxels  $\mathcal{M}_v$  from  $\mathcal{M}$  for a more comprehensive understanding of human geometry. To this end, we design a spatial interactive implicit function (c.f. Sec. 3.2) that leverages spatial information from local and global voxelized SMPL-X to predict the occupancy field O. Finally, we use the Marching Cubes ([29]) to obtain the 3D mesh of the clothed avatar from  $\hat{\mathcal{O}}$ .

The overall optimization of our proposed HiLo minimizes the following objective function:

$$\mathcal{L}_{overall} = \mathcal{L}_a(\hat{\mathcal{O}}, \mathcal{O}), \tag{2}$$

where  $\mathcal{L}_a$  is the MSE loss and  $\mathcal{O}$  is a GT occupancy field.

#### **3.1. Progressive Growing High-Frequency SDF**

Given that SDFs can enhance 3D reconstruction quality as confirmed in Sec. 4.1, we will leverage this for more realistic avatar reconstruction. However, directly fitting input coordinates with SDFs may lead to subpar representation of *high-frequency* variation in geometry (see Sec. 4.3). This aligns with previous work [42] indicating neural networks prioritize learning *low-frequency* signals. We will explore effective strategies to mitigate this.

**Conventional high-frequency SDF.** To improve the ability to represent complicated 3D shapes robustly, a straightforward way is to apply periodic functions  $\mathcal{H}$  such as sine and cosine [47, 66] to extract high-frequency signals on SDF of each sampled point via

$$\mathcal{H}(s) = [s, \mathcal{H}_0(s), \mathcal{H}_1(s), \dots, \mathcal{H}_k(s), \dots, \mathcal{H}_L(s)],$$
  
$$\mathcal{H}_k(s) = [\sin(2^k \pi s), \cos(2^k \pi s)], k \in \{0, 1, \dots, L\}.$$
(3)



Figure 3. Overview of our proposed HiLo. Conditioned on a single-view image  $\mathcal{I}$  and the corresponding SMPL-X  $\mathcal{M}$ , we first prepare a signed distance field s and a low-resolution voxel grid  $\mathcal{M}_v^{3D}$  of the naked body. Then, our proposed progressive high-frequency signed distance function  $\mathcal{H}(s;\beta)$  enhances s for detailed geometry of the clothed human and alleviates convergence difficulties introduced by large gradients in a coarse-to-fine learning manner. Moreover, we design an implicit function  $\phi_{si}$  which leverages the complementary information of low-frequency voxels from  $\mathcal{M}_v^{3D}$  to mitigate various levels of noise. Finally, we combine the above HF and LF features to  $\phi_{si}$  to infer the occupancy field  $\hat{\mathcal{O}}$  of the clothed avatar.

In this way, we amplify the variation of adjacent points that share similar SDFs, allowing for better delineation and capturing of fine details in the 3D object.

Despite the positive characteristic of high-frequency SDF, effective updating for parameters is difficult. Specifically, the gradient of  $\mathcal{H}_k(s)$  w.r.t. s is calculated by

$$\frac{\partial \mathcal{H}_k(s)}{\partial s} = 2^k \pi [\cos\left(2^k \pi s\right), -\sin\left(2^k \pi s\right)]. \tag{4}$$

Eqn. (4) incorporated with the coefficient  $2^k \pi$  will significantly amplify the gradient signals regarding SDF, especially for larger k. Large gradients could lead to convergence difficulties and numerical instability, ultimately resulting in poor representation performance [3, 39].

**High-frequency SDF in a growing manner.** To address the above issue, we introduce a progressively growing approach as shown in Fig. 3 (2), initially emphasizing lowfrequency signal learning and gradually focusing on learning the high-frequency geometry. Specifically, in the early stage of training, we reduce the weight of high-frequency signals (*e.g.*,  $\mathcal{H}_k(s)$ ) which have higher k and progressively increase their importance during training. We formulate this schedule as  $\mathcal{H}_k(s; \beta)$  with a weight  $\omega_k(\beta)$  following [38]:

$$\mathcal{H}_k(s;\beta) = \omega_k(\beta) [\sin(2^k \pi s), \cos(2^k \pi s)],$$
  

$$\omega_k(\beta) = \begin{cases} 0 & \text{if } \beta - k < 0; \\ \frac{1 - \cos((\beta - k)\pi)}{2} & \text{if } 0 \le \beta - k < 1; \\ 1 & \text{if } \beta - k > 1, \end{cases}$$
(5)

where  $\beta$  is proportional to the iteration of the optimization process, see Fig. 4 for the relationship between  $\omega_k(\beta)$  and  $\beta$ . With  $\omega_k(\beta)$ , the gradient of  $\mathcal{H}_k(s;\beta)$  becomes

$$\frac{\partial \mathcal{H}_k(s;\beta)}{\partial s} = \omega_k(\beta) 2^k \pi [\cos\left(2^k \pi s\right), -\sin\left(2^k \pi s\right)].$$
(6)

Then, during the beginning of the optimization,  $\beta$  is set so small that only frequency components with a smaller value of k will be assigned a non-zero weight, while the frequency components with a higher value of k will be omitted. Throughout the optimization, the higher-frequency components are progressively activated. This manner allows HiLo to explore the low-frequency part and later focus on the fine-grained geometry of 3D humans.



Figure 4. Illustration of the relationship between progressive weights  $\omega$  and  $\beta$  during the training process.

#### **3.2.** Low-Frequency Information for Robustness

Most recent methods [53–55, 64] are based on SMPL-X. However, SMPL-X estimation often faces **misalignment is**-

sues with the corresponding image, especially when facing a challenging human pose. Thus, it is crucial to **achieve robust reconstruction** against misaligned SMPL-X. Our results (*c.f.* Fig. 2) show that the low-frequency information, represented by low-resolution voxel grids  $\mathcal{M}_v$  of SMPL-X  $\mathcal{M}$ , enhances robustness against noise. We leverage this insight to incorporate local and global information of  $\mathcal{M}_v$  for improved regularization in reconstruction.

Local voxels for 3D feature preparation. Motivated by that point-wise local 3D features from  $\mathcal{M}_v$  are robust to out-of-distribution pose and shape [64], we voxelize the estimated SMPL-X  $\mathcal{M}$  and query it by p. Specifically, to obtain the voxelization features, we convert the corresponding SMPL-X  $\mathcal{M}$  into a low-resolution voxel grid  $\mathcal{M}_{\mathcal{V}}$  by mesh voxelization operation  $\mathcal{V}$  [64] and encode it via a 3D CNN  $f_{3D}$  for a 3D feature volume  $\mathcal{M}_{\mathcal{V}}^{3D}$ . To obtain point-wise 3D features, we use trilinear interpolation to sample  $\mathcal{M}_v^{3D}$ based on coordinate p of sampled 3D points, resulting in  $\mathcal{M}_{\mathcal{V}}^{3D}(\mathbf{p})$ . We empirically find that by combining  $\mathcal{M}_{\mathcal{V}}^{3D}(\mathbf{p})$ , HiLo is robust to SMPL-X noise (c.f. Sec. 4.4). As shown in Fig. 3 (1), in addition to  $\mathcal{H}(s;\beta)$  and  $\mathcal{M}_{\mathcal{V}}^{3D}$ , we follow ICON [53] to use a normal feature  $\mathbf{F}_{n}(\mathbf{p})$  to provide detailed texture information. Then, we concatenate them into one final feature  $\mathbf{F}_{c}^{1}$ =[ $\mathcal{H}(s;\beta), \mathcal{M}_{\mathcal{V}}^{3D}(\mathbf{p}), \mathbf{F}_{n}(\mathbf{p})$ ] and then fed  $\mathbf{F}_{c}^{1}$  to our designed implicit function to reconstruct the clothed avatar.

Global voxels for spatial interaction implicit function. To reconstruct clothed avatars, a typical solution is to map 3D features  $\mathbf{F}$  to a continuous occupancy field that represents the interior and exterior of a clothed human. To this end, numerous literature [44, 53, 55, 64] uses an implicit function parameterized by a memory-efficient multi-layer perceptron (MLP)  $\mathcal{T}$  to map  $\mathbf{F}$  into an occupancy field  $\hat{\mathcal{O}}$ .

However, the potential issue of the existing implicit function lies in its underutilization of the global information inherent in 3D data. Previous research [6, 22, 35] has shown that the representation ability of features can be improved by capturing the global correlation between the features. For 3D clothed human reconstruction, different human body parts contain distinct yet complementary spatial information. For instance, as shown in Fig. 5 (a), the voxels located on the shoulder (the red point) may offer valuable topological cues to constrain the prediction of geometry near the elbow (the blue point).

To leverage global information from the voxel grid of SMPL-X of  $\mathcal{M}_v^{3D}$ , we design a spatial interaction module  $\mathcal{A}$  into  $\phi$  to infer the 3D occupancy, denoted by  $\phi_{si}$ , see detail in Appx B.2. As shown in Fig. 5 (b),  $\phi_{si}$  injects global-scale features of  $\mathcal{M}_v^{3D}$  to the local 3D feature with the aim of introducing whole-body awareness to  $\phi_{si}$ :

$$\phi_{si}(\mathbf{F}_c^1) \to \hat{\mathcal{O}}, \ \phi_{si}(\cdot) = \mathcal{A}^{N+1} \circ T^{(N+1)} \circ \cdots \circ \mathcal{A}^1(\cdot) \circ T^{(1)}.$$
(7)

**Optimization.** We optimize parameters of  $\phi_{si}$  and  $f_{3D}$ 



Figure 5. (a) Complementarity of voxel ① and ②. (b) Illustration of the spatial interaction module A.

via minimizing the MSE loss in Eqn. (2) between the predicted occupancy field  $\hat{O}$  and the ground-truth occupancy field O. With  $\hat{O}$ , we reconstruct the triangular mesh of the 3D clothed avatar via marching cubes algorithm [29].

#### 4. Experiments

**Datasets:** We conduct experiments on two open-source datasets, *i.e.*, Thuman2.0 [63] and CAPE [30], which both contain various human shapes with different human poses and diverse clothes. Specifically, following ICON [53], the CAPE dataset is divided into the "CAPE-FP" and "CAPE-NFP" sets, which have "fashion" and "non-fashion" poses, respectively, to further analyze the generalization to complex body poses.

Moreover, to evaluate our HiLo on in-the-wild images, we follow ICON to collect 200 diverse images from Pinterest<sup>2</sup>. The images contain humans performing dramatic movements or wearing diverse clothes.

**Metrics:** We evaluate our HiLo and baseline methods in terms of three metrics: **Chamfer distance** and **P2S distance** mainly measure coarse geometry error, while **Normals** mainly captures high-frequency differences. See details in Appx. C.2.

**Baselines:** We compare our proposed HiLo with mainstream state-of-the-art methods, including PIFu [44], PaMIR [64], ICON [53], ECON [54] and D-IF [55], refer to the Appx. C.4 for the detailed description. To demonstrate the necessity of naked 3D body regularization, we first conduct a toy experiment to study the effect of SDF on different baselines. To this end, we design three variants based on PIFu, PaMIR, and ICON, which incorporate SDF into existing methods (PIFu and PaMIR) or remove SDF from the existing method (ICON), namely PIFu<sub>w SDF</sub>, PaMIR<sub>w SDF</sub> and ICON<sub>w/o SDF</sub>, respectively.

For ablation studies, we construct several variants of our HiLo: 1) HiLo<sub> $w/o \phi_{si}$ </sub> replaces our spatial interaction implicit function with the vanilla implicit function; 2) HiLo<sub> $w/o M_v^{3D}$ </sub> is constructed by removing the voxelized SMPL-X; 3) HiLo<sub> $w/o H(s;\beta)$ </sub> is constructed by replacing our progressive HF SDF with vanilla SDF.

<sup>&</sup>lt;sup>2</sup>https://www.pinterest.com/



Figure 6. Visualization results of 3D clothed avatar reconstruction with our HiLo from in-the-wild images, which present various clothing and challenging poses. We show the front (blue) and rotated (red) views.

Table 1. Toy experiments about the impact of SDF on 3D clothed human reconstruction, on seen (*i.e.*, training and test on the same dataset) and unseen (*i.e.*, training on Thuman2.0 and test on CAPE) settings.

Dataset	CAPE-FP			CAPE-NFP			CAPE			Thuman2			CAPE		
Methods	Chamfer $(\downarrow)$	P2S $(\downarrow)$	Normals $(\downarrow)$	Chamfer	P2S	Normals	Chamfer	P2S	Normals	Chamfer	P2S	Normals	Chamfer	P2S	Normals
PIFu [44]	2.1000	2.0930	0.0910	2.9730	2.9400	0.1110	2.6820	2.6580	0.1040	2.6880	2.5730	0.0970	7.1244	2.7633	0.3902
$\mathrm{PIFu}_{w \ \mathrm{SDF}}$	0.8908	0.8637	0.0676	0.9848	0.9545	0.0698	0.9437	0.9178	0.0707	1.6659	1.7934	0.1360	1.3078	1.4306	0.0980
PaMIR [64]	1.2250	1.2060	0.0550	1.4130	1.3210	0.0630	1.3500	1.2830	0.0600	1.4388	1.5613	0.1071	0.9339	0.9444	0.0659
$\mathrm{PaMIR}_{\mathrm{w}~\mathrm{SDF}}$	0.9188	0.8788	0.0565	1.1132	1.0729	0.0611	1.0112	0.9725	0.0601	1.4073	1.5624	0.1174	0.8438	0.8179	0.0572
ICON [53]	0.7475	0.7488	0.0508	0.8656	0.8639	0.0545	0.8055	0.8084	0.0539	1.1431	1.3020	0.0923	0.8610	0.8878	0.0606
$\rm ICON_{w/o \ SDF}$	1.0243	0.9478	0.0741	1.4862	1.3313	0.0919	1.2736	1.1538	0.0850	1.3114	1.2116	0.1015	7.4892	1.7708	0.3780



Figure 7. Reconstruction results with or without our progressive high-frequency SDF  $\mathcal{H}(s;\beta)$ . The geometry details of clothes, hands, faces are enhanced by introducing  $\mathcal{H}(s;\beta)$ .

**Implementation details:** We implement our approach using PyTorch<sup>3</sup> [40] and train our networks with RM-Sprop [50] optimizer. For a fair comparison, we follow all common hyper-parameter settings same as ICON [53]. See more implementation details in the Appx.

### 4.1. Toy Experiments

Our motivation stems from the idea that HF information and LF information improve details and robustness, respectively. To verify this, we employ two tools, *i.e.*, *SDF* and *voxelized SMPL-X*  $\mathcal{M}_n^{3D}$  to establish this constraint.

**Impact of SDF.** We build upon three representative methods, *i.e.*, PIFu, PaMIR, and ICON for the experiments. Specifically, PIFu<sub>w SDF</sub> adds SDF following equations:  $\phi(\mathcal{F}_s(\mathbf{p}), f_{2D}(\mathcal{I})(\mathbf{p})) \rightarrow \hat{\mathcal{O}}$ . PaMIR<sub>w SDF</sub> incoorporates SDF following:  $\phi(\mathcal{F}_s(\mathbf{p}), f_{2D}(\mathcal{I})(\mathbf{p}), \mathcal{V}(\mathcal{M})(\mathbf{p})) \rightarrow \hat{\mathcal{O}}$ . ICON<sub>w/o SDF</sub> replaces SDF with the z coordinate of  $\mathbf{p}$  following:  $\phi(\mathcal{F}_n^b(\mathbf{p}), \mathcal{F}_n^c(\mathbf{p}), \mathbf{p}_z) \rightarrow \hat{\mathcal{O}}$ , where  $\phi$  is the vanilla implicit function and  $f_{2D}$  is a 2D CNN. Experimental results in Fig. 2, and Tab. 1 demonstrate the SDF improves geometry details compared with variant methods without it.

**Impact of voxelized SMPL-X**  $M_v^{3D}$ . Our empirical verification is based on ICON, D-IF and HiLo that requires SMPL-X for reconstruction. we design three variants named ICON<sub>w</sub>  $\mathcal{M}_v^{3D}(p)$ , D-IF<sub>w</sub>  $\mathcal{M}_v^{3D}(p)$  and HiLo<sub>w/o</sub>  $\mathcal{M}_v^{3D}(p)$  that add or remove  $M_v^{3D}$ . From the experimental results in Fig. 2 and Tab. 3, we find that incorporating  $\mathcal{M}_v^{3D}$  helps achieve a more robust reconstruction even faces various levels of noise in SMPL-X shape and pose. The reason is that the low-resolution voxel grid of the

<sup>&</sup>lt;sup>3</sup>We will release our code.

Table 2. (A) Comparison experiments and (B) ablation studies on seen (*i.e.*, training and test on the same dataset) and unseen (*i.e.*, training on Thuman2.0 and test on CAPE) settings. The **bold** and the <u>underlined</u> numbers indicate the best and second-best results, respectively. "-" denotes that PIFuHD and ECON does not provide a training code.

			Train on Thuman2.0 and test on CAPE									Train and test on the same dataset (Thuman2.0 or CAPE)					
Group	Dataset	set CAPE-FP		þ	CAPE-NFP			CAPE			Thuman2.0			CAPE			
	Methods	Chamfer $(\downarrow)$	P2S $(\downarrow)$	Normals $(\downarrow)$	Chamfer	P2S	Normals	Chamfer	P2S	Normals	Chamfer	P2S	Normals	Chamfer	P2S	Normals	
	PIFu [44]	2.1000	2.0930	0.0910	2.9730	2.9400	0.1110	2.6820	2.6580	0.1040	2.6880	2.5730	0.0970	7.1244	2.7633	0.3902	
	PIFuHD [45]	2.3020	2.3350	0.0900	3.7040	3.5170	0.1230	3.2370	3.1230	0.1120	2.4613	2.3605	0.0924	-	-	-	
	PaMIR [64]	1.2250	1.2060	0.0550	1.4130	1.3210	0.0630	1.3500	1.2830	0.0600	1.4388	1.5613	0.1071	0.9339	0.9444	0.0659	
A	ICON [53]	0.7475	0.7488	0.0508	0.8656	0.8639	0.0545	0.8055	0.8084	0.0539	1.1431	1.3020	0.0923	0.8610	0.8878	0.0606	
	ECON [54]	0.9651	0.9175	0.0412	0.9983	0.9694	0.0415	0.9872	0.9521	0.0414	-	-	-	-	-	-	
	D-IF [55]	0.8038	0.7766	0.0546	0.9877	0.9491	0.0611	0.8878	0.8574	0.0589	<u>1.0305</u>	<u>1.0864</u>	<u>0.0830</u>	<u>0.8332</u>	<u>0.8489</u>	0.0597	
	$\operatorname{HiLo}_{w/o \mathcal{H}(s;\beta)}$	0.7564	0.7449	0.0514	0.8697	0.8658	0.0553	0.8118	0.8045	0.0547	0.9442	1.0323	0.0785	0.7971	0.7999	0.0551	
В	$HiLo_{w/o \phi_{si}}$	0.7996	0.7860	0.0569	0.9112	0.9042	0.0601	0.8555	0.8449	0.0468	0.9886	1.0836	0.0850	0.7999	0.7948	0.0547	
	$\mathrm{HiLo}_{\mathrm{w/o}}\mathrel{_{\mathcal{H}(s;\beta)}}_{\mathrm{w/o}} \mathrel{_{\phi_{si}}}$	0.8662	0.8970	0.0647	1.0201	1.0525	0.0706	0.9362	0.9720	0.0690	1.1220	1.2544	0.0954	0.8125	0.8224	0.0588	
	HiLo (Ours)	0.6954	0.6876	0.0471	0.7830	0.7876	0.0499	0.7430	0.7428	0.0499	0.9230	0.9855	0.0732	0.7861	0.7729	0.0544	

Table 3. Toy experiments on 3D reconstruction with different levels of SMPL-X noise in terms of chamfer distance on unseen CAPE dataset. The voxel grid of naked body  $\mathcal{M}_v^{3D}(p)$  improves the robustness of reconstruction.

		SMI	PL-X Noise=0.	1	SMI	PL-X Noise=0.	2	SMPL-X Noise=0.5			
Methods	$\mathcal{M}_v^{3D}(p)$	CAPE-FP	CAPE-NFP	CAPE	CAPE-FP	CAPE-NFP	CAPE	CAPE-FP	CAPE-NFP	CAPE	
ICON [53]	×	1.7949	2.0537	1.9284	2.6695	3.0917	2.9365	4.2181	4.8266	4.6716	
$\mathrm{ICON}_{\mathrm{w}} \mathcal{M}_{v}^{3D}(p)$	1	1.4381	1.5380	1.3411	2.1950	2.3067	2.0723	2.2129	2.3760	2.1188	
D-IF [55]	×	1.6078	1.7881	1.7037	2.6853	3.1002	2.9962	4.3591	4.8006	4.7310	
$\text{D-IF}_{\le \mathcal{M}_v^{3D}(p)}$	1	1.3557	1.7877	1.6064	1.8022	2.4110	2.1959	3.0307	4.4190	4.0807	
$\operatorname{HiLo}_{w/o \mathcal{M}^{3D}_{v}(p)}$	×	1.9518	2.1966	2.0877	2.9315	3.4561	3.2661	4.6382	5.0606	4.9844	
HiLo	1	1.0517	1.3210	1.2004	1.0893	1.5876	1.3427	1.0960	1.6156	1.3593	



Figure 8. Reconstructions w and w/o our spatial interaction implicit function  $\phi_{si}$ . Our  $\phi_{si}$  is able to perceive the global human body and is therefore able to remove non-human shapes.

naked body is insensitive to noise, and therefore provides robust low-frequency regularization in the training process.

### 4.2. Comparison Experiments

**Quantitative results.** We conducted comparative experiments in Tab. 2 under two settings. 1) *Setting1*: Following the setting of the previous methods, we train and test on the same datasets. 2) *Setting2*: To further evaluate the general-



Figure 9. Visualization comparisons on CAPE dataset. The model is training on Thuman2.0 dataset.

ization ability of our HiLo on unseen datasets, we train and test HiLo using different datasets. Our approach achieves the best results in the seen and the unseen settings due to the high- and low-frequency paradigm.

**Visualization Results.** We compare our HiLo with baselines on in-the-wild images and CAPE dataset in Fig. 1, and Fig. 9, respectively. The results show that our HiLo is able to reconstruct 3D clothed avatars with more realistic details. Although ECON obtains detailed fingers by replacing the hand of the SMPL-X model, there exists misalignment on the connection wrist when the corresponding SMPL-X is inaccurate. We put more visualization results of our HiLo on in-the-wild images in Fig. 6. The results demonstrate the ef-

Table 4. Ablations of our progressive high-frequency SDF on 3D reconstruction with different levels of SMPL-X noise in terms of chamfer distance on unseen CAPE dataset. In addition to  $\mathcal{M}_v^{3D}$ , our progressive high-frequency SDF is also to handle SMPL-X noise due to the coarse-to-fine learning manner.

					SM	PL-X Noise=0.	1	SM	PL-X Noise=0.	2	SMPL-X Noise=0.5			
Methods	$\mathcal{M}_v^{3D}(p)$	$\mathcal{H}_s(p;\beta)$	$\mathcal{H}_s(p)$	SDF	CAPE-FP	CAPE-NFP	CAPE	CAPE-FP	CAPE-NFP	CAPE	CAPE-FP	CAPE-NFP	CAPE	
$\operatorname{HiLo}_{W/o} \mathcal{H}_{s}(p;\beta)$	1	X	X	1	1.1435	1.4700	1.3124	1.3401	1.9339	1.6909	1.2861	1.8200	1.5620	
$HiLo_{W} \mathcal{H}_{s}(p)$	<ul> <li>✓</li> </ul>	×	1	×	1.1932	1.5541	1.3904	1.1243	1.5575	1.3701	1.2794	1.8973	1.6652	
HiLo	<ul> <li>Image: A set of the set of the</li></ul>	1	×	×	1.0517	1.3210	1.2004	1.0893	1.5876	1.3427	1.0960	1.6156	1.3593	

fectiveness and generalization ability of our HiLo in recovering detailed geometry (such as hairs and cloth wrinkles). We put more visualization results in the Appx.

### 4.3. Ablation Studies

How does  $\mathcal{H}(\mathbf{s};\beta)$  improve geometry details? We quantitatively demonstrate the necessity of  $\mathcal{H}(s;\beta)$  in Tab. 2. The results demonstrate that  $\operatorname{HiLo}_{w/o \mathcal{H}(s;\beta)}$ , the variant method that replaces  $\mathcal{H}(s;\beta)$  with standard SDF, achieves inferior performance than HiLo. To further study the impact of  $\mathcal{H}(s;\beta)$  on common (-FP) and challenging poses (-NFP), we evaluate HiLo on cape dataset that contains both categories. Tab. 2 demonstrates that  $\mathcal{H}(s;\beta)$  improves the performance of avatar reconstruction more on challenging poses (9.54% improvement) than in fashion (5.21%)improvement in terms of Chamfer distance. Furthermore, Fig. 7 demonstrates that  $\mathcal{H}(s;\beta)$  leads to more detailed reconstruction, resulting in clearer arms and more realistic cloth wrinkles. From the results, we observe that incorporating the power of high frequency with SDF helps in capturing detailed geometry.

How does  $\phi_{si}$  improve body topology of the reconstructed avatar? As shown in Fig. 8 and Tab. 2, our  $\phi_{si}$  removes the non-human shape and boosts reconstruction performance. The reason is that our  $\phi_{si}$  leverages a cross-scale attention module  $\mathcal{A}$  that builds topological signals between different spatial points in the body model.

## 4.4. Further Discussions

Is our  $\mathcal{H}_{s}(\mathbf{p};\beta)$  able to help HiLo be robust to SMPL-X noise? We study the impact of  $\mathcal{H}_{s}(p;\beta)$  on the robustness ability of our HiLo by replacing it with conventional high frequency SDF  $\mathcal{H}_{s}(p)$  and vanilla SDF. We perturb the SMPL-X model with various levels of noise to compare the robustness of the SDF variants. As illustrated in Tab. 4, our proposed  $\mathcal{H}_{s}(p;\beta)$  outperforms SDF and highfrequency SDF variants under multiple noise scales due to the progressive manner. See more results in the Appx.

Is HiLo able to converge faster? In the comparison of validation accuracy depicted in Fig. 10, it is evident that our HiLo exhibits a remarkable ability to rapidly converge and attain superior performance. Specifically, HiLo swiftly reaches a commendable accuracy of 0.90 at approximately iteration 100. In contrast, the second best



Figure 10. Convergence curves of different methods on CAPE dataset. Our HiLo is able to converge faster than existing methods.

method, *i.e.*, ICON, takes significantly longer, around iteration 200, to reach the same level of accuracy, underscoring the efficiency and efficacy of our approach.

### 5. Conclusion

In this paper, we propose a high-frequency and lowfrequency paradigm by exploiting high-frequency and low-frequency information from parametric body models. Based on the paradigm, we design clothed human reconstruction with high- and low-frequency information, namely HiLo that contains: 1) a progressive high-frequency SDF to improve geometry details and alleviate large gradients that hinder model convergence; 2) a spatial interaction implicit function that utilizes the low-frequency complementary information from the voxelized naked body to improve robustness against noise. Experimental results demonstrate the superiority of our HiLo. In the future, we will apply our method to more 3D reconstruction tasks such as 3D face reconstruction, and indoor scene 3D reconstruction.

### 6. Acknowledgement

This work was partially supported by National Natural Science Foundation of China (NSFC) 62072190, Program for Guangdong Introducing Innovative and Entrepreneurial Teams 2017ZT07X183, and TCL Science and Technology Innovation Fund.

# References

- [1] Thiemo Alldieck, Marcus Magnor, Bharat Lal Bhatnagar, Christian Theobalt, and Gerard Pons-Moll. Learning to reconstruct people in clothing from a single rgb camera. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 1175–1186, 2019. 1
- [2] Thiemo Alldieck, Gerard Pons-Moll, Christian Theobalt, and Marcus Magnor. Tex2shape: Detailed full human body geometry from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2293–2303, 2019. 1
- [3] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. On the convergence rate of training recurrent neural networks. *Advances in neural information processing systems*, 32, 2019.
   4
- [4] Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis. Scape: shape completion and animation of people. In ACM SIGGRAPH 2005 Papers, pages 408–416. 2005. 3, 1
- [5] Sahar Aseeri and Victoria Interrante. The influence of avatar representation on interpersonal communication in virtual social environments. *IEEE transactions on visualization and computer graphics*, 27(5):2608–2617, 2021. 1
- [6] Sean Bell, C Lawrence Zitnick, Kavita Bala, and Ross Girshick. Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2874–2883, 2016. 5
- [7] Bharat Lal Bhatnagar, Garvita Tiwari, Christian Theobalt, and Gerard Pons-Moll. Multi-garment net: Learning to dress 3d people from images. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5420– 5430, 2019. 1
- [8] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Combining implicit function learning and parametric models for 3d human reconstruction. In *The European Conference on Computer Vision*, pages 311–329, 2020. 3, 1
- [9] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. Advances in Neural Information Processing Systems, 33:12909–12922, 2020. 3, 1
- [10] Qingwen Bu, Dong Huang, and Heming Cui. Towards building more robust models with frequency bias. In *Proceedings* of the IEEE/CVF International Conference on Computer Vision, pages 4402–4411, 2023. 2
- [11] Yiting Chen, Qibing Ren, and Junchi Yan. Rethinking and improving robustness of convolutional neural networks: a shapley value-based approach in frequency domain. In Advances in Neural Information Processing Systems, pages 324–337, 2022. 2
- [12] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019. 1

- [13] Enric Corona, Albert Pumarola, Guillem Alenya, Gerard Pons-Moll, and Francesc Moreno-Noguer. Smplicit: Topology-aware generative model for clothed people. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 11875–11885, 2021. 1
- [14] Zijian Dong, Chen Guo, Jie Song, Xu Chen, Andreas Geiger, and Otmar Hilliges. Pina: Learning a personalized implicit neural avatar from a single rgb-d video sequence. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 20470–20480, 2022. 1
- [15] Chantat Eksombatchai, Pranav Jindal, Jerry Zitao Liu, Yuchen Liu, Rahul Sharma, Charles Sugnet, Mark Ulrich, and Jure Leskovec. Pixie: A system for recommending 3+ billion items to 200+ million users in real-time. In *Proceedings of the 2018 world wide web conference*, pages 1775– 1784, 2018. 2
- [16] Quentin Galvane, Rémi Ronfard, Christophe Lino, and Marc Christie. Continuity editing for 3d animation. In *Proceedings* of the AAAI Conference on Artificial Intelligence, 2015. 1
- [17] Kaiwen Guo, Peter Lincoln, Philip Davidson, Jay Busch, Xueming Yu, Matt Whalen, Geoff Harvey, Sergio Orts-Escolano, Rohit Pandey, Jason Dourgarian, et al. The relightables: Volumetric performance capture of humans with realistic relighting. ACM Transactions on Graphics, 38(6): 1–19, 2019. 1
- [18] Tong He, John Collomosse, Hailin Jin, and Stefano Soatto. Geo-pifu: Geometry and pixel aligned implicit functions for single-view human reconstruction. *Advances in Neural Information Processing Systems*, 33:9276–9287, 2020. 1, 2
- [19] Tong He, Yuanlu Xu, Shunsuke Saito, Stefano Soatto, and Tony Tung. Arch++: Animation-ready clothed human reconstruction revisited. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11046–11056, 2021. 3
- [20] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 603–612, 2019. 2
- [21] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. Arch: Animatable reconstruction of clothed humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3093–3102, 2020. 3, 1
- [22] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. Advances in neural information processing systems, 28, 2015. 5
- [23] Boyi Jiang, Juyong Zhang, Yang Hong, Jinhao Luo, Ligang Liu, and Hujun Bao. Bcnet: Learning body and cloth shape from a single image. In *The European Conference on Computer Vision*, pages 18–35. Springer, 2020. 1
- [24] Dongsik Jo, Ki-Hong Kim, and Gerard Jounghyun Kim. Effects of avatar and background representation forms to copresence in mixed reality (mr) tele-conference systems. In SIGGRAPH ASIA 2016 virtual reality meets physical reality: modelling and simulating virtual humans and environments, pages 1–4. 2016. 1

- [25] Daiyuan Li, Guo Chen, Xixian Wu, Zitong Yu, and Mingkui Tan. Cross-stage relation extraction and presentation attack material perception for face anti-spoofing. *Neural Networks*, page 106275, 2024. 2
- [26] Qiufu Li, Linlin Shen, Sheng Guo, and Zhihui Lai. Wavelet integrated cnns for noise-robust image classification. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7245–7254, 2020. 2
- [27] Ren Li, Benoît Guillard, Edoardo Remelli, and Pascal Fua. Dig: Draping implicit garment over the human body. In *Proceedings of the Asian Conference on Computer Vision*, pages 2780–2795, 2022. 1
- [28] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multiperson linear model. *Acm Transactions on Graphics*, 34, 2015. 2, 3, 1
- [29] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. In Seminal graphics: pioneering efforts that shaped the field, pages 347–353, 1998. 3, 5
- [30] Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades, Gerard Pons-Moll, Siyu Tang, and Michael J Black. Learning to dress 3d people in generative clothing. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6469–6478, 2020. 1, 5
- [31] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 5
- [32] 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. 1
- [33] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *The European Conference on Computer Vision*, 2020. 2
- [34] Gyeongsik Moon, Hyeongjin Nam, Takaaki Shiratori, and Kyoung Mu Lee. 3d clothed human reconstruction in the wild. In *European conference on computer vision*, pages 184–200, 2022. 1
- [35] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *The European Conference on Computer Vision*, pages 483–499, 2016. 5
- [36] Seung-Tak Noh, Hui-Shyong Yeo, and Woontack Woo. An hmd-based mixed reality system for avatar-mediated remote collaboration with bare-hand interaction. In Proceedings of the 25th International Conference on Artificial Reality and Telexistence and 20th Eurographics Symposium on Virtual Environments, pages 61–68, 2015. 1
- [37] 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. 2, 3,

- [38] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021. 4
- [39] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pages 1310–1318, 2013. 4
- [40] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems, 32, 2019.
- [41] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 10975–10985, 2019. 2, 3, 1, 4
- [42] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pages 5301–5310. PMLR, 2019. 3
- [43] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pages 5301–5310, 2019. 2
- [44] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2304–2314, 2019. 1, 2, 5, 6, 7
- [45] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 84–93, 2020. 1, 7, 2
- [46] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *The European Conference* on Computer Vision, 2016. 1
- [47] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in neural information processing systems*, 33:7462–7473, 2020. 3
- [48] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 5459– 5469, 2022. 2

- [49] Lin Sun. Research on the application of 3d animation special effects in animated films: Taking the film avatar as an example. *Scientific Programming*, 2022. 1
- [50] Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012. 6
- [51] Ea Christina Willumsen. Is my avatar my avatar? character autonomy and automated avatar actions in digital games. In *DiGRA Conference*, 2018. 1
- [52] Donglai Xiang, Fabian Prada, Chenglei Wu, and Jessica Hodgins. Monoclothcap: Towards temporally coherent clothing capture from monocular rgb video. In 2020 International Conference on 3D Vision (3DV), pages 322–332. IEEE, 2020. 1
- [53] Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J. Black. ICON: Implicit Clothed humans Obtained from Normals. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 13296– 13306, 2022. 1, 2, 4, 5, 6, 7, 3
- [54] Yuliang Xiu, Jinlong Yang, Xu Cao, Dimitrios Tzionas, and Michael J. Black. ECON: Explicit Clothed humans Optimized via Normal integration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 2, 5, 7, 1, 6
- [55] Xueting Yang, Yihao Luo, Yuliang Xiu, Wei Wang, Hao Xu, and Zhaoxin Fan. D-if: Uncertainty-aware human digitization via implicit distribution field. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9122–9132, 2023. 1, 2, 4, 5, 7, 6
- [56] Yifan Yang, Shuhai Zhang, Zixiong Huang, Yubing Zhang, and Mingkui Tan. Cross-ray neural radiance fields for novelview synthesis from unconstrained image collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15901–15911, 2023. 5
- [57] Boram Yoon, Hyung-il Kim, Gun A Lee, Mark Billinghurst, and Woontack Woo. The effect of avatar appearance on social presence in an augmented reality remote collaboration. In 2019 IEEE Conference on Virtual Reality and 3D User Interfaces (VR), pages 547–556. IEEE, 2019. 1
- [58] Hongwen Zhang, Yating Tian, Xinchi Zhou, Wanli Ouyang, Yebin Liu, Limin Wang, and Zhenan Sun. Pymaf: 3d human pose and shape regression with pyramidal mesh alignment feedback loop. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11446–11456, 2021. 2
- [59] Hongwen Zhang, Yating Tian, Yuxiang Zhang, Mengcheng Li, Liang An, Zhenan Sun, and Yebin Liu. Pymaf-x: Towards well-aligned full-body model regression from monocular images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 2
- [60] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2023. 5
- [61] Tianyuan Zhang and Zhanxing Zhu. Interpreting adversarially trained convolutional neural networks. In *International conference on machine learning*, pages 7502–7511, 2019. 2

- [62] Enliang Zheng, Enrique Dunn, Vladimir Jojic, and Jan-Michael Frahm. Patchmatch based joint view selection and depthmap estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1510–1517, 2014. 1
- [63] Zerong Zheng, Tao Yu, Yixuan Wei, Qionghai Dai, and Yebin Liu. Deephuman: 3d human reconstruction from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7739–7749, 2019. 5
- [64] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *IEEE transactions on pattern analysis and machine intelligence*, 44(6): 3170–3184, 2021. 1, 2, 4, 5, 6, 7
- [65] Hao Zhu, Xinxin Zuo, Sen Wang, Xun Cao, and Ruigang Yang. Detailed human shape estimation from a single image by hierarchical mesh deformation. In *Proceedings of* the IEEE/CVF conference on computer vision and pattern recognition, pages 4491–4500, 2019. 1
- [66] Huang Zixiong, Chen Qi, Sun Libo, Yang Yifan, Wang Naizhou, Tan Mingkui, and Wu Qi. G-nerf: Geometryenhanced novel view synthesis from single-view images. arXiv preprint arXiv:2404.07474, 2024. 3