

# GENERALIZABLE HUMAN GAUSSIANS FROM SINGLE-VIEW IMAGE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In this work, we tackle the task of learning 3D human Gaussians from a single image, focusing on recovering detailed appearance and geometry including unobserved regions. We introduce a single-view generalizable Human Gaussian Model (HGM), which employs a novel generate-then-refine pipeline with the guidance from human body prior and diffusion prior. Our approach uses a ControlNet to refine rendered back-view images from coarse predicted human Gaussians, then uses the refined image along with the input image to reconstruct refined human Gaussians. To mitigate the potential generation of unrealistic human poses and shapes, we incorporate human priors from the SMPL-X model as a dual branch, propagating image features from the SMPL-X volume to the image Gaussians using sparse convolution and attention mechanisms. Given that the initial SMPL-X estimation might be inaccurate, we gradually refine it with our HGM model. We validate our approach on several publicly available datasets. Our method surpasses previous methods in both novel view synthesis and surface reconstruction. Our approach also exhibits strong generalization for cross-dataset evaluation and in-the-wild images.

## 1 INTRODUCTION

Automatic 3D human reconstruction from single image is crucial in augmented and virtual reality (AR/VR), game industry, filmmaking, *etc.* Previous works rely on strong 3D supervision such as the signed distance value or occupancy (Saito et al., 2019; 2020; Zhang et al., 2023c; Xiu et al., 2022; 2023; Zhang et al., 2024; Ho et al., 2024) and focus on surface reconstruction, neglecting novel view synthesis quality, resulting in smoothed and blurred textures. With the development of neural radiance fields (Mildenhall et al., 2020), novel view rendering quality has been greatly improved with for human appearance modelling (Hu et al., 2023; Kwon et al., 2021; Gao et al., 2022). However, due to the ill-posed nature of single view reconstruction, the back and side views are always blurry and lack details without additional prior. Furthermore, these methods needs large amounts of query points sampled for volume rendering, which hinders practical real-time application in the industries. Some other methods optimize underlying appearance and geometry from scratch by introducing score distillation sampling during the optimization (Tang et al. (2024b); Cao et al. (2024)). Although effective, these methods still suffer from long time optimization and over-saturation problems. Recent 3D Gaussians generation methods (Tang et al., 2024a; Yinghao et al., 2024) combine multi-view diffusion models (Liu et al., 2023; Wang & Shi, 2023; Shi et al., 2023; Li et al., 2024; Xu et al., 2023) with generalizable multi-view Gaussians prediction models to generate 3D Gaussians with high quality and efficiency. We aim to extend this on human reconstruction. However, directly employ such methods to human reconstruction with complex texture and poses gives unsatisfying results due to: (1) Inconsistency across multiple views: The multi-view images generated from diffusion model lacks consistency in appearance and pose across different viewpoints. This inconsistency stems from the inherent complexity of human body structure and movement, which leads to low-quality reconstruction results. (2) Quality loss in front view reconstruction: multi-view diffusion process involves down-sampling and changing the original input image. This step results in significant quality degradation when reconstructing the front view image, compromising the fidelity to the original input. (3) Estimating SMPL-X parameters from single view input is ill-posed, directly applying initially estimated SMPL-X can lead to bending legs and wrong elevation issues in previous method (Xiu et al. (2022; 2023); Zhang et al. (2024); Ho et al. (2024)).



Figure 1: Our method reconstructs detailed and geometrically consistent human Gaussian models from single view images, including loosing clothes, challenging pose and in-the-wild images.

To address the above-mentioned problems, we introduce a novel Human Gaussians Model (HGM), which supports fast and high quality rendering from single view input, and generalize well to loosing clothes, challenging poses and in-the-wild images as shown in Fig. 1. We do not use multi-view diffusion models due to the multiview inconsistency and resolution degradation problem. Instead, we propose a coarse-to-fine framework, where the diffusion model is adapted to refine back-view images rendered from our predicted coarse human Gaussians. In this way, we can keep the resolution and content of the original input image for high-fidelity reconstruction. In order to model the complex structure of a human, we inject the human prior into the Gaussian prediction process. Specifically, our model consists of two branches: 1) The first branch is a UNet to directly predict Gaussians from the input image, as inspired by image splatter (Szymanowicz et al., 2024). 2) The second branch uses learnable tokens attached to SMPL-X vertices for structural feature extraction with attention layers, and then combined with UNet features with SparseConvGraham et al. (2018) and a transformer for Gaussian enhancement. Recognizing the inaccurate estimate of SMPL-X from the pre-trained model, during inference, we iteratively refine the initial SMPL-X parameters with our HGM pre-trained with ground truth SMPL-X. Given the loss of details of the back view by directly predicting the Gaussians from a single view, we further apply a ControlNet to refine the back view with the control signal from the back-view image rendered from the coarse stage. We then input the original front view and refined back view images to our HGM model to get the final refined Gaussians. Meshes can be extracted from densely rendered depth map and TSDF fusion. Our model can be trained with only posed multiview images without 3D supervision and generalizes well to untrained datasets and in the wild images.

In summary, our contributions are:

- We introduce a generate-then-refine pipeline for single view human Gaussian reconstruction that leverages diffusion priors for back view refinement, avoiding the multi-view inconsistencies commonly observed in multiview diffusion models.
- Our proposed dual-branch reconstruction pipeline incorporates human priors by attaching learnable tokens to the SMPL-X vertices for structural feature extraction. We then fuse these features from the SMPL-X branch with the U-Net branch using Sparse Convolution and transformer.
- To address potential inaccuracies in initial SMPL-X estimations, we employ our Human Gaussian Model (HGM) to iteratively refine the estimated SMPL-X parameters, resulting in better alignment.
- Through extensive experimentation, we demonstrate the efficacy of our method in both novel view synthesis and 3D reconstruction tasks. Our approach consistently achieves state-of-the-art performance on various metrics and benchmarks.

## 2 RELATED WORKS

**Single-view Human Reconstruction.** PIFu (Saito et al., 2019), PIFuHD (Saito et al., 2020), PaMIR (Zheng et al., 2021), and GTA (Zhang et al., 2023c) are capable of inferring full textures from a single image. Techniques such as PHORHUM (Alldieck et al., 2022) and S3F (Corona et al., 2023) go further by separating albedo and global illumination. However, these methods lack information

108 from other views or prior knowledge, such as diffusion models, often resulting in unsatisfactory  
 109 textures. TeCH (Huang et al., 2024) utilizes diffusion-based models to visualize unseen areas, pro-  
 110 ducing realistic results. However, it requires time-intensive optimization per subject and is depen-  
 111 dent on accurate SMPL-X. The emergence of Neural Radiance Fields (NeRF) has led to methods  
 112 (Hu et al., 2023; Huang et al., 2023; Gao et al., 2022; Kwon et al., 2021) using videos or multi-  
 113 view images to optimize NeRF for the capture of human forms. Recent advancements like SHERF  
 114 (Hu et al., 2023) and ELICIT (Huang et al., 2023) aim to generate human NeRFs from single im-  
 115 ages. Although NeRF-based approaches are effective in creating high-quality images from various  
 116 perspectives, they often struggle with detailed 3D mesh generation from single images and require  
 117 extensive optimization time. More recently, SiTH (Ho et al., 2024) proposes to combine a back-view  
 118 hallucination model with an SDF-based mesh reconstruction model. Similarly, SIFU (Zhang et al.,  
 119 2024) employs a text-to-image diffusion-based prior to generating consistent textures for invisible  
 120 views. However, these methods require 3D annotations such as the SDF of the meshes and texture  
 121 maps as strong supervision and still fail to generate renderings with high fidelity due to the limited  
 122 3D training data and representation. In addition, these methods suffer from SMPL estimation errors,  
 123 leading to bending legs and wrong elevation of the reconstructed 3D humans. Compared to these  
 124 methods, our approach can be trained solely on multi-view images and achieves much better novel  
 view synthesis quality.

125 **Human Gaussians.** 3D Gaussians (Kerbl et al., 2023) and differentiable splatting (Szymanowicz  
 126 et al., 2024) have gained broad popularity due to their efficiency in reconstructing high-fidelity 3D  
 127 scenes from posed images using only a moderate number of 3D Gaussians. This representation has  
 128 been quickly adopted for various applications, including image or text-conditioned 3D generation  
 129 and avatar reconstruction. Among these methods, GauHuman and HUGS (Hu & Liu, 2024; Kocabas  
 130 et al., 2024) are the first to propose optimizing human Gaussians from monocular human videos.  
 131 However, they are not applicable to single static human images. GPS-Gaussian (Zheng et al., 2024)  
 132 propose a generalizable multi-view human Gaussian model with high quality rendering; however, it  
 133 needs dense views 16 or 8, which cannot be directly applied to single-view human images. Our hu-  
 134 man model achieves strong generalization in generating human Gaussians from single-view images,  
 135 complementing concurrent work such as Pan et al. (2024).

136 **Generalizable Gaussians with Multi-view Diffusion.** The Large Reconstruction Model (LRM)  
 137 (Hong et al., 2024) scales up both the model and the dataset to predict a neural radiance field (NeRF)  
 138 from single-view images. Although LRM is primarily a reconstruction model, it can be combined  
 139 with Diffusion Models (DMs) to achieve text-to-3D and image-to-3D generation as demonstrated  
 140 by extensions such as Zero123 (Liu et al., 2023), Image Dream (Wang & Shi, 2023) Instant3D (Li  
 141 et al., 2024) and DMV3D (Xu et al., 2023). Our method also builds on a strong reconstruction model  
 142 and uses pre-trained 2D DMs to provide input images missing information in a feedforward manner.  
 143 Some concurrent works, such as LGM (Tang et al., 2024a), AGG (Xu et al., 2024), and Splat-  
 144 ter Image (Szymanowicz et al., 2024), also utilize 3D Gaussians in a feed-forward model. LGM (Tang  
 145 et al., 2024a) combines novel view generation diffusion models with generalizable Gaussians in a  
 146 feedforward manner, while GRM (Yinghao et al., 2024) replaces the U-Net architecture with a pure-  
 147 transformer one and scales up to large resolution. However, these methods face two main challenges  
 148 when using pre-trained diffusion models. Firstly, the generated input view image becomes blurry  
 149 compared to the original input, which affects the subsequent generalizable Gaussian model. Sec-  
 150 ondly, diffusion models can introduce multiview inconsistency, especially for human images with  
 151 different poses, making direct adaptation unfeasible. We solve these problems by using Control-  
 152 Net as the refinement tools without damaging the input image quality or introducing multi-view  
 153 inconsistency.

## 154 3 OUR METHOD

### 155 3.1 PRELIMINARIES

156 **3D Gaussian Splatting (3DGS).** Introduced by (Kerbl et al., 2023), 3D Gaussian splatting repre-  
 157 sents 3D assets or scenes using a collection of 3D Gaussians. Each Gaussian is characterized by its  
 158 center  $x \in \mathbb{R}^3$ , scaling factor  $s \in \mathbb{R}^3$ , rotation  $r \in \mathbb{R}^3$ , opacity  $\alpha \in \mathbb{R}$ , and color features  $c \in \mathbb{R}^c$ .  
 159 View-dependent effects can be modeled with spherical harmonics. 3D scenes can be explicitly rep-  
 160 resented by a set of Gaussians  $G = \{G_i\}$ , where  $G_i = \{x_i, s_i, r_i, \alpha_i, c_i\}$  represents the attributes for  
 161

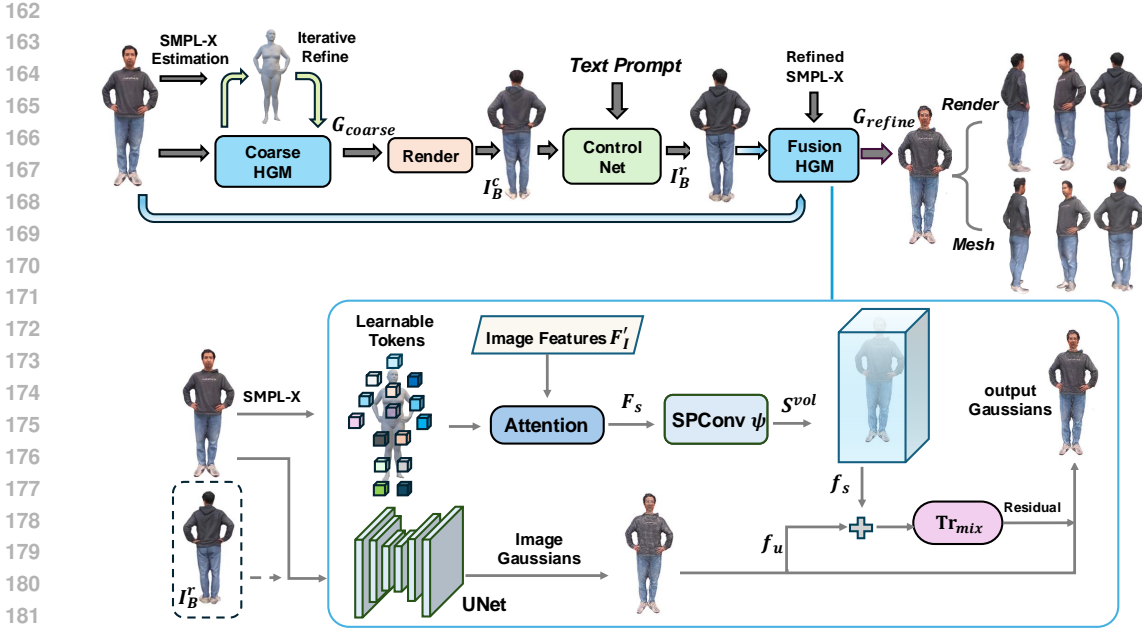


Figure 2: Our framework and HGM model. (Top) Our framework consists of three steps: 1) Coarse Gaussians prediction with iterative SMPL-X refinement. 2) Back view refinement with ControlNet. 3) Two view reconstruction to get the refined  $G_{refine}$ . (Bottom) Our HGM model consists of two branches: Image Gaussians prediction by UNet and adding additional structural features extracted from SMPL-X branch.  $f_{smpl}$  are sampled by the Gaussian centers from the SMPL-X volume  $S^{vol}$  and fused with  $f_u$  to the fusion transformer  $Tr_{mix}$  to obtain the Gaussian output.

the  $i$ -th Gaussian. Compared with NeRF (Mildenhall et al., 2020), 3DGS performs fast rendering by first projecting Gaussians onto the image plane as 2D Gaussians and performing alpha-blending for each pixel in front-to-back depth order. Building on this, Image Splatter (Szymanowicz et al., 2024) proposes predicting Gaussians from a single image through image-to-image translation. Specifically, each pixel is converted to a Gaussian with corresponding attributes, supervised by multi-view images. Our model builds on this representation by directly predicting XYZ coordinates from the image instead of the depth.

### 3.2 OVERVIEW

Fig. 2 shows an overview of our framework. Given a single input human image  $I$ , our aim is to predict the corresponding human Gaussians, which can be further rendered for novel view synthesis and mesh extraction. As shown in the upper part of Fig. 2, our proposed method consists of three parts: 1) **Coarse Gaussians prediction with SMPL-X refinement** (cf. Sec. 3.3) 2) **Back-view refinement with ControlNet**(cf. Sec. 3.4). 3) **Two-view reconstruction** (cf. Sec. 3.5).

### 3.3 COARSE GAUSSIANS PREDICTION WITH SMPL-X REFINEMENT

#### 3.3.1 OUR HGM MODEL

The lower part of Fig. 2 shows our proposed Human Gaussian Model (HGM). The direct prediction of Gaussians from image pixels with UNet (Szymanowicz et al., 2024; Tang et al., 2024a) lacks human shape and pose prior, thus leading to unsatisfactory results. We therefore introduce a dual branch that utilizes SMPL-X to enforce human shape and pose prior into the Gaussian prediction process. Specifically, for UNet branch, the collection of the RGB value and ray embedding for each pixel are concatenated into a 9-channel feature map as the input  $F_I = \{c_i, o_i \times d_i, d_i \mid i = 1, 2, \dots, N\}$ . Our HGM model predicts Gaussians from the U-Net as:

$$G_u = \text{UNet}(F_I), \tag{1}$$

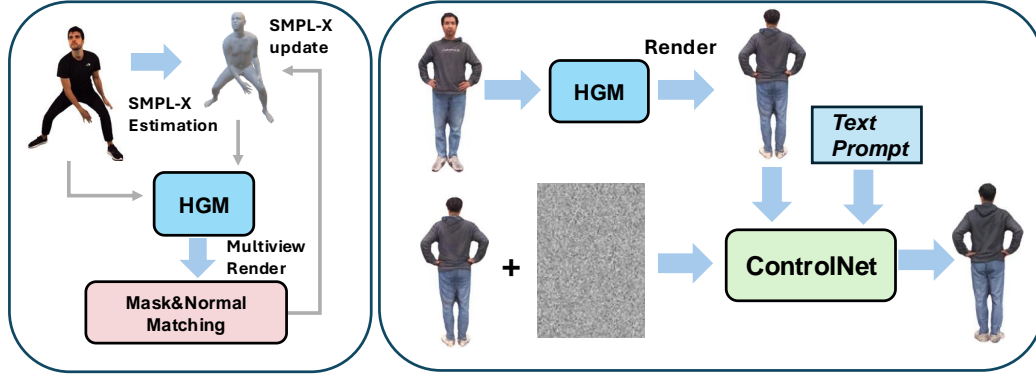


Figure 3: Left: Our SMPL-X refinement pipeline. Right: Our back-view refinement ControlNet.

We call it as Image Gaussians. For the SMPL-X branch, we attach learnable tokens to each of the SMPL-X vertices and extract the patch features of the image, denoted  $F'_i$ . We use cross-attention between these learnable tokens and the patch features to obtain  $F_S$ . This approach takes advantage of the fact that SMPL-X vertices are defined in semantically similar areas across different identities. Consequently, the learned tokens can memorize the mapping from the training dataset to unseen identities during inference, effectively providing structural human priors. This mechanism enables our model to capture and utilize semantic consistent features across diverse identities, enhancing its ability to generalize to new subjects. We apply SparseConvNet (Graham et al., 2018) (SPConv  $\Psi$ ) to propagate the SMPL-X features to the predefined whole bounding box, and we denote this feature as the SMPL-X volume feature:  $S^{vol} = \Psi(F_S)$ , where  $F_S$  are the SMPL-X features. The volume feature reconstructed from the SMPL-X vertex feature provides geometric cues of the target human body. The centers of the Image Gaussians from  $G_u$  are then used to sample the propagated SMPL-X volume features, denoted as  $f_s = S^{vol}(C_u)$ , where  $C_u$  are the centers of  $G_u$ . SMPL-X  $f_s$  features are then concatenated with the features of UNet  $f_u$  for each Gaussian. This concatenated feature is fed into a transformer  $\text{Tr}_{mix}$  to obtain the coarse Gaussians shown in Eq. 2. Specifically, we predict the xyz coordinates residuals for the Image Gaussians and all the other updated Gaussian features.  $\text{Tr}_{mix}$  is a transformer that contains multiple self-attention blocks among Gaussians to ensure that each Gaussian is aware of the other.

$$G_{coarse} = \text{Tr}_{mix}([f_u, f_s]). \quad (2)$$

### 3.3.2 SMPL-X REFINEMENT

Given initial estimated SMPL-X is not accurate, we leverage our pre-trained HGM to iteratively refine the SMPL-X parameters based on the mask and normal matching as shown in Fig. 3 left part. Specifically, we use the initial estimated SMPL-X from pre-trained SMPL-X estimator to reconstruct the initial Gaussians and rendered the side-view masks and normals. We minimize the mask and normal difference between the Gaussian and SMPL-X renderings and back-propagate the loss to SMPL-X parameters. Then we iteratively input the updated SMPL-X to our HGM model, so the Gaussian is also updated to give more accurate masks. For normal matching, we use pre-trained normal estimator from Xiu et al. (2022) which also needs SMPL-X as input and iteratively update the SMPL-X parameters. Specifically, we compute the  $\mathcal{L}_{SMPL-X}$  as shown in 3 and back propagate it to SMPL-X parameters.  $\mathcal{L}_{normal}$ ,  $\mathcal{L}_{front}$  and  $\mathcal{L}_{side}$  are  $\mathcal{L}_1$  losses between the rendered masks of SMPL-X model and coarse Gaussians by HGM.  $\mathcal{L}_{normal}$  is the  $\mathcal{L}_1$  loss between the rendered SMPLX normal and the pre-trained normal estimator (Xiu et al., 2022). Note that both of the normal and side mask supervision are updated as the input of the HGM and normal estimator also contains the updated SMPL-X in each iteration. The whole process takes 100 iterations. We provide more analysis on our SMPL-X refinement in the appendix.

$$\mathcal{L}_{SMPL-X} = \lambda_{front} \mathcal{L}_{m_{front}} + \lambda_{side} \mathcal{L}_{m_{side}} + \lambda_n \mathcal{L}_{normal} \quad (3)$$

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323



Figure 4: Our back-view refinement can generate more realistic back-view images, compared with back-view hallucination of SiTH Ho et al. (2024).

### 3.4 BACK-VIEW REFINEMENT

Back-view hallucination poses a significant challenge in single-view human reconstruction. As shown in Fig. 4, directly using diffusion models to generate a back-view image can result in incorrect perspective projection with the front-view image as well as unrealistic texture by (Ho et al., 2024). The reason is that their diffusion is conditioned on the back-view mask, and thus can only be applied for orthogonal projection where back-view mask can be directly flipped with the front-view image. However, the back-view mask is not available during our inference stage since our prediction is based on perspective projection. To address this issue, we adopt a generate-then-refine strategy that leverages the diffusion prior to produce a perspective-fitted and realistic back-view image that is suitable for the subsequent two-view fusion stage. We train a ControlNet (Zhang et al., 2023b) to enhances realistic details based on our coarse results as shown in Fig. 3 right part. Specifically, we generate coarse back-view rendering by our HGM for the training dataset and only train the ControlNet and keep the base Stable Diffusion model as fixed.

$$\mathcal{L}_{CN} = \mathbb{E}_{\mathbf{z}_0, t, \mathbf{c}_t, \epsilon_t \sim \mathcal{N}(0,1)} \left[ \|\epsilon - \epsilon_\theta(\mathbf{z}_t, t, \mathbf{c}_t, \mathbf{y})\|_2^2 \right] \quad (4)$$

where  $y$  is the text prompt. We set it as ‘Best quality’ and the negtive prompt as ‘blur, bad anatomy, bad hands, cropped, worst quality’ during inference,  $\mathbf{c}_t$  is the coarse back-view image from our HGM rendering  $I_B^c$ , which is the ControlNet condition. We carefully design the reversing process by adding small amount of noise to the VAE encoded latent of  $I_B^c$  to keep the original content as much as possible. The sampling process takes around 2 seconds. In Fig. 4, we showcase our generated and refined back-view results, comparing them with the back-view hallucination diffusion network in SiTH (Zhang et al., 2024). Our results maintain high resolution and generate details such as the hair for the first woman and the wrinkles in the clothes, whereas SiTH (Zhang et al., 2024) produces artifacts and unrealistic hallucination results and are also not perspectiveally fitted to the input image.

### 3.5 TWO-VIEW RECONSTRUCTION

We combine the refined perspective-fitted back-view image  $I_B^r$  with the front-view image  $I$  and input them into the fusion HGM model to get:

$$G_{refine} = \text{HGM}(I, I_B^r). \quad (5)$$

Specifically, our fusion HGM model retains the design of the coarse HGM model with the additional refined back-view image as the input. The coarse HGM and fusion HGM models are trained separately with ground-truth one view and two views as input, as well as ground-truth SMPL-X. The objective function for HGM training includes  $\mathcal{L}_2$  color loss,  $\mathcal{L}_{rgb}$ , VGG-based LPIPS perceptual loss,  $\mathcal{L}_{lpiips}$  (Zhang et al., 2018), and  $\mathcal{L}_2$  background mask loss  $\mathcal{L}_{bg}$  with ground truth masks. Each of these losses has corresponding weights that are treated as hyperparameters:

$$\mathcal{L}_{HGM} = \lambda_{rgb} \mathcal{L}_{rgb} + \lambda_{lpiips} \mathcal{L}_{lpiips} + \lambda_{bg} \mathcal{L}_{bg}, \quad (6)$$

where  $\lambda_{rgb} = \lambda_{lpiips} = \lambda_{bg}=1.0$ .  $\mathcal{L}_{HGM}$  is applied for coarse HGM and fusion HGM.

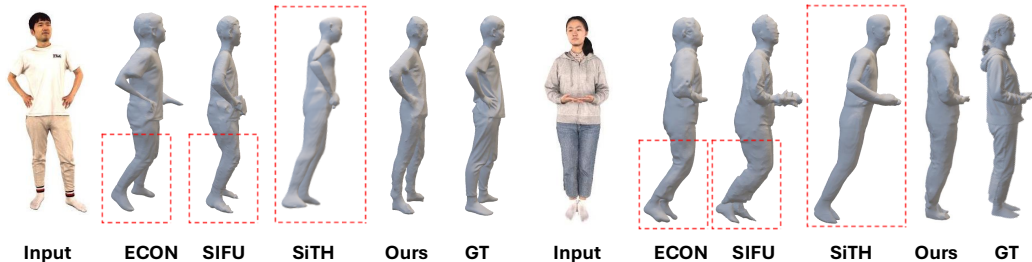


Figure 5: Leveraging our HGM model, SMPL-X parameters are iteratively refined to mitigate the issue of blended legs commonly seen in other approaches.

### 3.6 IMPLEMENTATION

Our model is trained on 4 NVIDIA RTX A6000 with batch size of 4 for 20 hours. Our input image size is 512x512, and the number of Gaussians for each view is 256x256, totaling 65,536 Gaussians per view. For SMPL-X estimation, we use PIXIE(Feng et al., 2021). Network structures and more implementation details are in the appendix.

## 4 EXPERIMENTS

We conduct experiments on the publicly available 3D human datasets THuman2.0 (Yu et al., 2021), CustomHumans(Ho et al., 2023) and HuMMan (Cai et al., 2022). Our method is compared with state-of-the-art (SOTA) methods in both novel view synthesis and 3D mesh reconstruction. We train our HGM on 500 human scans from THuman2.0 dataset following Zhang et al. (2024). We render the images with resolution of 512x512 and using weak perspective camera on 12 fixed cameras evenly distributed with the azimuths from 0 to 360 degree. During evaluation, all the methods are tested without the ground truth SMPL-X. We follow the train and test list from SIFU (Zhang et al., 2024) and SHERF (Hu et al., 2023) to evaluate our method on THuman2.0 and HuMMan dataset. For CustomHumans dataset we use 45 scans for cross-dataset evaluation containing losing clothes and challenging poses. For novel view synthesis, We use PSNR, SSIM, LPIPS as evaluation metrics. For 3D reconstruction, we use commonly used Chamfer Distance (CD), Point-to-Surface Distance (P2S), and Normal Consistency as the evaluation metrics.

### 4.1 NOVEL VIEW SYNTHESIS

Table 1: Novel view synthesis comparison with SOTA methods.

Method	THuman2.0			CustomHumans		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
GTA Zhang et al. (2023c)	19.09	0.882	0.113	19.59	0.887	0.125
SiTH Ho et al. (2024)	17.12	0.843	0.155	18.09	0.856	0.144
LGM Tang et al. (2024a)	18.34	0.856	0.134	19.87	0.877	0.132
SV3D Voleti et al. (2024)	19.11	0.892	0.117	20.86	0.902	0.112
SIFU Zhang et al. (2024)	22.10	0.924	0.0794	20.83	0.898	0.117
<b>Ours</b>	<b>23.54</b>	<b>0.938</b>	<b>0.0524</b>	<b>23.84</b>	<b>0.944</b>	<b>0.0514</b>

For novel view synthesis, we compare our method with mesh SOTA human reconstruction methods GTA Zhang et al. (2023c), ECON Xiu et al. (2023), SIFU Zhang et al. (2024) and SiTH (Ho et al., 2024), as well as multiview diffusion reconstruction method LGM(finetuned with the same training data) Tang et al. (2024a) and video diffusion method SV3D Voleti et al. (2024) on THuman2.0 Yu et al. (2021) and CustomHuman Ho et al. (2023). We also compare our method with state-of-the-art HumanNeRF methods: SHERF (Hu et al., 2023), MPS-NeRF (Gao et al., 2022), and NHP (Kwon et al., 2021) on the HuMMan dataset (Cai et al., 2022).

As shown in the Tab. 1 and Tab. 2, our method significantly surpasses state-of-the-art single-view human reconstruction methods in all evaluation metrics for the three datasets. As shown in Fig. 6, LGM Tang et al. (2024a) generates incorrect blue color and inconsistent content. SiTH Ho et al. (2024) fails to model loose clothes due to the high dependency of the SMPL-X model. Side views are blurry and unrealistic in SIFU’s results. SV3D Voleti et al. (2024) generates strange colors and wrong human pose. Compared with these methods, ours generates more realistic and consistent rendering especially for the unseen regions such as clothes wrinkles and hair that are well-fitted to the front views and more robust to initial SMPL-X estimation errors thanks to our iterative refinement. We provide rendering videos in 360 degree comparison with other methods in the Appendix. We also provide a visual comparison with the SOTA NeRF-based method SHERF Hu et al. (2023) on the HuMMan dataset in the appendix.



Figure 6: Novel view synthesis comparison with other approaches on THuman2.0 and CustomHumans dataset. **The details are highlighted in the red boxes.**

#### 4.2 3D RECONSTRUCTION

For mesh reconstruction we extract the 3D mesh by densely rendering the depth map with Gaussian render and using TSDFusion to extract the surface followed by a fast optimization based on the normal map obtained in section 3.3.2. We compare our results with SOTA human surface reconstruction methods GTA (Zhang et al., 2023c), ECON (Xiu et al., 2023), SIFU (Zhang et al., 2024) and SiTH



Table 2: Novel view synthesis comparison with SOTA HumanNeRF methods on HuMMan.

Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
NHP (Kwon et al., 2021)	18.99	0.845	0.182
MPS-NeRF (Gao et al., 2022)	17.44	0.824	0.193
SHERF (Hu et al., 2023)	20.83	0.891	0.125
Ours	<b>23.86</b>	<b>0.952</b>	<b>0.0591</b>

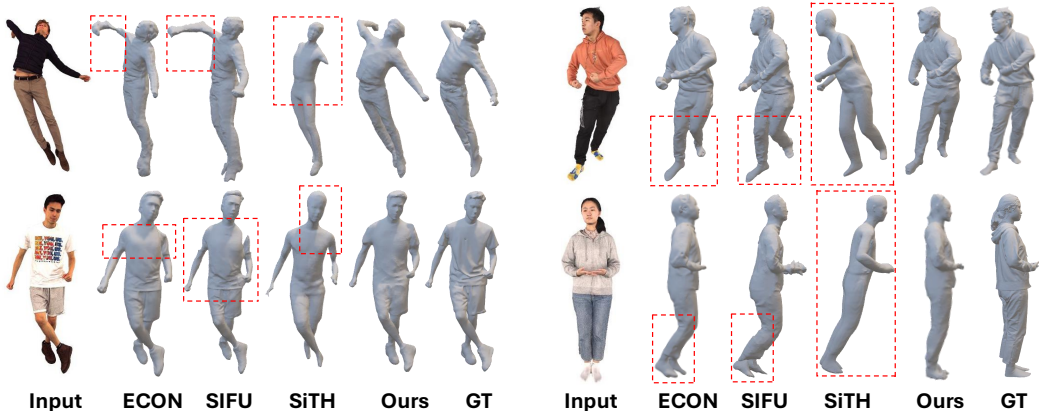


Figure 7: 3D reconstruction visualization compared with SOTA methods. **Details are highlighted in the red boxes.**

Ho et al. (2024). Note that our method do not use the 3D ground truth for supervision, but can also achieve best performance compared with all those fully supervised methods. Thanks to our iterative SMPL-X refinement. Our method alleviates commonly occurring problems of bent legs and incorrect postures found in previous methods, as shown in Fig.5. Our HGM model reconstructs more accurate geometry with the prior learned from our dual branch Gaussian reconstruction model as well as our innovative SMPL-X refinement. We provide a qualitative 3D reconstruction comparison in Fig.7. As shown in the figure, ECON and SIFU suffer from bending legs and wrong arms problems. SiTH generates an over-smoothed surface and missing parts. Our method can reconstruct more accurate poses while preserving geometric details.

Table 3: 3D reconstruction comparison with SOTA methods. Only our method trained *without* 3D supervision.

Method	THuman2.0			CustomHumans		
	Chamfer $\downarrow$	P2S $\downarrow$	Normal $\uparrow$	Chamfer $\downarrow$	P2S $\downarrow$	NC $\uparrow$
ECON Xiu et al. (2023)	2.342	2.431	0.765	2.107	2.355	0.771
GTA Zhang et al. (2023c)	2.201	2.314	0.773	1.987	2.115	0.769
SiTH Ho et al. (2024)	2.519	2.442	0.786	2.223	2.584	0.785
SIFU Zhang et al. (2024)	<b>2.063</b>	2.205	0.792	1.864	1.976	0.778
Ours	2.134	<b>2.118</b>	<b>0.823</b>	<b>1.729</b>	<b>1.835</b>	<b>0.834</b>

### 4.3 ABLATION STUDIES

We conduct ablation studies to evaluate the effectiveness of our SMPL-X dual branch Gaussian prediction model, coarse-to-fine refinement strategy, and back-view refine ControlNet, as well as our SMPL-X refinement. We show the quantitative ablation results in Table 4. The performance decrease when any component is removed. SMPL-X dual branch plays an important role in adding human priors through structural features to Image Gaussians predicted by UNet. We visualize the rendered images using our model without SMPL-X dual branch as the guidance and those produced by our full model, as shown in Fig. 8. Without SMPL-X dual branch as guidance, the side view im-

ages exhibit significant artifacts, such as misaligned arms and unnatural shapes of clothes and heads, highlighted in the red boxes. This demonstrates the effectiveness of our SMPL-X dual branch Gaussian prediction design. The predicted Gaussians lack human shape and pose priors without SMPL-X guidance, resulting in unnatural shapes and poses. The initial SMPL-X prediction is not accurate, we ablate the effectiveness of our iterative SMPL-X refinement with the 3D prior learned in our HGM model. We also visualize the SMPL-X refinement in the appendix. Two-view refinement double the number of the Gaussians to improve the reconstruction quality. Gaussians tend to concentrate more on the front view without the two-view refinement strategy, leading to poorer rendering of the back part. Additionally, the diffusion-based refinement is crucial for improving the novel view synthesis quality, especially for the back-view images as shown in Fig. 4. The best performance is achieved with all three components.

#### 4.4 DISCUSSIONS

Our method improves upon previous approaches like ECON Xiu et al. (2023) and SIFU Zhang et al. (2024) by leveraging our pre-trained HGM model to incorporate 3D priors, specifically side view masks, for enhanced SMPL-X refinement. Unlike earlier techniques that only align 2D front-view masks and normals, our approach achieves better reconstruction accuracy and alignment. The refined SMPL-X also benefits our human Gaussians reconstruction by providing structural information encoded in learnable tokens. Our 3D Gaussians-based method offers significant advantages in rendering speed, achieving 300 FPS compared to NeRF-based methods like SHERF Hu et al. (2023), which manages only 2 FPS. Furthermore, the mesh extracted from our 3D Gaussians with normal refinement attains high 3D reconstruction quality. In summary, our method excels in both high-quality rendering and accurate 3D reconstruction, offering a comprehensive solution.

Table 4: Ablation study for each component.

Components	NVS			3D reconstruction		
	PSNR( $\uparrow$ )	SSIM( $\uparrow$ )	LPIPS( $\downarrow$ )	Chamfer( $\downarrow$ )	P2S( $\downarrow$ )	NC( $\uparrow$ )
w/o SMPL-X dual branch	21.85	0.908	0.769	2.421	2.543	0.776
w/o SMPL-X refine	22.86	0.921	0.656	2.245	2.346	0.798
w/o two-view refine	22.95	0.924	0.641	2.301	2.343	0.781
w/o Diffusion refine	23.11	0.921	0.637	2.145	2.176	0.812
Full model	<b>23.54</b>	<b>0.938</b>	<b>0.524</b>	<b>2.134</b>	<b>2.118</b>	<b>0.823</b>

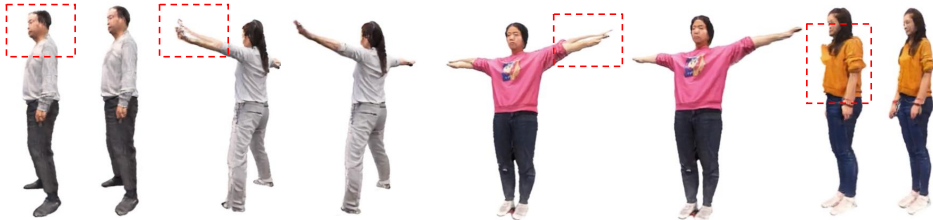


Figure 8: Ablation studies in terms of SMPL-X dual branch guidance. For each pair of images, left one is the results of our model w/o SMPL-X guidance. **Details are highlighted in the red boxes.**

## 5 CONCLUSION

In this paper, we introduce a novel generalizable single-view human Gaussian reconstruction framework. By incorporating human priors through a SMPL-X dual branch Gaussian prediction and diffusion priors using a refinement ControlNet, our method effectively handles invisible parts and varying poses. By incorporating our pre-trained HGM, inaccurate SMPL-X is iteratively refined, which benefits the Gaussian reconstruction quality. Combining all of these techniques, our method can generalize well to unseen subjects for high-quality and view-consistent reconstruction. We validate the proposed method on several benchmarks and demonstrate that it achieves state-of-the-art performance in both novel view synthesis and 3D reconstruction.

## REFERENCES

- 540  
541  
542 Thiemo Alldieck, Mihai Zanfir, and Cristian Sminchisescu. Photorealistic monocular 3d reconstruction  
543 of humans wearing clothing. In *CVPR*, 2022.
- 544  
545 Zhongang Cai, Daxuan Ren, Ailing Zeng, Zhengyu Lin, Tao Yu, Wenjia Wang, Xiangyu Fan, Yang  
546 Gao, Yifan Yu, Liang Pan, Fangzhou Hong, Mingyuan Zhang, Chen Change Loy, Lei Yang, and  
547 Ziwei Liu. HuMMan: Multi-modal 4d human dataset for versatile sensing and modeling. In  
548 *ECCV*, 2022.
- 549  
550 Yukang Cao, Yan-Pei Cao, Kai Han, Ying Shan, and Kwan-Yee K. Wong. Dreamavatar: Text-and-  
551 shape guided 3d human avatar generation via diffusion models. In *CVPR*, 2024.
- 552  
553 Enric Corona, Mihai Zanfir, Thiemo Alldieck, Eduard Gabriel Bazavan, Andrei Zanfir, and Cristian  
554 Sminchisescu. Structured 3d features for reconstructing relightable and animatable avatars. In  
555 *CVPR*, 2023.
- 556  
557 Yao Feng, Vasileios Choutas, Timo Bolkart, Dimitrios Tzionas, and Michael J. Black. Collaborative  
558 regression of expressive bodies using moderation. In *3DV*, 2021.
- 559  
560 Xiangjun Gao, Jiaolong Yang, Jongyoo Kim, Sida Peng, Zicheng Liu, and Xin Tong. Mps-nerf:  
561 Generalizable 3d human rendering from multiview images. In *PAMI*, 2022.
- 562  
563 Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with  
564 submanifold sparse convolutional networks. In *CVPR*, 2018.
- 565  
566 Hsuan-I Ho, Lixin Xue, Jie Song, and Hilliges Otmar. Learning locally editable virtual humans. In  
567 *CVPR*, 2023.
- 568  
569 Hsuan-I Ho, Jie Song, and Otmar Hilliges. Sith: Single-view textured human reconstruction with  
570 image-conditioned diffusion. In *CVPR*, 2024.
- 571  
572 Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli,  
573 Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d, 2024.
- 574  
575 Shoukang Hu and Ziwei Liu. Gauhuman: Articulated gaussian splatting from monocular human  
576 videos. In *CVPR*, 2024.
- 577  
578 Shoukang Hu, Fangzhou Hong, Liang Pan, Haiyi Mei, Lei Yang, and Ziwei Liu. Sherf: Generaliz-  
579 able human nerf from a single image. In *ICCV*, 2023.
- 580  
581 Yangyi Huang, Hongwei Yi, Weiyang Liu, Haofan Wang, Boxi Wu, Wenxiao Wang, Binbin Lin,  
582 Debing Zhang, and Deng Cai. One-shot implicit animatable avatars with model-based priors. In  
583 *ICCV*, 2023.
- 584  
585 Yangyi Huang, Hongwei Yi, Yuliang Xiu, Tingting Liao, Jiaxiang Tang, Deng Cai, and Justus Thies.  
586 TeCH: Text-guided Reconstruction of Lifelike Clothed Humans. In *3DV*, 2024.
- 587  
588 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splat-  
589 ting for real-time radiance field rendering. *TOG*, 2023.
- 590  
591 Muhammed Kocabas, Jen-Hao Rick Chang, James Gabriel, Oncel Tuzel, and Anurag Ranjan.  
592 HUGS: Human gaussian splatting. In *CVPR*, 2024.
- 593  
594 Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural human performer: Learn-  
595 ing generalizable radiance fields for human performance rendering. *NIPS*, 2021.
- 596  
597 Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan  
598 Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3d: Fast text-to-3d with sparse-view gen-  
599 eration and large reconstruction model. In *ICLR*, 2024.
- 600  
601 Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick.  
602 Zero-1-to-3: Zero-shot one image to 3d object. In *ICCV*, 2023.

- 594 Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and  
595 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.  
596
- 597 Panwang Pan, Zhuo Su, Chenguo Lin, Zhen Fan, Yongjie Zhang, Zeming Li, Tingting Shen, Yadong  
598 Mu, and Yebin Liu. Humansplat: Generalizable single-image human gaussian splatting with  
599 structure priors. *arXiv preprint arXiv:2406.12459*, 2024.
- 600 Shunsuke Saito, , Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao  
601 Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *ICCV*,  
602 2019.  
603
- 604 Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned  
605 implicit function for high-resolution 3d human digitization. In *CVPR*, 2020.
- 606 Yichun Shi, Peng Wang, Jianglong Ye, Long Mai, Kejie Li, and Xiao Yang. Mvdream: Multi-view  
607 diffusion for 3d generation. *arXiv:2308.16512*, 2023.  
608
- 609 Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Splatter image: Ultra-fast  
610 single-view 3d reconstruction. In *CVPR*, 2024.
- 611 Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm:  
612 Large multi-view gaussian model for high-resolution 3d content creation. *arXiv preprint*  
613 *arXiv:2402.05054*, 2024a.  
614
- 615 Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative  
616 gaussian splatting for efficient 3d content creation. In *ICLR*, 2024b.
- 617 Vikram Voleti, Chun-Han Yao, Mark Boss, Adam Letts, David Pankratz, Dmitrii Tochilkin, Chris-  
618 tian Laforte, Robin Rombach, and Varun Jampani. SV3D: Novel multi-view synthesis and 3D  
619 generation from a single image using latent video diffusion. In *ECCV*, 2024.  
620
- 621 Peng Wang and Yichun Shi. Imagedream: Image-prompt multi-view diffusion for 3d generation,  
622 2023.  
623
- 624 Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J. Black. ICON: Implicit Clothed hu-  
625 mans Obtained from Normals. In *CVPR*, 2022.
- 626 Yuliang Xiu, Jinlong Yang, Xu Cao, Dimitrios Tzionas, and Michael J. Black. ECON: Explicit  
627 Clothed humans Optimized via Normal integration. In *CVPR*, 2023.  
628
- 629 Dejie Xu, Ye Yuan, Morteza Mardani, Sifei Liu, Jiaming Song, Zhangyang Wang, and Arash Vahdat.  
630 Agg: Amortized generative 3d gaussians for single image to 3d. *arXiv preprint 2401.04099*, 2024.
- 631 Yinghao Xu, Hao Tan, Fujun Luan, Sai Bi, Peng Wang, Jiahao Li, Zifan Shi, Kalyan Sunkavalli,  
632 Gordon Wetzstein, Zexiang Xu, and Kai Zhang. Dmv3d: Denoising multi-view diffusion using  
633 3d large reconstruction model, 2023.  
634
- 635 Xu Yinghao, Shi Zifan, Yifan Wang, Chen Hansheng, Yang Ceyuan, Peng Sida, Shen Yujun, and  
636 Wetzstein Gordon. Grm: Large gaussian reconstruction model for efficient 3d reconstruction and  
637 generation, 2024.
- 638 Tao Yu, Zerong Zheng, Kaiwen Guo, Pengpeng Liu, Qionghai Dai, and Yebin Liu. Function4d:  
639 Real-time human volumetric capture from very sparse consumer rgb-d sensors. In *CVPR*, 2021.  
640
- 641 Hongwen Zhang, Yating Tian, Yuxiang Zhang, Mengcheng Li, Liang An, Zhenan Sun, and Yebin  
642 Liu. Pymaf-x: Towards well-aligned full-body model regression from monocular images. *IEEE*  
643 *Transactions on Pattern Analysis and Machine Intelligence*, 2023a.
- 644 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
645 diffusion models. In *ICCV*, 2023b.  
646
- 647 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.

648 Zechuan Zhang, Li Sun, Zongxin Yang, Ling Chen, and Yi Yang. Global-correlated 3d-decoupling  
649 transformer for clothed avatar reconstruction. In *NIPS*, 2023c.

651 Zechuan Zhang, Zongxin Yang, and Yi Yang. Sifu: Side-view conditioned implicit function for  
652 real-world usable clothed human reconstruction. In *CVPR*, 2024.

653 Shunyuan Zheng, Boyao Zhou, Ruizhi Shao, Boning Liu, Shengping Zhang, Liqiang Nie, and Yebin  
654 Liu. Gps-gaussian: Generalizable pixel-wise 3d gaussian splatting for real-time human novel view  
655 synthesis. In *CVPR*, 2024.

657 Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. Pamir: Parametric model-conditioned implicit  
658 representation for image-based human reconstruction. In *PAMI*, 2021.

## 661 A APPENDIX

662 We introduce the following content in the appendix: SMPL-X optimization and mesh optimization  
663 details, **SMPL-X evaluation**, **back-view details and evaluation**, **additional comparison**, experimental  
664 environment, network structures, limitations, and more visualizations.

665 **SMPL-X optimization and mesh optimization details** During optimization, we render SMPL-X  
666 side views and compute the side view mask loss and normal loss for a total of 45 iterations. SMPL-  
667 X parameters are updated at each iteration. The updated SMPL-X parameters are fed into HGM to  
668 update GS every 15 iterations (3 times in total), reducing the overall optimization time. The initial  
669 SMPL-X are not input to HGM for GS prediction until after the first 15 iterations because poorly-  
670 aligned SMPL-X can lead to degraded GS. For the front view, we utilize the original image mask  
671 rather than the one rendered from GS to stabilize the process. We simultaneously render 12 views  
672 (one front view and all the other views are considered as side views) to compute the mask loss. The  
673 loss weights are set as follows:  $\lambda_{front} = 10$ ,  $\lambda_{side} = 1$ , and  $\lambda_n = 0.5$ . For the normal loss, we  
674 only utilize the front and back views with a pre-trained normal estimator from ICON. Throughout  
675 the optimization process, HGM remains fixed while only SMPL-X parameters are updated. The  
676 elegance of our method lies in its iterative nature: GS refines SMPL-X and better-aligned SMPL-X  
677 estimates feed back into the HGM model to generate improved 3D Gaussians, which in turn enhance  
678 the reconstruction. We show the visualization of our side view mask rendered from iteratively recon-  
679 structed Gaussians by HGM, initial SMPL-X, refined SMPL-X, and our final-extracted meshes in  
680 Fig. 11, as shown in the figure, the side view masks effectively help refine the initial SMPL-X error  
681 for accurate reconstruction. For mesh refinement, we minimize the  $\mathcal{L}_1$  loss between the predicted  
682 normal map and the rendered normal map. We also add Laplacian loss for the preservation of the  
683 local structure. The whole Gaussian reconstruction takes around 40s and the mesh optimization and  
684 extraction takes another 30s.

685 **Additional comparison with TeCH** We compare our method with TeCH on CustomHumans  
686 dataset quantitatively in Tab. 5. TeCH needs 4-5 hours for each sample, so we use 10 samples  
687 from CustomHumans dataset for comparison. TeCH Huang et al. (2024) has several obvious limita-  
688 tions compared with us: (1) The geometry refinement from SDS is not stable and surface is broken  
689 as shown in the left part of Fig 9 even though capturing more high-frequency geometric details. (2)  
690 Long time optimization: it needs 4-5 hours optimization, while ours use only 85s. (3) The caption  
691 guidance can sometimes be incorrect. For example, as shown in the left part of Fig 9, the wrong  
692 caption of the gender resulting wrong face reconstruction. (4) SMPL-X error leading to bending legs  
693 and wrong geometry, which is the same issue in SIFU, SiTH, ECON and GTA as shown in SIFU,  
694 SiTH, ECON and GTA.

695 Table 5: Additional evaluation with TeCH.

	PSNR( $\uparrow$ )	SSIM( $\uparrow$ )	LPIPS( $\downarrow$ )	CD( $\downarrow$ )	P2S( $\downarrow$ )	NC( $\uparrow$ )
Ours	24.56	0.949	0.051	1.715	1.844	0.833
TeCH	23.87	0.927	0.079	2.232	2.432	0.778

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

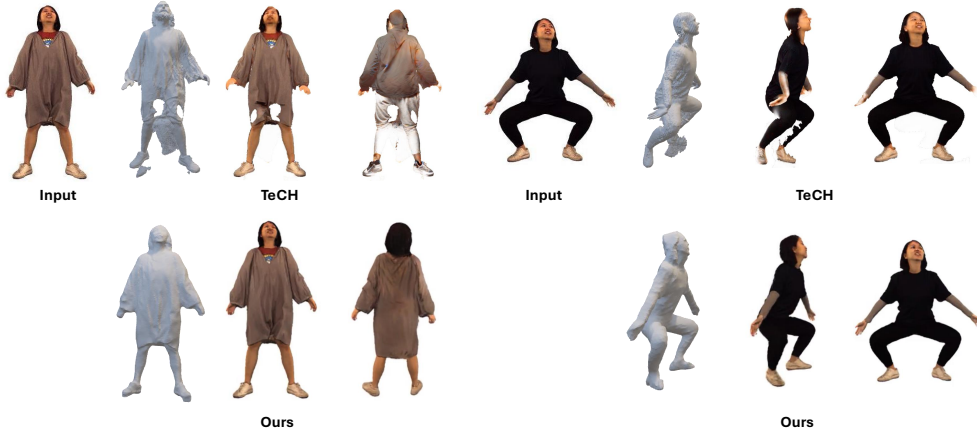


Figure 9: Visual comparison with TeCH on loose cloth and challenging pose cases.

**SMPL-X refinement evaluation** We evaluated it using SMPL-X initializations from PyMAF-X Zhang et al. (2023a) and PIXIE Feng et al. (2021). We calculated the MPJPE (mm) using the first 22 body joints defined in SMPL-X on both THuman2.0 Yu et al. (2021) and CustomHumans Ho et al. (2023) datasets. In the Tab. 6, ‘Initial’ means the direct estimation from SMPL-X predictors. ‘w/o side-views’ represents optimization without side-views mask loss. ‘ours’ refers to our optimization with all the loss including side-views mask loss. From the table we can see, our method successfully refines the initial SMPL-X estimates using side view priors from our HGM, which significantly reduces the error compared with without using side-view masks. Although PyMAF-X provides better initial SMPL-X estimates than PIXIE, both methods achieve comparable MPJPE scores after optimization, as the side-view mask loss guides them toward similar convergence points. This also demonstrates our method is robust to diverse SMPL-X initial estimators and can effectively improve the initial SMPL-X estimation.

Table 6: SMPL-X refinement evaluation in terms of MPJPE.

Dataset	PIXIE			PyMAF-X		
	Initial(↓)	w/o side-views(↓)	Ours(↓)	Initial(↓)	w/o side-views(↓)	Ours(↓)
CustomHumans	75.79	65.33	39.11	65.20	58.12	39.78
Thuman2.0	80.11	72.36	44.30	71.18	65.65	44.84

**Backview refinement details and evaluation** We apply original ControlNet architecture and initialized the ControlNet with ControlNet-tile model. ControlNet-tile is originally trained as a image super-resolution model, we finetune the ControlNet part with our constructed data pair at learning rate of 1e-5, with the base SD1.5 keep fixed. Data pair construction involves first training our HGM using single-view input without full convergence. Subsequently, we perform inference, render the back view and down-sample it. The resulting back-view renderings, intentionally designed to have lower quality, served as conditioning inputs for our ControlNet training. In order to validate the effectiveness of our proposed back-view refinement strategy, we do quantitative evaluation with SiTH Ho et al. (2024) and Huang et al. (2024) for back-view quality on CustomHumans dataset. We use SSIM, LPIPS and KID as evaluation metrics between the ground truth and generated back view images. SiTH generates back-view using pure hallucination, which always generate unrealistic image as shown in Fig. 4 and Tab. 7. TeCH use SDS loss to optimize the back view, however, the back view always fits to wrong SMPL-X pose and imprecise text description, which leads to lower generation quality.

**Experimental environment** We conduct all the experiments on NVIDIA RTX A6000 GPU. The experimental environment is PyTorch 2.2.1 and CUDA 12.2.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

Table 7: Backview evaluation.

	SSIM( $\uparrow$ )	LPIPS( $\downarrow$ )	KID( $\times 10^{-3}$ $\downarrow$ )
Ours	0.949	0.079	9.26
SiTH	0.855	0.123	29.8
TeCH	0.876	0.118	20.3

**Network structures** Our UNet model consists of 6 down blocks, 1 middle block and 5 up blocks, with an input image size of 512×512 and an output Gaussian feature map size of 256×256. We use 2 input views, resulting in a total of 256×256×4 = 131,072 output Gaussians. Each block contains several residual layers and an optional down-sample or up-sample layer. For the last 3 down blocks, the middle block, and the first 3 up blocks, we insert cross-view self-attention layers after the residual layers. The final feature maps are processed by a 1×1 convolution layer to produce 14-channel pixel-wise Gaussian features. We adopt SiLU activation and group normalization for the UNet. Our  $Tr_{mix}$  share the same structure, consisting of multiple self-attention blocks. Specifically, they each have 5 up blocks and 5 down blocks. The down-sample channels are [64, 128, 256, 512, 1024], and the up-sample channels are [1024, 512, 256, 128, 64]. The input dimensions for  $Tr_{mix}$  is 256, respectively. The output dimension for both is 14, which matches the dimension of the Gaussian features. For the attention blocks, we use a memory-efficient attention implementation.



Figure 10: Novel view synthesis comparison SHERF on HuMMAN dataset.

**Limitations** Currently, our method is hard to generate high-quality hands and faces, which could potentially be solved by utilizing the SMPL-X model and regional diffusion guidance such as SDS loss for further refinement.

**Additional results** We provide visual comparison with the SOTA NeRF-based method: SHERF Hu et al. (2023) shown in 10. While SHERF predicts blurry results and loses fidelity, our method preserves high-frequency details and generates realistic back views such as wrinkles and hair that fit well to the front views. We also provide more mesh reconstruction results in Fig. 12. We also provide more visual comparison with other methods for loosening clothes in Fig. 13 and Fig. 14 and extreme poses in Fig. 15. As shown in the figures, our method robustly recovers the realistic details of the loosening clothes, while SiTH shrinks to the SMPL-X body. LGM generates unnatural colors. SIFU and PIFu generate blurry rendering with artifacts.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

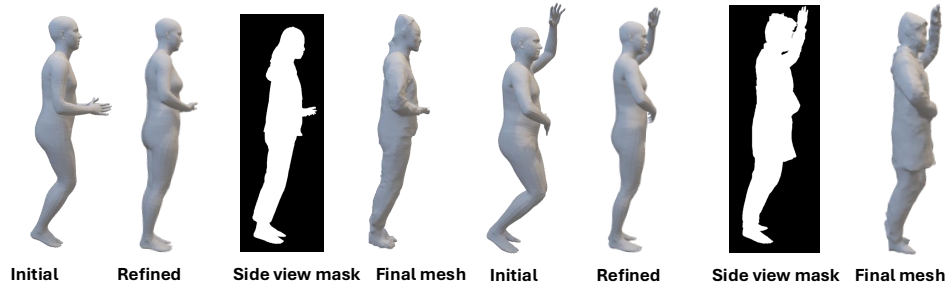


Figure 11: SMPL-X refinement visualization.



Figure 12: Mesh reconstruction visualization for THuman 2.1 data and in the wild images.



864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917



Figure 13: Novel view synthesis comparison for loosing clothes. **The details are highlighted in the red boxes.**

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971



Figure 14: Novel view synthesis comparison for loosing clothes. **The details are highlighted in the red boxes.**

972  
 973  
 974  
 975  
 976  
 977  
 978  
 979  
 980  
 981  
 982  
 983  
 984  
 985  
 986  
 987  
 988  
 989  
 990  
 991  
 992  
 993  
 994  
 995  
 996  
 997  
 998  
 999  
 1000  
 1001  
 1002  
 1003  
 1004  
 1005  
 1006  
 1007  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015  
 1016  
 1017  
 1018  
 1019  
 1020  
 1021  
 1022  
 1023  
 1024  
 1025

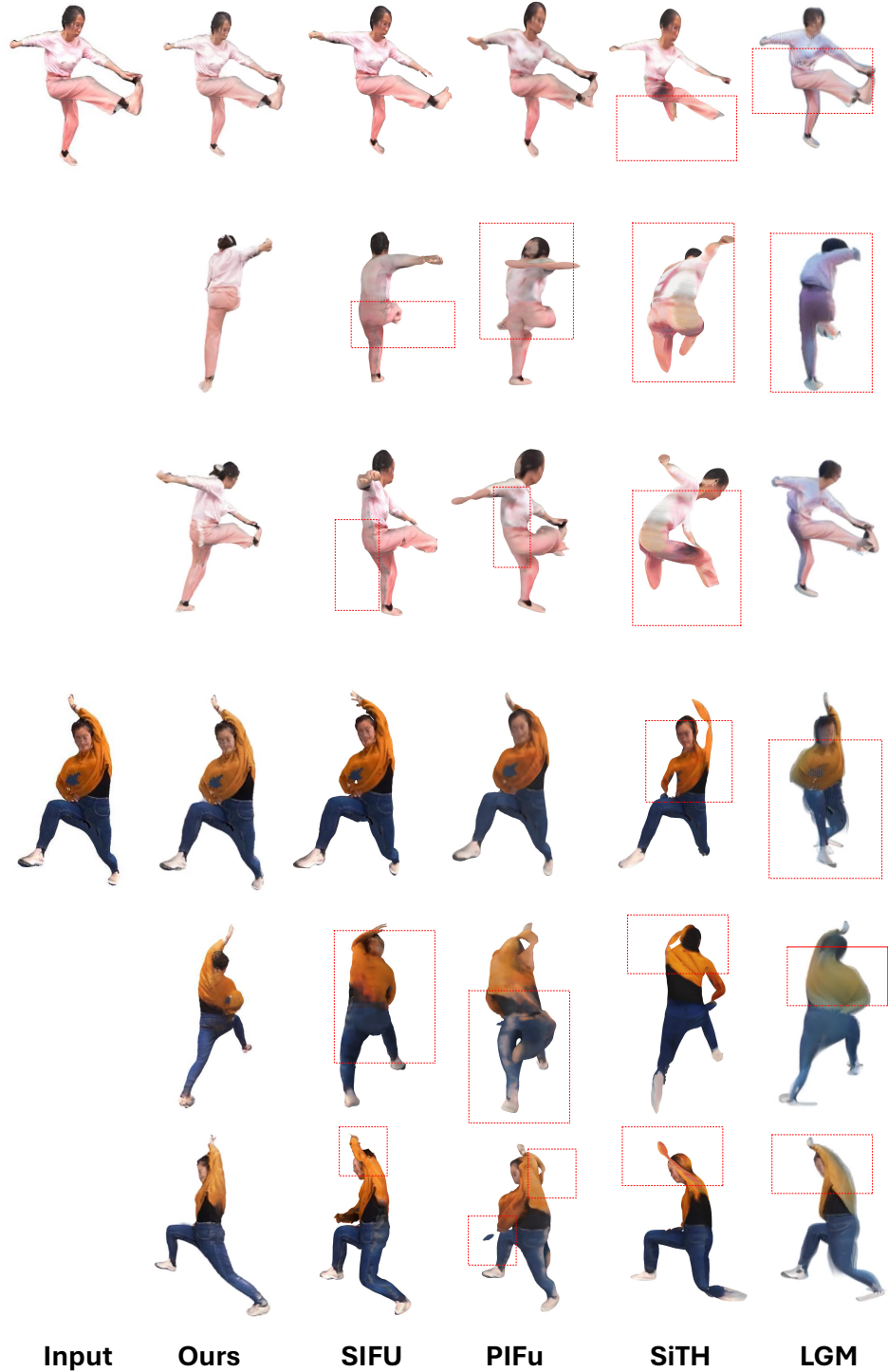


Figure 15: Novel view synthesis comparison for extreme poses from HuMMan dataset. **The details are highlighted in the red boxes.**