

# GBCN:3D-Guided Brownian Bridge Diffusion for Clothed Human Reconstruction via Normal Integration

Anonymous ICCV submission

Paper ID 8220

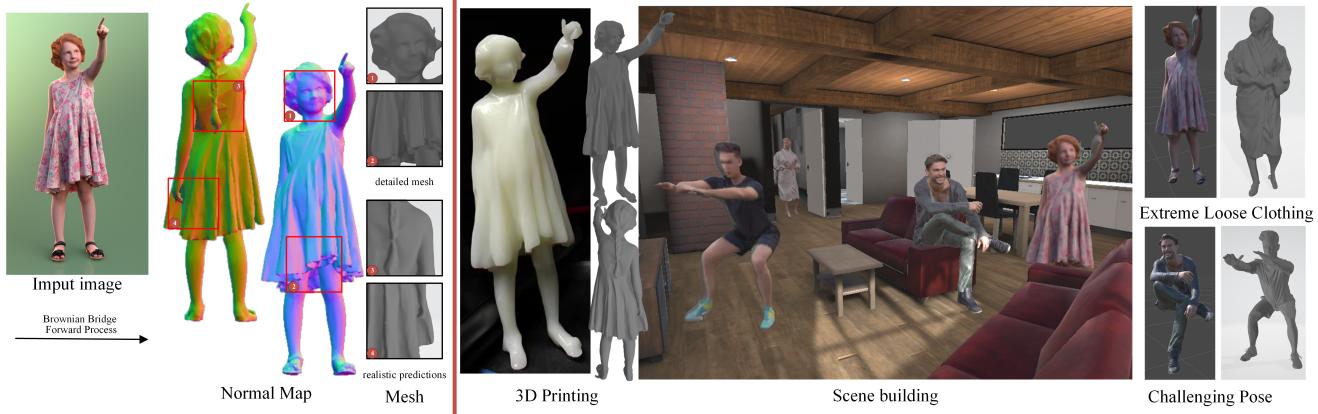


Figure 1. GBCN is a novel human body reconstruction network based on a diffusion model. Given an image, GBCN utilizes the Brownian Bridge diffusion process to map the original image to 3D features in the latent space. It refines normal maps to reconstruct a detailed clothed human mesh. By employing a consistency loss function and an SMPL-X prior decoupling module, GBCN achieves an effective extension of the Brownian Bridge diffusion process from 2D to 3D, bridging large domain gaps. Our method demonstrates excellent performance on in-the-wild images and shows significant potential in real-world applications such as 3D printing and scene construction.

## Abstract

Creating high-quality 3D models of clothed humans from a single image is crucial for real-world applications. Existing models often make insufficient assumptions about the geometric details of invisible areas, leading to less refined outputs. Despite recent progress, 3D diffusion models still treat diffusion as a conditional generation process, significantly impacted by domain gaps. We addressed this by proposing GBCN (3D-Guided Brownian Bridge Diffusion for Clothed Human Reconstruction via Normal Integration) as a novel solution. Specifically, we utilized the Brownian Bridge diffusion process to directly map the image domain to the spatial domain through bidirectional diffusion, avoiding conditional information leverage. Additionally, we progressively guided the diffusion process with 3D pose priors to bridge the gap between 2D and 3D domains. To better learn fine 3D geometric details in the latent space, we employed a consistency loss function to enhance the model's robustness while maintaining sensitivity to fine geometric details, ul-

timately regressing detailed front and back normal maps. Finally, we used these detailed normal maps to reconstruct a fine 3D clothed model through bidirectional normal integration and Poisson surface reconstruction. Our method significantly improves geometric details, including invisible areas, and shows robustness under various poses and loose clothing. Experiments demonstrate that due to fine normal reconstruction, GBCN exhibits impressive performance in geometric reconstruction pipelines, showing enhanced robustness in complex scenarios. We also applied our method to virtual scene construction and 3D printing, showcasing its broad potential for practical applications.

## 1. Introduction

High-quality 3D human avatar models have crucial applications in fields such as virtual reality, education, and gaming. Creating a detailed clothed human model traditionally requires experienced artists and expensive scanning equipment [19, 53]. Generally, single RGB images of humans

019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030

031  
032  
033  
034  
035  
036

are more accessible. Nevertheless, reconstructing clothed humans from a single RGB image still faces three key technical challenges: **(i)** Single RGB images suffer from severe self-occlusion, making inference about large invisible areas a major challenge. **(ii)** In practical applications, people in images may pose with large angles or wear extremely loose clothing, demanding high robustness from the model. **(iii)** Images may contain many details such as clothing wrinkles or accessories, requiring the model to accurately reconstruct these high-frequency details.

Existing models [5, 9, 31, 65, 66, 76, 77, 83, 89–91] typically use various methods to map 2D images to 3D space. Most methods [76, 77, 83, 91] focus on improving the transition from 2D features to 3D using SMPL [59] guidance, but they neglect 3D information during 2D feature extraction. Some methods [89] use fixed learnable embeddings for 3D feature learning. SIFU [90] improves on these methods by using SMPL priors through an attention mechanism to learn 3D information from images and predict the overall mesh. Despite success in handling poses and loose clothing, these methods lack exploration of surface details, limiting the generation of high-precision geometric models, as shown in Figure 2. Therefore, a new approach is needed that maintains robustness to clothing and poses while achieving good performance in geometric details, including invisible areas.

Diffusion models [24, 57, 60, 64, 70, 70, 74], compared to traditional GAN-based models, exhibit better stability and controllability [11, 12]. These methods improve the quality of generated samples by avoiding information loss through a gradual diffusion process from the source domain to the target domain. The Brownian Bridge Diffusion Model (BBDM) [44] models image-to-image translation [32] as a stochastic Brownian Bridge process, achieving progress by directly mapping between two image domains. SiTH [23] successfully applied diffusion models to the field of 3D human reconstruction and achieved impressive results. SiTH first hallucinates back-view appearances through an image-conditioned diffusion model, followed by the reconstruction of full-body textured meshes using both the front and back-view images. However, in essence, the diffusion module in SiTH is still essentially a 2D image-to-image conversion task. While recent diffusion models show promise in image synthesis, their application to 3D reconstruction [23, 30, 38], faces fundamental limitations: **(i)** As BBDM states, most diffusion models treat diffusion as conditional generation and lack direct mapping methods [44]. **(ii)** Unlike image translation tasks, transitioning from 2D images to 3D space often faces domain gaps, making effective mapping difficult without prior conditions. **(iii)** Some text-guided conditional diffusion models also require more computation.

Our key insight is that bidirectional diffusion between

two domains, rather than conditional generation, more effectively guides the model in learning the mapping process. This motivates us to extend BBDM to 3D reconstruction through a prior decoupling module to address significant domain gaps. Specifically, we perform bidirectional diffusion using the original image and normal maps in the model and progressively refine geometric details in the latent space. To better match the mapping between the two domains, we gradually introduce SMPL-X [59] priors into the diffusion process using a cross-attention mechanism and use a consistency loss to guide the model in maintaining sensitivity to geometric details. Through this approach, we restored detailed front and back normal maps from the latent space and used normal integration and Poisson reconstruction [77] to generate the corresponding 2.5D surfaces. We then used SMPL-X to merge the human body mesh, achieving detailed reconstruction of the clothed human mesh without additional parameter training. Experiments show that our model benefits significantly from 3D information, proving it is not just a 2D image conversion task, as detailed in our experimental section.

Through extensive experiments, GBCN has demonstrated superior geometric reconstruction quality. It surpasses existing state-of-the-art (SOTA) methods. It shows excellent robustness in handling complex poses and loose clothing. GBCN excels in preserving detailed geometric features on surfaces, including invisible areas. We have extended the reconstruction results to practical applications. These include 3D printing and scene construction, showcasing its wide applicability. Our main contributions include:

- We improved the Brownian Bridge diffusion process and extended it to human reconstruction, establishing a mapping from 2D images to the 3D latent space and achieving richer reconstruction details.
- The normal map prediction module in GBCN is both simple and portable. Applied to recent works, it showed varying degrees of performance improvement.
- The GBCN pipeline achieved significant progress, surpassing previous methods across all metrics. Its strong performance on real-world images enabled applications in 3D printing and scene construction, highlighting its potential for state-of-the-art (SOTA) reconstruction.

## 2. Related Work

**Diffusion models.** Recent advances in diffusion models [10, 24, 64, 70] have demonstrated superior performance in high-fidelity image synthesis, particularly in preserving geometric details and handling multi-modal distributions. However, their application to 3D human reconstruction poses unique challenges. Most generative models [28, 61, 69, 73, 79] focus on object generation in the 3D creation process but are limited in generating models with realistic details due to the difficulty in collecting large-

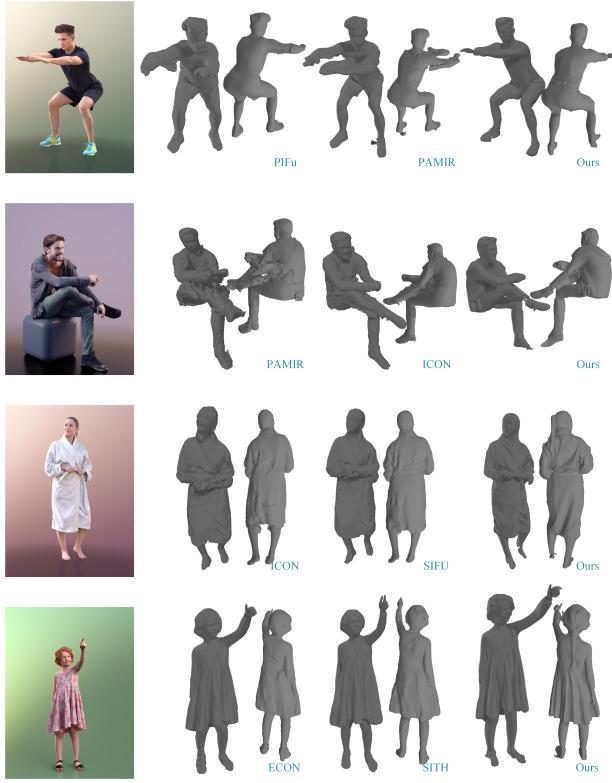


Figure 2. Comparison between GBCN and State-of-the-Art (SOTA) Methods for 3D Human Inference on In-the-Wild Images. Existing SOTA methods consistently underperform when handling complex poses or loose clothing, often producing inference results that lack sufficient detail. In contrast, GBCN effectively overcomes these challenges, yielding high-quality, highly detailed results. We use slightly different angles to fairly show the shortcomings of each model.

scale ground truth. Recent works [8, 23, 27, 30, 38, 67, 80] still treat diffusion models as conditional generative models, which lack direct mapping methods. Some methods integrate image-to-text models [8, 30] or directly use text [38, 51, 52, 55, 71, 80, 90] as conditional input in the U-Net during the reverse process to guide diffusion towards the target domain, which introduces high computational costs. BBDM [44] is a novel diffusion model that first models the image-to-image translation process as a stochastic Brownian Bridge process, making progress by directly mapping between two image domains. Our method extends the Brownian Bridge diffusion to the 3D domain, addressing the significant domain gap and achieving good performance.

**Implicit-function-based Reconstruction.** Implicit representations, such as occupancy fields and SDF, effectively describe 3D clothed humans in various scenarios, including challenging poses and loose clothing. Recent methods [1, 20, 65, 66] employ neural implicit functions to directly

regress 3D surfaces from single images by learning a continuous mapping from 2D features to occupancy/SDF values, achieving topology-aware reconstruction. To address the inherent ambiguity in single-view reconstruction, several methods [6, 7, 9, 21, 22, 30, 31, 50, 76, 77, 83, 89, 91] incorporate SMPL [59] priors to provide strong geometric constraints and anatomical knowledge, thereby improving the accuracy of both 2D feature extraction and 3D surface prediction. For instance, GTA [89] uses a transformer with fixed learnable embeddings to convert image features into 3D planar features. SIFU [90] goes further by incorporating SMPL-X priors into 2D image extraction and using a transformer to extract multi-view features and regress the complete mesh. Despite progress, these methods lack fine surface details and fail to assume reasonably about invisible surfaces.

**Explicit-shape-based Reconstruction.** In contrast to implicit methods that learn continuous functions, explicit shape-based approaches [14, 15, 37, 39–42, 45–48, 68, 86, 87] directly estimate discrete representations (e.g., meshes [36, 54, 58, 59, 78], depth maps [85], or point clouds [85]) and utilize clothing offsets [2–4, 43, 75] to model surface details. While these methods offer direct geometric control, their fixed topology assumption fundamentally limits their ability to handle significant topology changes in loose clothing, model complex self-intersections, and represent fine-scale surface details. ECON [77] constructs partial surfaces by integrating front and back normal maps and stitching them with the implicit IF-Net and SMPL meshes using Poisson surface reconstruction. Although ECON’s hybrid strategy enhances surface reconstruction, the use of less detailed normal maps and the polishing process of IF-Net can result in meshes lacking detail.

**NeRF-based Reconstruction.** Neural Radiance Fields (NeRF) [17, 18, 25, 34, 35, 49, 56, 74, 88] revolutionize novel view synthesis by learning continuous volumetric representations through implicit MLPs and differentiable volume rendering. This enables high-fidelity reconstruction of both geometry and appearance. SHER [26] and ELICIT [29] aim to reconstruct 3D clothed humans from a single image. However, applying NeRF to single-image human reconstruction faces several fundamental challenges: computational complexity due to volumetric optimization requiring extensive resources, view dependency causing significant quality degradation with limited input views, geometry ambiguity resulting in imprecise recovered geometry in occluded regions without multi-view constraints, and training efficiency issues with per-scene optimization making it difficult to learn generalizable priors.

Inspired by previous work, our method achieved impressive results in 3D clothed human reconstruction across various scenarios. We used the Brownian Bridge diffusion process, employing bidirectional diffusion with the original im-

age and ground-truth normal maps to guide the network in refining geometric details in the latent space. By progressively incorporating SMPL-X priors into the diffusion process, we successfully addressed large domain gaps. We also applied a combined consistency loss to constrain the model, ensuring robustness while maintaining sensitivity to fine geometric details.

### 221 3. Method

222 Our pipeline consists of three key components: **(i)** A Brownian  
 223 Bridge diffusion framework that learns the mapping between  
 224 2D and 3D spaces(Section 3.1). **(ii)** A novel consistency  
 225 loss that ensures geometric detail preservation (Section  
 226 3.2). **(iii)** A reconstruction pipeline that leverages  
 227 normal maps for high-fidelity surface generation (Section 3.3).  
 228 Specifically, GBCN first uses the Brownian Bridge  
 229 diffusion process to learn 3D information in the latent space  
 230 through bidirectional diffusion from the original image and  
 231 normal maps, while attempting to refine these details. To  
 232 address large domain gaps, we use SMPL-X priors and  
 233 progressively integrate them into the diffusion process. Once  
 234 the model learns the mapping relationship during diffusion,  
 235 we can directly regress detailed front and back normal maps  
 236 from the latent space without intermediate steps. To en-  
 237 hance the model’s sensitivity to fine details while maintain-  
 238 ing overall perception, we apply a consistency loss in the  
 239 latent space to constrain the network and guide it in refi-  
 240 ning geometric details. Finally, we use the regressed detailed  
 241 normal maps to reconstruct the front and back 2.5D sur-  
 242 faces, completing them with human priors. For the com-  
 243 plete pipeline, refer to Fig.3.

#### 244 3.1. 3D-Guided Brownian Bridge Diffusion.

245 Given an RGB image input  $I$  and its corresponding nor-  
 246 mal maps  $N$ , we formulate the reconstruction process as  
 247 a bidirectional mapping problem. Unlike traditional dif-  
 248 fusion models that treat the endpoint as noise, we innovatively  
 249 use the clean input image  $x_t$  as the endpoint  $x_t$ , while set-  
 250 ting the target normal maps as the starting point  $x_0$ . This  
 251 design allows us to: **(i)** Maintain high-fidelity image fea-  
 252 tures throughout the diffusion process. **(ii)** Explicitly learn  
 253 the geometric mapping in the latent space. **(iii)** Better pre-  
 254 serve structural consistency during diffusion. To accelerate  
 255 training and inference, we retrained the popular VQGAN  
 256 and performed diffusion in the latent space. For details on  
 257 traditional diffusion models, please refer to the supple-  
 258 mentary materials. We represent the state distribution at each  
 259 time step of a Brownian bridge process starting from point  
 260  $(\mathbf{x}_0) \sim q_{data}(\mathbf{x}_0)$  at  $t = 0$  and ending at point  $x_t$  at  $t = T$   
 261 can be formulated as:

$$262 p(x_t | x_0, x_T) = \mathcal{N}\left((1 - \frac{t}{T})\mathbf{x}_0 + \frac{t}{T}\mathbf{x}_T, \frac{t(T-t)}{T}\mathbf{I}\right) \quad (1)$$

The two ends of the process are anchored by  $x_0$  and  $x_t$ ,  
 263 forming a bridge spanning two domains.

When predicting the front normal maps, we applied the  
 264 diffusion process to the original image pairs and made  
 265 progress. However, for the back normal maps and original  
 266 image pairs, the model failed to estimate the mapping be-  
 267 tween the two domains due to the large domain gap (see ex-  
 268 periment section for details). Although back normal maps  
 269 differ in viewpoint, they maintain several critical corre-  
 270 lations with the input image: **(i) Geometric consistency:** The  
 271 overall body structure remains coherent. **(ii) Feature cor-  
 272 respondence:** Local geometric details (e.g., clothing wrin-  
 273 kles) follow consistent patterns. **(iii) Topological relation-  
 274 ships:** The spatial arrangement of body parts preserves rel-  
 275 ative positioning. We leverage these inherent correlations  
 276 through a cross-attention mechanism to guide the genera-  
 277 tion of occluded regions. The specific formula is as follows:  
 278

$$279 CrossAttn(z_i, h) = \text{Softmax}\left(\frac{(W^Q h)(W^K z_i)^T}{\sqrt{d}}\right)(W^V z_i) \quad (2)$$

where  $h$  represents the embedding extracted from the  
 280 SMPL-X priors via a feature encoder.  $z_i$  is the inter-  
 281 mediate feature sampled during the diffusion process, and  
 282  $W^Q$ ,  $W^K$ , and  $W^V$  are learnable parameters. We believe  
 283 this design allows the model to selectively integrate prior  
 284 knowledge and adaptively handle viewpoint changes while  
 285 maintaining geometric consistency. We follow the tradi-  
 286 tional cross-attention design, to add norm layers and resi-  
 287 dual connections after each layer. Our feature extractor is  
 288 multi-scale and performs cross-attention operations at dif-  
 289 ferent stages of the diffusion process. This design enables  
 290 the model to effectively learn spatial information, including  
 291 back-view perspectives. It also allows the model to decou-  
 292 ple features from the original image, successfully mapping  
 293 front-view and back-view 3D information.

#### 294 3.2. Normal consistency loss.

295 The key challenge in normal map prediction lies in pre-  
 296 serving geometric accuracy while maintaining spatial con-  
 297 sistency. Unlike regular RGB images with purely vi-  
 298 sual color variations, normal maps encode crucial geomet-  
 299 ric information through their RGB channels, specifically  
 300 representing surface orientations in tangent space. This  
 301 special texture pattern makes the lighting effects resem-  
 302 blle those of real geometric shapes. By learning this pat-  
 303 tern, the model can refine spatial details in the 3D latent  
 304 space and produce logical, detailed textures, avoiding frag-  
 305 mented reconstruction errors. An intuitive method is to  
 306 constrain the model’s learning of texture patterns through  
 307 color consistency loss, but this faces challenges in bidi-  
 308 rectional diffusion latent space. Despite perceptual con-  
 309 sistency, VQGAN-encoded high-dimensional features sta-  
 310 tistically deviate from the original input, leading to uncer-  
 311

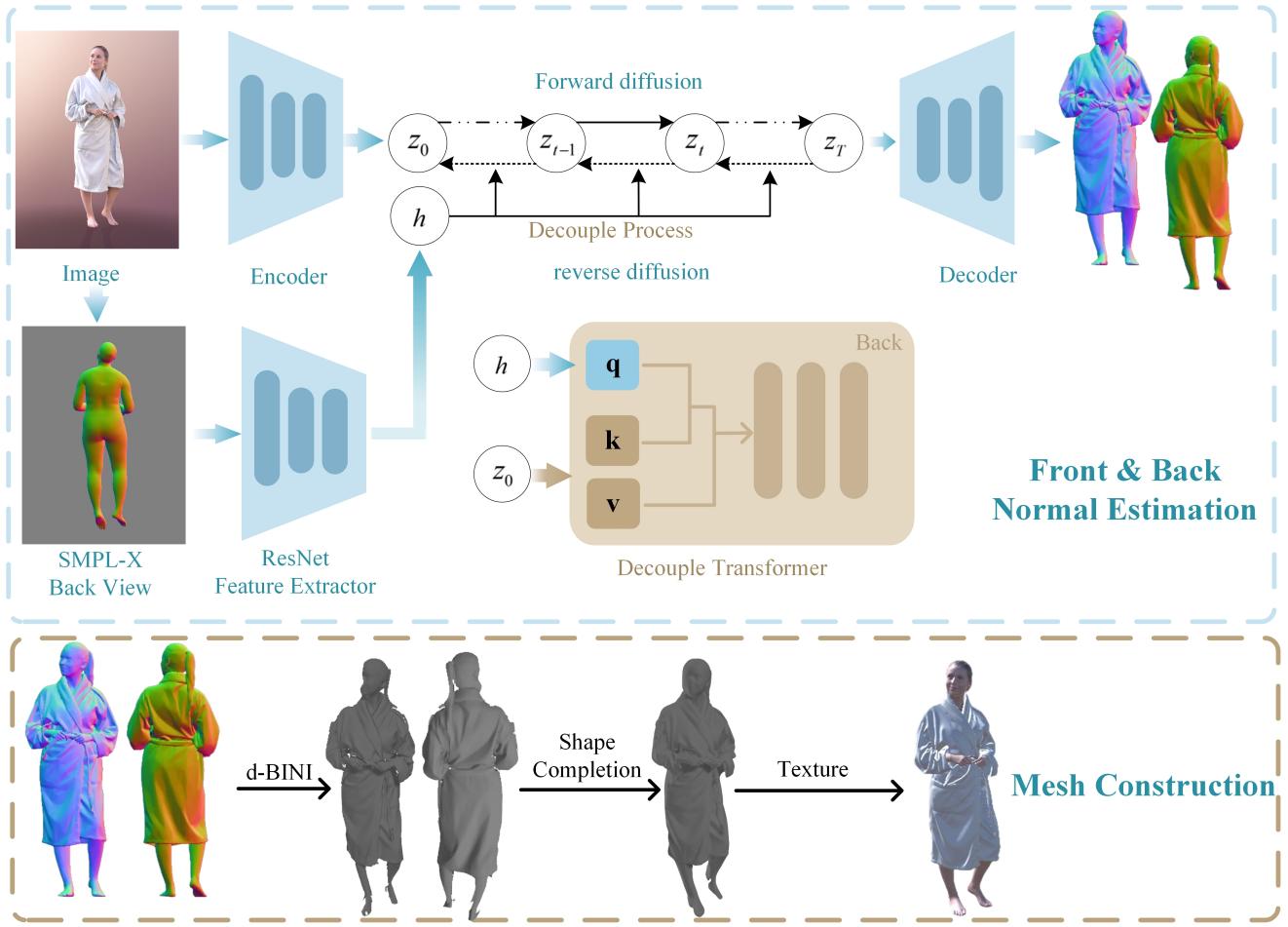


Figure 3. Given an image, GBCN utilizes the Brownian bridge diffusion process (§3.1) to directly predict the front and back normal maps of a clothed human body. During the back normal map prediction, a cross-attention mechanism is employed to integrate SMPL back-view priors, where a prior-mixing strategy incorporates these priors as query inputs, thereby enhancing the model’s capability to infer occluded regions from the image. To facilitate more accurate learning of normal map color patterns, a series of hybrid consistency losses (§3.2) are applied to guide the network’s training process. The subsequent step involves 3D surface reconstruction (§3.3): specifically, under the guidance of SMPL-X estimates, GBCN converts the normal maps into incomplete 2.5D front and back surfaces and reconstructs the missing geometric details. Finally, a pre-trained texture prediction network is used to produce a high-quality textured mesh that can be directly employed in practical applications.

313      tainties in conventional losses within the latent space. Our  
 314      method uses cosine similarity loss in the latent space to  
 315      constrain the direction of high-dimensional vectors. We have  
 316      the following equation:

$$317 \quad \mathcal{L}_{\cos}(\mathbf{a}, \mathbf{b}) = 1 - \cos_{\text{sim}}(\mathbf{a}, \mathbf{b}) \quad (3)$$

318      where  $\mathbf{a}$  and  $\mathbf{b}$  represent the predicted vector and the ground  
 319      truth vector, respectively.  $\cos_{\text{sim}}(\mathbf{a}, \mathbf{b})$  is the cosine similar-  
 320      ity of vectors  $\mathbf{a}$  and  $\mathbf{b}$ , defined as:

$$321 \quad \cos_{\text{sim}}(\mathbf{a}, \mathbf{b}) = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \cdot \sqrt{\sum_{i=1}^n b_i^2}} \quad (4)$$

To capture both local and global geometric patterns, we were inspired by the traditional Discrete Fourier Transform (DFT) [16, 33, 62, 81, 82, 92] to convert high-dimensional features in the latent space into frequency domain components to describe vector patterns. We perform the Discrete Fourier Transform on the original latent space vectors using a two-dimensional DFT:

$$322 \quad F(u, v) = \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} f(x, y) \cdot e^{-i2\pi(\frac{ux}{M} + \frac{vy}{N})} \quad (5)$$

330      where  $f(x, y)$  represents the pixel value at coordinate  
 331       $(x, y)$ , and  $F(u, v)$  represents the complex frequency value



Figure 4. Qualitative results on in-the-wild images. We present nine examples of 3D humans with detailed clothing reconstructed from wild images. These examples include loose clothing and challenging poses. For each example, we show the input image, the predicted front and back normal maps, and the reconstruction results. Our method shows impressive results.

332 at coordinate  $(u, v)$  in the frequency domain. The amplitude  
333 and phase of the frequency are given by:

$$\begin{cases} |F(u, v)| = \sqrt{R(u, v)^2 + I(u, v)^2}, \\ \angle F(u, v) = \arctan\left(\frac{I(u, v)}{R(u, v)}\right) \end{cases} \quad (6)$$

335 where  $R(u, v)$  and  $I(u, v)$  are the real and imaginary components  
336 of  $F(u, v)$ . To compute the frequency distance between the real and generated images, we use the Euclidean  
337 distance:  
338

$$\mathcal{L}_{ffl} = d(Fr, Ff) = \frac{1}{MN} \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} |Fr(u, v) - Ff(u, v)|^2 \quad (7)$$

339 In addition to the aforementioned losses, our research found  
340 that MSE loss corrects larger discrepancies but neglects  
341 some details. This results in many reconstruction artifacts  
342 and poor performance during the reconstruction process.  
343 Specific experimental results are discussed in the experiment  
344 section. Therefore, we only used L1 loss for constraints.  
345

$$\mathcal{L}_{total} = \lambda_1 \mathcal{L}_{11} + \lambda_2 \mathcal{L}_{ffl} + \lambda_3 \mathcal{L}_{cos} \quad (8)$$

346 where  $\lambda_1, \lambda_2, \lambda_3$  are the weights attributed to each loss.  
347

### 3.3. 3D Surface Reconstruction

348 For the detailed normal maps, we use variational normal integration for explicit reconstruction of the front and back surfaces. Specifically, we apply customized bilateral normal integration to form 2.5D surfaces and integrate them using Poisson reconstruction and SMPL-X surfaces [77]. Previous methods added an implicit function to polish and complete the surfaces to address artifacts caused by various poses or incomplete photos. However, this often made the predicted surfaces too smooth. Our method uses the detailed and complete normal maps. This allows us to use Poisson reconstruction to directly reconstruct the 2.5D surfaces without implicit polishing, achieving good results.

## 4. Experiments

### 4.1. Dataset and Metrics

349 We trained our model on Thuman2.0 and tested it on the CAPE dataset. The 526 human scans in Thuman2.0 were  
350 divided into 490 for training, 15 for validation, and 21 for testing. Ground-truth  
351

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

Method	Publication	CAPE-NFP			CAPE-FP			THuman2.0		
		Chamfer ↓	P2S ↓	Normal ↑	Chamfer ↓	P2S ↓	Normal ↑	Chamfer ↓	P2S ↓	Normal ↑
<b>w/o SMPL-X body prior</b>										
<b>PIFu[65]</b>	ICCV 2019	3.2551	2.5492	0.7651	1.8391	1.7612	0.8596	1.2097	1.1324	0.7707
<b>PIFuHD[66]</b>	CVPR 2020	2.9766	2.3700	0.7680	1.5228	1.4860	0.8741	0.9953	0.9676	0.7912
<b>PaMIR[91]</b>	TPAMI 2021	7.1599	3.3852	0.6372	6.0132	3.2895	0.6759	1.0898	1.0169	0.7965
<b>SITH[23]</b>	CVPR 2024	2.8753	2.1244	0.7822	2.1158	1.6769	0.8353	0.9677	0.9052	0.7857
<b>GBCN-PIFuHD</b>	-	2.6909	2.1533	0.8014	1.2622	1.2803	0.9035	0.7485	0.6627	0.8149
<b>GBCN-SITH</b>	-	2.6617	1.9239	0.8164	1.8256	1.4044	0.8533	0.8098	0.8041	0.8113
<b>w/ SMPL-X body prior</b>										
<b>ICON[76]</b>	CVPR 2022	1.5521	1.1997	0.8584	0.9962	0.8881	0.9209	0.6173	0.5956	0.8518
<b>ECON[77]</b>	CVPR 2023	1.8541	1.5721	0.8414	1.1779	1.1369	0.8989	0.6741	0.6349	0.8377
<b>D-IF[83]</b>	ICCV 2023	1.8324	1.5637	0.8634	1.1493	1.1435	0.8893	0.6725	0.6331	0.8364
<b>GTA[89]</b>	NeurIPS 2023	1.8866	1.4929	0.8278	1.1496	0.9937	0.9021	0.5809	0.5601	0.8412
<b>SIFU[90]</b>	CVPR 2024	1.5758	1.2803	0.8545	1.0549	0.9702	0.9040	0.5774	0.5595	0.8512
<b>GBCN-ICON</b>	-	1.2654	1.0980	0.8860	0.8424	0.8251	0.9277	0.5844	0.5670	0.8497
<b>GBCN(Ours)</b>	-	<b>1.2189</b>	<b>0.9793</b>	<b>0.9011</b>	<b>0.7935</b>	<b>0.7083</b>	<b>0.9337</b>	<b>0.5654</b>	<b>0.5476</b>	<b>0.8702</b>

Table 1. Quantitative evaluation against SOTA (Section 4.1). All models use a resolution of 256 for marching cubes and ground-truth SMPL-X models are used during testing.

SMPL-X models were used during training, and PIXIE was employed for inference. To enhance robustness, we added random noise to the original images during training to simulate real environment occlusions. All training and testing were conducted on an NVIDIA GeForce RTX 3090 GPU. We adhered to the common design of previous approaches, using Chamfer and Point-to-Surface (P2S) distances to compute larger errors and measuring normal errors to assess reconstruction quality.

## 4.2. Evaluation

We compared GBCN with state-of-the-art (SOTA) models across varying data volumes. To evaluate our model, we split CPAE into CAPE-FP and CAPE-NFP subsets following established conventions. Ground-truth SMPL models were provided for methods requiring SMPL-X body priors. The results are presented in Tab.1. To further validate our approach, we applied GBCN to other frameworks, including PIFuHD, ICON, and SiTH. By replacing their normal prediction modules, we observed improved performance metrics despite constraints from other internal modules. GBCN effectively mapped input images to 3D space and generated detailed normal maps. Qualitative results are displayed in Fig.4. Tests on in-the-wild data confirmed that our model produced accurate normal maps and reconstructed detailed clothed human meshes, performing particularly well with loose clothing and complex poses.

## 4.3. Ablation Studies

The quantitative results of the ablation study are summarized in Table 2. Without the SMPL prior, even the SOTA image-to-image translation model BBDM shows only limited improvement. Incorporating the SMPL prior through our decoupling module significantly enhances performance, highlighting the importance of 3D information in gener-

Method	Chamfer ↓	P2S ↓	Normal ↑
ECON	0.6741	0.6349	0.8377
GBCN w/o SMPL + L1	0.6332	0.6054	0.8448
GBCN w/ L2	0.6354	0.6048	0.8478
GBCN w/ L1	0.6046	0.5995	0.8574
GBCN w/ L1 + FFL	0.5725	0.5537	0.8633
GBCN w/ L1 + Cos	0.5833	0.5638	0.8637
GBCN(ours)	<b>0.5654</b>	<b>0.5476</b>	<b>0.8702</b>

Table 2. Ablation experiments of GBCN on Thuman2.0 dataset. The experiments showed that each module and loss in our method improved the model’s performance.

ating normal maps, which cannot be treated as a simple 2D transformation task. Qualitative results are presented in Fig.5. Without the decoupling module and human pose priors, the network struggles to infer detailed back normal maps from front images. The large domain gap between back normal maps and original images hinders effective mapping to the 3D latent space. Integrating SMPL-X priors into the diffusion model’s information flow improves the network’s understanding of back pose structures and prevents distortions.

During network training, we used different loss mixing strategies to supervise the network and compared the cross-entropy (MSE) loss values between the predicted and ground-truth normal maps, as shown in Tab.2 and Fig.6. Our qualitative results show that while L2 supervision slightly reduced loss values, it caused large artifacts in the images, leading to instability in subsequent surface reconstruction. Considering the impact of different losses, we adopted the mixed loss strategy shown in Equation 8.



Figure 5. We conducted ablation experiments to evaluate the conditional decoupling module by selecting a subset of images from the training set and feeding them into the converged model. Qualitative results indicate that, without the SMPL-X prior, the network is unable to effectively learn the normal map for unseen regions.

422

#### 4.4. Applications

423 GBCN uses a pipeline similar to ECON during reconstruction,  
424 allowing quick adaptation to other scenarios. Our  
425 method leverages existing image-text texture generation  
426 models to recover model textures. Thanks to GBCN's de-  
427 tailed reconstruction, we can directly 3D print the output  
428 files without processing. Additionally, we use reconstructed  
429 results with texture edits for scene building, as shown in  
430 Fig.1 and Fig.3. In summary, GBCN's results can easily  
431 extend to real-world applications, offering significant pot-  
432 tential in education, gaming, and virtual reality.

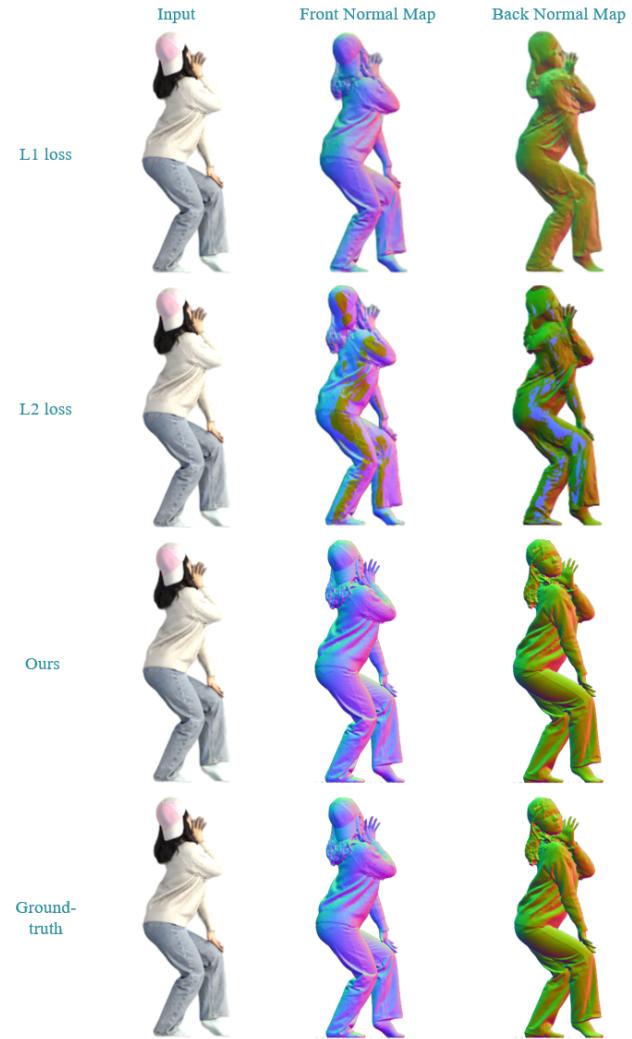


Figure 6. Qualitative evaluation of different loss combinations demonstrates the effectiveness of our chosen combined consistency loss. Notably, although the L2 loss achieved lower reconstruction error in quantitative evaluation, it introduced artifacts in the images, particularly noticeable in the front view.

## 5. Conclusion

We propose GBCN, an innovative diffusion model that refines the Brownian Bridge diffusion process and introduces it to human reconstruction. This method delivers notable results, effectively addressing the substantial domain gap between 2D and 3D. The normal map prediction module in GBCN is simple, efficient, and highly portable. Applied to recent works, it demonstrated state-of-the-art performance across multiple metrics and impressive results on natural images. Additionally, we further extend the outputs of GBCN to practical applications, showcasing its strong potential in tasks such as 3D printing and scene construction.

433

434

435

436

437

438

439

440

441

442

443

444

445

446 

## References

- 447 [1] Badour AlBahar, Shunsuke Saito, Hung-Yu Tseng, Changil  
448 Kim, Johannes Kopf, and Jia-Bin Huang. Single-image 3d  
449 human digitization with shape-guided diffusion. In *SIG-  
450 GRAPH Asia 2023 Conference Papers*, pages 1–11, 2023.  
451 3
- 452 [2] Thimo Alldieck, Marcus Magnor, Weipeng Xu, Christian  
453 Theobalt, and Gerard Pons-Moll. Detailed human avatars  
454 from monocular video. In *2018 International Conference on  
455 3D Vision (3DV)*, pages 98–109. IEEE, 2018. 3
- 456 [3] Thimo Alldieck, Marcus Magnor, Weipeng Xu, Christian  
457 Theobalt, and Gerard Pons-Moll. Video based reconstruction  
458 of 3d people models. In *Proceedings of the IEEE Conference  
459 on Computer Vision and Pattern Recognition*, pages 8387–  
460 8397, 2018.
- 461 [4] Thimo Alldieck, Marcus Magnor, Bharat Lal Bhatnagar,  
462 Christian Theobalt, and Gerard Pons-Moll. Learning to re-  
463 construct people in clothing from a single rgb camera. In  
464 *Proceedings of the IEEE/CVF Conference on Computer Vi-  
465 sion and Pattern Recognition*, pages 1175–1186, 2019. 3
- 466 [5] Thimo Alldieck, Mihai Zanfir, and Cristian Sminchisescu.  
467 Photorealistic monocular 3d reconstruction of humans wear-  
468 ing clothing. In *Proceedings of the IEEE/CVF Conference  
469 on Computer Vision and Pattern Recognition*, pages 1506–  
470 1515, 2022. 2
- 471 [6] Yukang Cao, Guanying Chen, Kai Han, Wenqi Yang, and  
472 Kwan-Yee K Wong. Jiff: Jointly-aligned implicit face func-  
473 tion for high quality single view clothed human reconstruc-  
474 tion. In *Proceedings of the IEEE/CVF Conference on Com-  
475 puter Vision and Pattern Recognition*, pages 2729–2739,  
476 2022. 3
- 477 [7] Yukang Cao, Kai Han, and Kwan-Yee K Wong. Sesdf: Self-  
478 evolved signed distance field for implicit 3d clothed human  
479 reconstruction. In *Proceedings of the IEEE/CVF Conference  
480 on Computer Vision and Pattern Recognition*, pages 4647–  
481 4657, 2023. 3
- 482 [8] Jinnan Chen, Chen Li, Jianfeng Zhang, Hanlin Chen, Buzhen  
483 Huang, and Gim Hee Lee. Generalizable human gaussians  
484 from single-view image. *arXiv preprint arXiv:2406.06050*,  
485 2024. 3
- 486 [9] Enric Corona, Mihai Zanfir, Thimo Alldieck, Eduard  
487 Gabriel Bazavan, Andrei Zanfir, and Cristian Smin-  
488 chisescu. Structured 3d features for reconstructing control-  
489 lable avatars. In *Proceedings of the IEEE/CVF Conference  
490 on Computer Vision and Pattern Recognition*, pages 16954–  
491 16964, 2023. 2, 3
- 492 [10] Florinel-Alin Croitoru, Vlad Hondu, Radu Tudor Ionescu,  
493 and Mubarak Shah. Diffusion models in vision: A survey.  
494 *IEEE Transactions on Pattern Analysis and Machine Intelli-  
495 gence*, 45(9):10850–10869, 2023. 2, 1
- 496 [11] Prafulla Dhariwal and Alexander Nichol. Diffusion models  
497 beat gans on image synthesis. *Advances in neural informa-  
498 tion processing systems*, 34:8780–8794, 2021. 2
- 499 [12] Prafulla Dhariwal and Alexander Nichol. Diffusion models  
500 beat gans on image synthesis. *Advances in neural informa-  
501 tion processing systems*, 34:8780–8794, 2021. 2
- 502 [13] Ben Fei, Jingyi Xu, Rui Zhang, Qingyuan Zhou, Weidong  
503 Yang, and Ying He. 3d gaussian splatting as new era: A  
504 survey. *IEEE Transactions on Visualization and Computer  
505 Graphics*, 2024. 3
- 506 [14] Yao Feng, Vasileios Choutas, Timo Bolkart, Dimitrios  
507 Tzionas, and Michael J Black. Collaborative regression of  
508 expressive bodies using moderation. In *2021 International  
509 Conference on 3D Vision (3DV)*, pages 792–804. IEEE,  
510 2021. 3, 2
- 511 [15] Yao Feng, Weiyang Liu, Timo Bolkart, Jinlong Yang, Marc  
512 Pollefeys, and Michael J Black. Learning disentangled  
513 avatars with hybrid 3d representations. *arXiv preprint  
514 arXiv:2309.06441*, 2023. 3
- 515 [16] Dario Fuoli, Luc Van Gool, and Radu Timofte. Fourier space  
516 losses for efficient perceptual image super-resolution. In  
517 *Proceedings of the IEEE/CVF International Conference on  
518 Computer Vision*, pages 2360–2369, 2021. 5
- 519 [17] Chen Geng, Sida Peng, Zhen Xu, Hujun Bao, and Xiaowei  
520 Zhou. Learning neural volumetric representations of dy-  
521 namic humans in minutes. In *Proceedings of the IEEE/CVF  
522 Conference on Computer Vision and Pattern Recognition*,  
523 pages 8759–8770, 2023. 3
- 524 [18] Chen Guo, Tianjian Jiang, Xu Chen, Jie Song, and Otmar  
525 Hilliges. Vid2avatar: 3d avatar reconstruction from videos  
526 in the wild via self-supervised scene decomposition. In *Pro-  
527 ceedings of the IEEE/CVF Conference on Computer Vision  
528 and Pattern Recognition*, pages 12858–12868, 2023. 3
- 529 [19] Kaiwen Guo, Peter Lincoln, Philip Davidson, Jay Busch,  
530 Xueming Yu, Matt Whalen, Geoff Harvey, Sergio Orts-  
531 Escolano, Rohit Pandey, Jason Dou�arian, et al. The re-  
532 lightables: Volumetric performance capture of humans with  
533 realistic relighting. *ACM Transactions on Graphics (ToG)*,  
534 38(6):1–19, 2019. 1
- 535 [20] Sang-Hun Han, Min-Gyu Park, Ju Hong Yoon, Ju-Mi Kang,  
536 Young-Jae Park, and Hae-Gon Jeon. High-fidelity 3d human  
537 digitization from single 2k resolution images. In *Pro-  
538 ceedings of the IEEE/CVF Conference on Computer Vision and  
539 Pattern Recognition*, pages 12869–12879, 2023. 3
- 540 [21] Tong He, John Collomosse, Hailin Jin, and Stefano Soatto.  
541 Geo-pifu: Geometry and pixel aligned implicit functions for  
542 single-view human reconstruction. *Advances in Neural In-  
543 formation Processing Systems*, 33:9276–9287, 2020. 3
- 544 [22] Tong He, Yuanlu Xu, Shunsuke Saito, Stefano Soatto, and  
545 Tony Tung. Arch++: Animation-ready clothed human recon-  
546 struction revisited. In *Proceedings of the IEEE/CVF interna-  
547 tional conference on computer vision*, pages 11046–11056,  
548 2021. 3
- 549 [23] I Ho, Jie Song, Otmar Hilliges, et al. Sith: Single-view tex-  
550 tured human reconstruction with image-conditioned diffu-  
551 sion. In *Proceedings of the IEEE/CVF Conference on Com-  
552 puter Vision and Pattern Recognition*, pages 538–549, 2024.  
553 2, 3, 7
- 554 [24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising dif-  
555 fusion probabilistic models. *Advances in neural information  
556 processing systems*, 33:6840–6851, 2020. 2, 1
- 557 [25] Liangxiao Hu, Hongwen Zhang, Yuxiang Zhang, Boyao  
558 Zhou, Boning Liu, Shengping Zhang, and Liqiang Nie.

- 559 Gaussianavatar: Towards realistic human avatar modeling  
560 from a single video via animatable 3d gaussians. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
561 and Pattern Recognition, pages 634–644, 2024. 3
- 562 [26] Shoukang Hu, Fangzhou Hong, Liang Pan, Haiyi Mei, Lei  
563 Yang, and Ziwei Liu. Sherf: Generalizable human nerf from  
564 a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9352–9364,  
565 2023. 3
- 566 [27] Shuo Huang, Zongxin Yang, Liangting Li, Yi Yang, and  
567 Jia Jia. Avatarfusion: Zero-shot generation of clothing-decoupled 3d avatars using 2d diffusion. In *Proceedings of the 31st ACM International Conference on Multimedia*,  
568 pages 5734–5745, 2023. 3
- 569 [28] Shuo Huang, Shikun Sun, Zixuan Wang, Xiaoyu Qin, Yan-  
570 min Xiong, Yuan Zhang, Pengfei Wan, Di Zhang, and Jia  
571 Jia. Placiddreamer: Advancing harmony in text-to-3d genera-  
572 tion. *arXiv preprint arXiv:2407.13976*, 2024. 2
- 573 [29] Yangyi Huang, Hongwei Yi, Weiyang Liu, Haofan Wang,  
574 Boxi Wu, Wenxiao Wang, Binbin Lin, Debing Zhang, and  
575 Deng Cai. One-shot implicit animatable avatars with model-  
576 based priors. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8974–8985, 2023. 3
- 577 [30] Yangyi Huang, Hongwei Yi, Yuliang Xiu, Tingting Liao, Ji-  
578 axiang Tang, Deng Cai, and Justus Thies. Tech: Text-guided  
579 reconstruction of lifelike clothed humans. In *2024 Interna-  
580 tional Conference on 3D Vision (3DV)*, pages 1531–1542.  
581 IEEE, 2024. 2, 3
- 582 [31] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and  
583 Tony Tung. Arch: Animatable reconstruction of clothed hu-  
584 mans. In *Proceedings of the IEEE/CVF Conference on Com-  
585 puter Vision and Pattern Recognition*, pages 3093–3102,  
586 2020. 2, 3
- 587 [32] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A  
588 Efros. Image-to-image translation with conditional adver-  
589 sarial networks. In *Proceedings of the IEEE conference on com-  
590 puter vision and pattern recognition*, pages 1125–1134,  
591 2017. 2
- 592 [33] Liming Jiang, Bo Dai, Wayne Wu, and Chen Change Loy.  
593 Focal frequency loss for image reconstruction and synthesis.  
594 In *Proceedings of the IEEE/CVF international conference on  
595 computer vision*, pages 13919–13929, 2021. 5
- 596 [34] Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. In-  
597 stantavatar: Learning avatars from monocular video in 60  
598 seconds. In *Proceedings of the IEEE/CVF Conference on  
599 Computer Vision and Pattern Recognition*, pages 16922–  
600 16932, 2023. 3
- 601 [35] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel,  
602 and Anurag Ranjan. Neuman: Neural human radiance field  
603 from a single video. In *European Conference on Computer  
604 Vision*, pages 402–418. Springer, 2022. 3
- 605 [36] Hanbyul Joo, Tomas Simon, and Yaser Sheikh. Total cap-  
606 ture: A 3d deformation model for tracking faces, hands, and  
607 bodies. In *Proceedings of the IEEE conference on computer  
608 vision and pattern recognition*, pages 8320–8329, 2018. 3
- 609 [37] Angjoo Kanazawa, Michael J Black, David W Jacobs, and  
610 Jitendra Malik. End-to-end recovery of human shape and  
611 pose. In *Proceedings of the IEEE conference on computer  
612 vision and pattern recognition*, pages 7122–7131, 2018. 3
- 613 [38] Byungjun Kim, Patrick Kwon, Kwangho Lee, Myunggi Lee,  
614 Sookwan Han, Daesik Kim, and Hanbyul Joo. Chupa:  
615 Carving 3d clothed humans from skinned shape priors us-  
616 ing 2d diffusion probabilistic models. In *Proceedings of the  
617 IEEE/CVF International Conference on Computer Vision*,  
618 pages 15965–15976, 2023. 2, 3
- 619 [39] Muhammed Kocabas, Nikos Athanasiou, and Michael J  
620 Black. Vibe: Video inference for human body pose and  
621 shape estimation. In *Proceedings of the IEEE/CVF con-  
622 ference on computer vision and pattern recognition*, pages  
623 5253–5263, 2020. 3
- 624 [40] Muhammed Kocabas, Chun-Hao P Huang, Otmar Hilliges,  
625 and Michael J Black. Pare: Part attention regressor for 3d  
626 human body estimation. In *Proceedings of the IEEE/CVF  
627 international conference on computer vision*, pages 11127–  
628 11137, 2021.
- 629 [41] Muhammed Kocabas, Chun-Hao P Huang, Joachim Tesch,  
630 Lea Müller, Otmar Hilliges, and Michael J Black. Spec:  
631 Seeing people in the wild with an estimated camera. In  
632 *Proceedings of the IEEE/CVF International Conference on  
633 Computer Vision*, pages 11035–11045, 2021.
- 634 [42] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and  
635 Kostas Daniilidis. Learning to reconstruct 3d human pose  
636 and shape via model-fitting in the loop. In *Proceedings of  
637 the IEEE/CVF international conference on computer vision*,  
638 pages 2252–2261, 2019. 3
- 639 [43] Verica Lazova, Eldar Insafutdinov, and Gerard Pons-Moll.  
640 360-degree textures of people in clothing from a single im-  
641 age. In *2019 International Conference on 3D Vision (3DV)*,  
642 pages 643–653. IEEE, 2019. 3
- 643 [44] Bo Li, Kaitao Xue, Bin Liu, and Yu-Kun Lai. Bbdm: Image-  
644 to-image translation with brownian bridge diffusion models.  
645 In *Proceedings of the IEEE/CVF conference on computer vi-  
646 sion and pattern Recognition*, pages 1952–1961, 2023. 2, 3,  
647 1
- 648 [45] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang,  
649 and Cewu Lu. Hybrik: A hybrid analytical-neural inverse  
650 kinematics solution for 3d human pose and shape estimation.  
651 In *Proceedings of the IEEE/CVF conference on computer vi-  
652 sion and pattern recognition*, pages 3383–3393, 2021. 3
- 653 [46] Jiefeng Li, Siyuan Bian, Qi Liu, Jia Sheng Tang, Fan Wang,  
654 and Cewu Lu. Niki: Neural inverse kinematics with invert-  
655 ible neural networks for 3d human pose and shape estima-  
656 tion. In *Proceedings of the IEEE/CVF Conference on Com-  
657 puter Vision and Pattern Recognition*, pages 12933–12942,  
658 2023.
- 659 [47] Jiahao Li, Zongxin Yang, Xiaohan Wang, Jianxin Ma, Chang  
660 Zhou, and Yi Yang. Jotr: 3d joint contrastive learning  
661 with transformers for occluded human mesh recovery. In  
662 *Proceedings of the IEEE/CVF International Conference on  
663 Computer Vision*, pages 9110–9121, 2023.
- 664 [48] Zhihao Li, Jianzhuang Liu, Zhensong Zhang, Songcen Xu,  
665 and Youliang Yan. Cliff: Carrying location information  
666 in full frames into human pose and shape estimation. In  
667 *European Conference on Computer Vision*, pages 590–606.  
668 Springer, 2022. 3

- 674 [49] Zhe Li, Zerong Zheng, Lizhen Wang, and Yebin Liu. Animatable gaussians: Learning pose-dependent gaussian maps 675 for high-fidelity human avatar modeling. In *Proceedings of 676 the IEEE/CVF Conference on Computer Vision and Pattern 677 Recognition*, pages 19711–19722, 2024. 3
- 678 [50] Tingting Liao, Xiaomei Zhang, Yuliang Xiu, Hongwei Yi, 679 Xudong Liu, Guo-Jun Qi, Yong Zhang, Xuan Wang, Xiangyu 680 Zhu, and Zhen Lei. High-fidelity clothed avatar 681 reconstruction from a single image. In *Proceedings of 682 the IEEE/CVF Conference on Computer Vision and Pattern 683 Recognition*, pages 8662–8672, 2023. 3
- 684 [51] Tingting Liao, Hongwei Yi, Yuliang Xiu, Jiaxiang Tang, 685 Yangyi Huang, Justus Thies, and Michael J Black. Tada! 686 text to animatable digital avatars. In *2024 International 687 Conference on 3D Vision (3DV)*, pages 1508–1519. IEEE, 2024. 688 3
- 689 [52] Hongyu Liu, Xuan Wang, Ziyu Wan, Yujun Shen, Yibing 690 Song, Jing Liao, and Qifeng Chen. Headartist: Text- 691 conditioned 3d head generation with self score distillation. 692 In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–12, 693 2024. 3
- 694 [53] Stephen Lombardi, Jason Saragih, Tomas Simon, and Yaser 695 Sheikh. Deep appearance models for face rendering. *ACM 696 Transactions on Graphics (ToG)*, 37(4):1–13, 2018. 1
- 697 [54] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard 698 Pons-Moll, and Michael J Black. Smpl: A skinned multi- 699 person linear model. In *Seminal Graphics Papers: Pushing 700 the Boundaries, Volume 2*, pages 851–866. 2023. 3
- 701 [55] Yiwei Ma, Zhekai Lin, Jiayi Ji, Yijun Fan, Xiaoshuai Sun, 702 and Rongrong Ji. X-oscar: A progressive framework for 703 high-quality text-guided 3d animatable avatar generation. 704 *arXiv preprint arXiv:2405.00954*, 2024. 3
- 705 [56] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, 706 Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: 707 Representing scenes as neural radiance fields for view 708 synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 709 3
- 710 [57] Alexander Quinn Nichol and Prafulla Dhariwal. Improved 711 denoising diffusion probabilistic models. In *International 712 conference on machine learning*, pages 8162–8171. PMLR, 713 2021. 2
- 714 [58] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, 715 Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and 716 Michael J Black. Expressive body capture: 3d hands, 717 face, and body from a single image. In *Proceedings of 718 the IEEE/CVF conference on computer vision and pattern 719 recognition*, pages 10975–10985, 2019. 3
- 720 [59] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, 721 Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and 722 Michael J Black. Expressive body capture: 3d hands, 723 face, and body from a single image. In *Proceedings of 724 the IEEE/CVF conference on computer vision and pattern 725 recognition*, pages 10975–10985, 2019. 2, 3, 1
- 726 [60] Ryan Po, Wang Yifan, Vladislav Golyanik, Kfir Aberman, 727 Jonathan T Barron, Amit Bermano, Eric Chan, Tali Dekel, 728 Aleksander Holynski, Angjoo Kanazawa, et al. State of the 729 art on diffusion models for visual computing. In *Computer 730 Graphics Forum*, page e15063. Wiley Online Library, 2024. 731 2
- 732 [61] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv 733 preprint arXiv:2209.14988*, 2022. 2
- 734 [62] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix 735 Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and 736 Aaron Courville. On the spectral bias of neural networks. In 737 *International conference on machine learning*, pages 5301– 738 5310. PMLR, 2019. 5
- 739 [63] Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, 740 and Daniel Cohen-Or. Texture: Text-guided texturing of 3d 741 shapes. In *ACM SIGGRAPH 2023 conference proceedings*, 742 pages 1–11, 2023. 3
- 743 [64] Robin Rombach, Andreas Blattmann, Dominik Lorenz, 744 Patrick Esser, and Björn Ommer. High-resolution image 745 synthesis with latent diffusion models. In *Proceedings of 746 the IEEE/CVF conference on computer vision and pattern 747 recognition*, pages 10684–10695, 2022. 2, 1
- 748 [65] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned 749 implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF international 750 conference on computer vision*, pages 2304–2314, 2019. 2, 3, 7
- 751 [66] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul 752 Joo. Pifuhd: Multi-level pixel-aligned implicit function for 753 high-resolution 3d human digitization. In *Proceedings of 754 the IEEE/CVF conference on computer vision and pattern 755 recognition*, pages 84–93, 2020. 2, 3, 7
- 756 [67] Akash Sengupta, Thiemo Alldieck, Nikos Kolotouros, Enric 757 Corona, Andrei Zanfir, and Cristian Sminchisescu. Diffhuman: Probabilistic photorealistic 3d reconstruction of 758 humans. In *Proceedings of the IEEE/CVF Conference on 759 Computer Vision and Pattern Recognition*, pages 1439–1449, 760 2024. 3
- 761 [68] Xiaolong Shen, Zongxin Yang, Xiaohan Wang, Jianxin Ma, 762 Chang Zhou, and Yi Yang. Global-to-local modeling for 763 video-based 3d human pose and shape estimation. In *Pro- 764 ceedings of the IEEE/CVF Conference on Computer Vision 765 and Pattern Recognition*, pages 8887–8896, 2023. 3
- 766 [69] J Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, 767 Jiajun Wu, and Gordon Wetzstein. 3d neural field generation 768 using triplane diffusion. In *Proceedings of the IEEE/CVF 769 Conference on Computer Vision and Pattern Recognition*, 770 pages 20875–20886, 2023. 2
- 771 [70] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, 772 and Surya Ganguli. Deep unsupervised learning using 773 nonequilibrium thermodynamics. In *International confer- 774 ence on machine learning*, pages 2256–2265. PMLR, 2015. 775 2, 1
- 776 [71] Xiaokun Sun, Zhenyu Zhang, Ying Tai, Qian Wang, Hao 777 Tang, Zili Yi, and Jian Yang. Barbie: Text to barbie-style 778 3d avatars. *arXiv preprint arXiv:2408.09126*, 2024. 3
- 779 [72] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang 780 Zeng. Dreamgaussian: Generative gaussian splatting for effi- 781 cient 3d content creation. *arXiv preprint arXiv:2309.16653*, 782 2023. 3

- 788 [73] Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany,  
789 Sanja Fidler, Karsten Kreis, et al. Lion: Latent point dif-  
790 fusion models for 3d shape generation. *Advances in Neural*  
791 *Information Processing Systems*, 35:10021–10039, 2022. 2
- 792 [74] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan  
793 Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and  
794 diverse text-to-3d generation with variational score distilla-  
795 tion. *Advances in Neural Information Processing Systems*,  
796 36, 2024. 2, 3
- 797 [75] Donglai Xiang, Fabian Prada, Chenglei Wu, and Jessica  
798 Hodgins. Monoclothcap: Towards temporally coherent  
799 clothing capture from monocular rgb video. In *2020 Inter-  
800 national Conference on 3D Vision (3DV)*, pages 322–332.  
801 IEEE, 2020. 3
- 802 [76] Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J  
803 Black. Icon: Implicit clothed humans obtained from nor-  
804 mals. In *2022 IEEE/CVF Conference on Computer Vi-  
805 sion and Pattern Recognition (CVPR)*, pages 13286–13296.  
806 IEEE, 2022. 2, 3, 7
- 807 [77] Yuliang Xiu, Jinlong Yang, Xu Cao, Dimitrios Tzionas, and  
808 Michael J Black. Econ: Explicit clothed humans optimized  
809 via normal integration. In *Proceedings of the IEEE/CVF con-  
810 ference on computer vision and pattern recognition*, pages  
811 512–523, 2023. 2, 3, 6, 7, 1
- 812 [78] Hongyi Xu, Eduard Gabriel Bazavan, Andrei Zanfir,  
813 William T Freeman, Rahul Sukthankar, and Cristian Smin-  
814 chisescu. Ghum & ghuml: Generative 3d human shape and  
815 articulated pose models. In *Proceedings of the IEEE/CVF  
816 Conference on Computer Vision and Pattern Recognition*,  
817 pages 6184–6193, 2020. 3
- 818 [79] Jiale Xu, Xintao Wang, Weihao Cheng, Yan-Pei Cao, Ying  
819 Shan, Xiaohu Qie, and Shenghua Gao. Dream3d: Zero-shot  
820 text-to-3d synthesis using 3d shape prior and text-to-image  
821 diffusion models. In *Proceedings of the IEEE/CVF Con-  
822 ference on Computer Vision and Pattern Recognition*, pages  
823 20908–20918, 2023. 2
- 824 [80] Yuanyou Xu, Zongxin Yang, and Yi Yang. Seeavatar: Photo-  
825 realistic text-to-3d avatar generation with constrained geo-  
826 metry and appearance. *arXiv preprint arXiv:2312.08889*,  
827 2023. 3
- 828 [81] Zhiqin John Xu. Understanding training and generaliza-  
829 tion in deep learning by fourier analysis. *arXiv preprint  
830 arXiv:1808.04295*, 2018. 5
- 831 [82] Zhi-Qin John Xu, Yaoyu Zhang, and Yanyang Xiao. Train-  
832 ing behavior of deep neural network in frequency domain. In  
833 *Neural Information Processing: 26th International Confer-  
834 ence, ICONIP 2019, Sydney, NSW, Australia, December 12–  
835 15, 2019, Proceedings, Part I 26*, pages 264–274. Springer,  
836 2019. 5
- 837 [83] Xuetong Yang, Yihao Luo, Yuliang Xiu, Wei Wang, Hao Xu,  
838 and Zhaoxin Fan. D-if: Uncertainty-aware human digitiz-  
839 ation via implicit distribution field. In *Proceedings of the  
840 IEEE/CVF International Conference on Computer Vision*,  
841 pages 9122–9132, 2023. 2, 3, 7
- 842 [84] Taoran Yi, Jiemin Fang, Junjie Wang, Guanjun Wu, Lingxi  
843 Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and Xinggang  
844 Wang. Gaussiandreamer: Fast generation from text to 3d  
gaussians by bridging 2d and 3d diffusion models. In *Pro-  
845 ceedings of the IEEE/CVF Conference on Computer Vision  
846 and Pattern Recognition*, pages 6796–6807, 2024. 3
- [85] Ilya Zakharkin, Kirill Mazur, Artur Grigorev, and Victor  
Lempitsky. Point-based modeling of human clothing. In  
*Proceedings of the IEEE/CVF International Conference on  
Computer Vision*, pages 14718–14727, 2021. 3
- [86] 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 interna-  
tional conference on computer vision*, pages 11446–11456,  
2021. 3
- [87] Hongwen Zhang, Yating Tian, Yuxiang Zhang, Mengcheng  
Li, Liang An, Zhenan Sun, and Yebin Liu. Pymaf-x: To-  
wards well-aligned full-body model regression from mono-  
cular images. *IEEE Transactions on Pattern Analysis and Ma-  
chine Intelligence*, 45(10):12287–12303, 2023. 3
- [88] Jingbo Zhang, Xiaoyu Li, Qi Zhang, Yanpei Cao, Ying Shan,  
and Jing Liao. Humanref: Single image to 3d human gen-  
eration via reference-guided diffusion. In *Proceedings of  
the IEEE/CVF Conference on Computer Vision and Pattern  
Recognition*, pages 1844–1854, 2024. 3
- [89] Zechuan Zhang, Li Sun, Zongxin Yang, Ling Chen, and  
Yi Yang. Global-correlated 3d-decoupling transformer for  
clothed avatar reconstruction. *Advances in Neural Infor-  
mation Processing Systems*, 36, 2024. 2, 3, 7
- [90] Zechuan Zhang, Zongxin Yang, and Yi Yang. Sifu:  
Side-view conditioned implicit function for real-world us-  
able clothed human reconstruction. In *Proceedings of the  
IEEE/CVF Conference on Computer Vision and Pattern  
Recognition*, pages 9936–9947, 2024. 2, 3, 7
- [91] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai.  
Pamir: Parametric model-conditioned implicit representa-  
tion for image-based human reconstruction. *IEEE transac-  
tions on pattern analysis and machine intelligence*, 44(6):  
3170–3184, 2021. 2, 3, 7
- [92] Yunliang Zhuang, Zhuoran Zheng, and Chen Lyu. Dpfnet: A  
dual-branch dilated network with phase-aware fourier con-  
volution for low-light image enhancement. *arXiv preprint  
arXiv:2209.07937*, 2022. 5