

Human101: Training 100+FPS Human Gaussians in 100s from 1 View

Mingwei Li Jiachen Tao Zongxin Yang Yi Yang[†]
 ReLER, CCAI, Zhejiang University

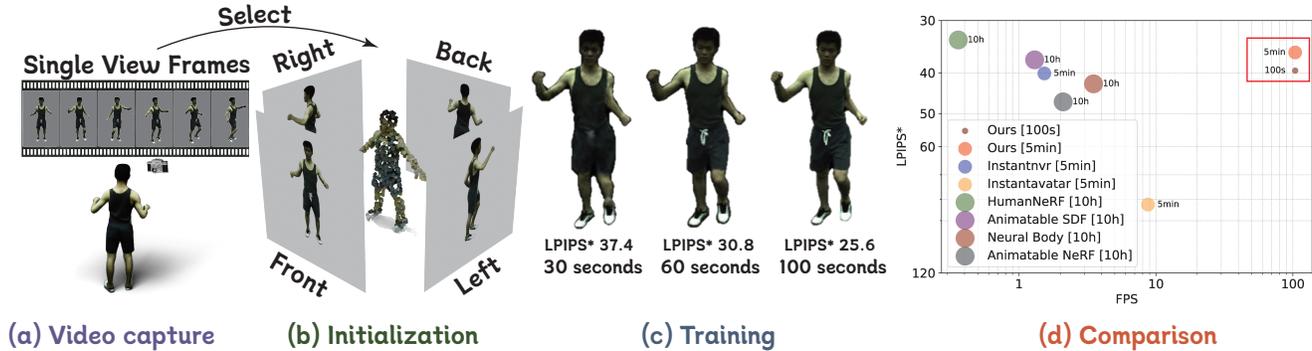


Figure 1. **One VR common use-case.** (a) The user captures a short monocular video and proceeds to upload it. (b) Our model automatically selects four frames from the monocular video and after an initialization process, we can obtain an initial point cloud. (c) Our model can learn in minutes to get a dynamic human representation. (d) Our model achieves comparable or better visual quality while rendering much faster than previous works. $LPIPS^* = LPIPS \times 10^3$. The area of each circle is proportional to the training time required, with larger areas representing longer training durations.

Abstract

Reconstructing the human body from single-view videos plays a pivotal role in the virtual reality domain. One prevalent application scenario necessitates the rapid reconstruction of high-fidelity 3D digital humans while simultaneously ensuring real-time rendering and interaction. Existing methods often struggle to fulfill both requirements. In this paper, we introduce Human101, a novel framework adept at producing high-fidelity dynamic 3D human reconstructions from 1-view videos by training 3D Gaussians in 100 seconds and rendering in 100+ FPS. Our method leverages the strengths of 3D Gaussian Splatting, which provides an explicit and efficient representation of 3D humans. Standing apart from prior NeRF-based pipelines, Human101 ingeniously applies a Human-centric Forward Gaussian Animation method to deform the parameters of 3D Gaussians, thereby enhancing rendering speed (i.e., rendering 1024-resolution images at an impressive 60+ FPS and rendering 512-resolution images at 100+ FPS). Experimental results indicate that our approach substantially eclipses current methods, clocking up to a $10 \times$ surge in frames per second and delivering comparable or superior rendering quality. Code and demos will

be released at <https://github.com/longxiang-ai/Human101>.

1. Introduction

In the realm of virtual reality, a prevalent use case involves rapidly crafting custom virtual avatars and facilitating interactions with them. Within this context, two significant technical challenges emerge: **(1)** How can we swiftly produce a digitized virtual avatar, preferably within a user’s acceptable waiting time (e.g., within 3 minutes), using readily available equipment (e.g., a single-camera setup)? **(2)** How can we achieve real-time rendering to cater to the interactive demands of users?

While previous methods [15, 22, 43–45, 70, 72] have made some progress, they still haven’t fully met the requirements of the application scenario described earlier. The limitations of these methods can be summarized in two main points: **(1) Slow rendering speed in implicit methods.** Methods based on implicit neural network [15, 22, 43–45] optimization using NeRF have slower rendering processes and challenges with inverse skinning deformation, preventing real-time rendering. **(2) Slow convergence speed in explicit methods.** Approaches such as those in [70, 72], capable of achieving real-time rendering, necessitate extensive data for training. This requirement results in slower optimization, thus hindering the rapid reconstruction of dynamic humans.

[†]: the corresponding author.

To address these challenges, a more practical and improved approach would fit these goals. **First**, to enhance rendering speed, we should choose a rasterization rendering pipeline, replacing the traditional volume rendering approach. **Second**, to speed up training, we should choose a better representation method that’s easier to optimize, ideally reducing optimization time to just a few minutes. Recently, a novel method [24] has employed 3D Gaussians to explicitly depict 3D scenes. With the integration of a differentiable tile rasterization method, it achieves superior visual quality and a much quicker rendering speed (over 100 FPS) compared to previous works [1, 39, 52]. The emergence of this method makes realizing the described application scenario (*i.e.*, achieving both fast reconstruction and real-time rendering) a tangible possibility.

Recognizing the advantages of [24], we introduce a novel framework for single-view human reconstruction. This framework not only accomplishes dynamic human reconstruction in less than one minute but also ensures real-time rendering capabilities. Merging the fast and straightforward methods of 3D GS with human body structures, we’ve created a new kind of **forward skinning process** for rendering. Different from the usual inverse skinning used by [15, 22, 43, 45] this forward skinning deformation method avoids searching for the corresponding canonical points of the target pose points but directly deform the canonical points into observation space. Cause [24] utilizes 3D Gaussians rather than just points, we use a **Human-centric Forward Gaussian Animation** method to deform the positions, rotations, and scales of Gaussians, and modify spherical coefficients by rotating their directions. For faster convergence, we design a **Canonical Human Initialization** method to initialize the original Gaussians.

To validate the effectiveness of the proposed pipeline, we conduct extensive experiments on ZJU-MoCap Dataset [44] and the Monocap Dataset [15]. Results show that Human101 could not only swiftly reconstruct a dynamic human, but also outperform incredible rendering speed together with better visual quality. With a single RTX 3090 GPU, our method can learn **in 100 seconds** to get comparable or better visual quality and maintain **100+ FPS rendering speed**, which makes it a tangible possibility for real-time interactive applications and immersive virtual reality experiences.

Our **contributions** can be summarised as follows:

- We introduce an innovative approach to dynamically represent 3D human bodies by employing 3D Gaussian Splatting [24], utilizing its efficient and explicit representation capabilities for detailed and accurate human modeling. We have proposed a **Canonical Human Initialization** method, which significantly enhances the model’s convergence rate and improves visual detail representation.
- We propose a deformation methodology, composed of **Human-centric Forward Gaussian Animation** and

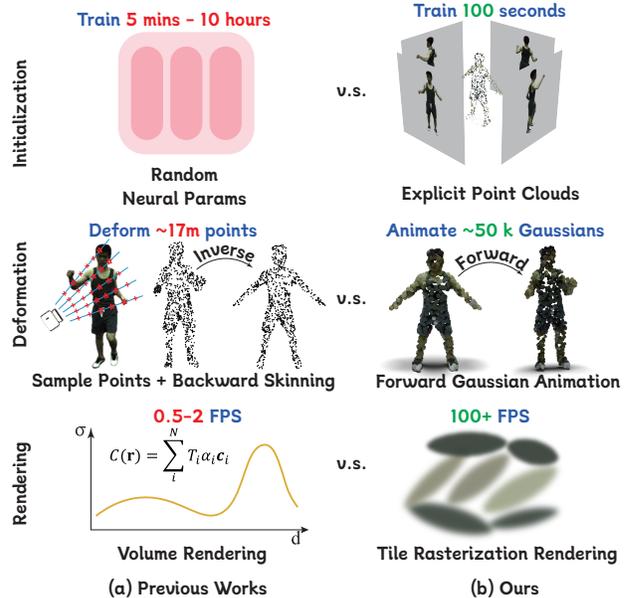


Figure 2. **Method difference with previous works [15, 43] (§ 3).** (1) Initialization. Previous works rely on an implicit representation, initialized with random neural parameters. In contrast, our method employs explicit colored point clouds for initialization, which accelerates convergence. (2) Deformation. Earlier approaches require deforming a larger number of points from the target pose to the canonical pose, consuming more computational resources. (3) Rendering. Compared to the NeRF-based ray marching technique, our approach utilizing rasterization rendering achieves significantly faster rendering speeds.

Human-centric Gaussian Refinement, which is distinct from the prevailing time-consuming inverse skinning frameworks, making it possible to fast reconstruct and render a dynamic human in real time.

- We achieve a $\sim 10.8 \times$ speed increase in rendering during inference (with **FPS 100+** for 512×512 resolution images) compared to previous neural human representations, while simultaneously ensuring comparable or superior rendering quality and higher image resolution.

2. Related Work

Human Reconstruction from Monocular Video. Reconstructing 3D humans from monocular videos is challenging, primarily due to complex human poses and the limited information from a single camera. Significant strides in this area have been made by [4, 5, 21, 26, 49, 50, 53, 56–58, 61, 64, 65, 71]. Approaches like [23, 43–45, 62] have excelled in high-quality reconstruction with precise pose and deformation adjustments. Yet, their lengthy convergence times, often exceeding 10 hours, limit practical utility.

Recent works like [15, 22, 39] have accelerated convergence to about 5 minutes, but rendering speeds remain a bottleneck, with [15] achieving only about 1.5 FPS on an RTX 3090. 3D GS [24], effective for static scenes with its

explicit 3D Gaussian representation and fast GPU sorting, offers a potential solution with over 100 FPS rendering speeds. However, its application to dynamic human reconstruction is not straightforward. Our work harnesses the principles of 3D GS for monocular human reconstruction, targeting high rendering speeds to bridge the gap towards practical implementation in real-world applications.

Human Deformation and Animation. The Skinned Multi-Person Linear model (SMPL) [3, 35, 42] is a prevalent framework for representing human structure, simplifying pose changes with Linear Blend Skinning (LBS). Various generative articulation methods have been explored [2, 7, 11, 19, 40, 67], alongside backward skinning techniques [20, 36, 51, 59] and forward skinning methods [6–8, 12, 21, 31, 33, 60]. For dynamic human reconstruction, studies like [15, 43, 62] use neural networks to enhance the deformation process, applying residuals to point coordinates or blending weights. **Unlike** these point-based approaches, our method employs 3D Gaussians [24] for spatial representation, accounting for position, rotation, and scale. Our human-centric forward skinning deformation approach successfully animates humans based on 3D Gaussians, effectively addressing challenges such as artifacts and jaggedness after deformation.

Accelerating Neural Rendering. Since the introduction of NeRF by [37], numerous studies have sought to accelerate neural scene rendering. Techniques utilize voxel grids [14, 18, 32, 39, 47, 52, 66], explicit surfaces [10, 28, 41], and point-based representations [27, 30, 46, 48, 68] or other methods to speed up rendering process. These methods effectively minimize the number of NeRF MLP evaluations required, thereby reducing computational costs. Focusing on human body, some approaches [34, 54, 72] utilize innovative processes like UV map prediction [9, 29, 70, 72] and voxel grids [13]. However, these techniques predominantly suffer from lengthy training durations and are mostly restricted to static scenes, hindering downstream applications.

A significant breakthrough in this area is the development of 3D Gaussian Splatting (3D GS) [24], utilizing anisotropic 3D Gaussians combined with spherical harmonics [38] to represent 3D scenes. This method effectively circumvents the slow ray marching operation, delivering high-fidelity and high-speed rendering. Nevertheless, its application has been primarily confined to static scenes. Our work pioneers the application of 3D GS to animatable human reconstruction. We extend the capabilities of 3D GS beyond static multi-view scenarios, addressing its limitations in dynamic monocular human movement reconstruction.

3. Method

Overview. In this work, our primary objective is to **swiftly** reconstruct dynamic human movements from **single-view videos**, simultaneously ensuring **real-time rendering** capabilities. Our approach builds upon the techniques introduced

in [15], with the underlying assumption that the cameras are pre-calibrated and each image is accompanied by provided human poses and foreground human masks. Fig. 3 shows the main training pipeline of our model. Within Sec. 3.1, we explore the foundational aspects of 3D Gaussian Splatting (3D GS) [24] and SMPL [35]. Sec. 3.2 delves into the process of canonical human initialization, which is the base of the initial 3D Gaussians, speeding up our training process. Sec. 3.3 presents our novel human-centric forward gaussian animation approach. This section details how we adapt and apply Gaussian models to accurately represent and animate human figures, focusing on achieving both high fidelity and efficiency in dynamic scenarios. Finally, Sec. 3.4 illustrates our refinements of Gaussians and spherical harmonics.

3.1. Preliminary

3D Gaussian Splatting. Our framework employs 3D Gaussian Splatting [24] to parameterize dynamic 3D shapes for 2D image transformation. Differing from NeRF-based methods, we define 3D Gaussians with a full 3D covariance matrix Σ centered at μ in world space. For 2D rendering, the projected covariance matrix Σ' is calculated as:

$$\Sigma' = JV\Sigma V^\top J^\top, \quad (1)$$

where J is the Jacobian of the affine approximation of the projective transformation, and V is the world-to-camera matrix. To simplify learning, Σ is decomposed into a quaternion r for rotation and a 3D-vector s for scaling, yielding rotation matrix R and scaling matrix S . Hence, Σ is expressed as:

$$\Sigma = RSS^\top R^\top. \quad (2)$$

Adapting to dynamic scenes. The method proposed in [24] excels in static scene representation from multi-view data. However, to extend this framework to dynamic human scenarios, we refine the 3D Gaussian model within the canonical space. We systematically optimize the Gaussians’ defining attributes: *position* x , *rotation* r , *scale* s , and *view direction* d of the radiance field, which is represented using spherical harmonics (SH). These refinements are carried out through an additional deformation field that enables the precise capture of the nuanced movements inherent to human dynamics. For clarity and conciseness, the 3D Gaussians in our framework are denoted as $\mathcal{G}(x, r, s, d)$, succinctly encapsulating the parameters critical to modeling dynamic human forms.

SMPL and LBS deformation. SMPL [35] is a widely used human skinning model, composed of 6890 vertices and 13776 triangular faces. Each of the SMPL vertices owns a 3D position v_i and a corresponding weight vector w_i . The deformation of the SMPL mesh \mathcal{M} is performed using the Linear Blend Skinning (LBS) technique, which deforms the mesh based on a set of pose parameters θ and shape parameters β . Specifically, Given θ and β , the LBS technique can

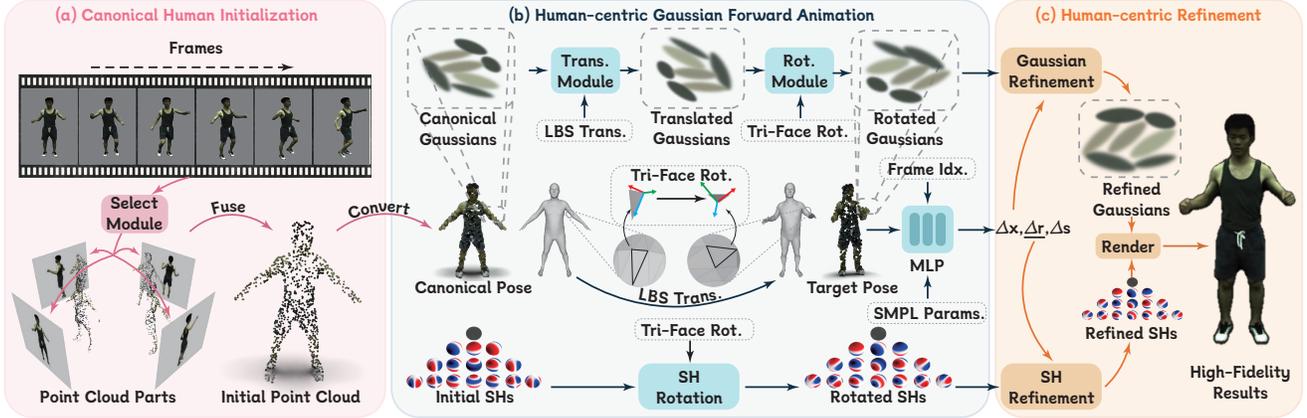


Figure 3. **Overview of Human101** (§ 3). (a) **Canonical Human Initialization** (§ 3.2). We use an offline model [65] to extract 4 point cloud parts from 4 selected frames, and fuse them into a canonical point cloud, which then is converted into canonical Gaussians. (b) **Human-centric Gaussian Forward Animation** (§ 3.3). We deform canonical 3D Gaussians into the target pose by modifying Gaussian positions x , rotations r , and scales s . And we rotate the spherical coefficients with triangle face rotation. (c) **Human-centric Gaussian Refinement** (§ 3.4). We refine positions x , rotations r and scales s of Gaussians and refine the view direction d of spherical harmonics.

transform vertex v_{can} from canonical space to observation space v_{ob} as follows:

$$v_{ob} = \sum_{i=1}^N w_i G_i(\theta, \beta) v_{can}, \quad (3)$$

where N is the joint number, w_i is the blend weight of v , and $G_i(\theta, \beta)$ is the transformation matrix of joint i .

3.2. Canonical Human Initialization

Previous Point Cloud Extraction. The initialization of the canonical Gaussian space is critically dependent on the point cloud data quality, with variations in initial point clouds substantially influencing the model’s convergence speed and the refinement of outcomes. While prior studies [24] have adeptly derived initial point clouds from multi-view data leveraging COLMAP techniques [55], these methods excel predominantly in static scenes. **However**, such approaches do not perform well with a single viewpoint and are incapable of estimating point clouds for dynamic data.

Human-centric Point Cloud Parts Extraction. Consequently, taking human structure into account, we consider the adoption of monocular reconstruction methods [64, 65, 71] to acquire the initial point cloud. The prior state-of-the-art [65] is capable of estimating a mesh from a single image input and projecting the image color onto the mesh. However, this approach only attaches the color of one image to the predicted mesh, resulting in a mesh that is only partially colored. Our goal is to select as few images as possible from the input sequence to achieve the best initial results. Our **Automatic Selection** strategy involves extracting multiple sets of images from the sequence, each set containing four images with the human subject’s angles as close to 90 degrees apart as possible. We then choose one set of images

where the poses most closely resemble the “A” pose (shown in the Fig. 3 (a)), which facilitates mesh deformation into the canonical pose in the subsequent steps. These images are labeled according to their orientation: *front F*, *back B*, *left L*, and *right R*. Using [65], we generate four meshes $\mathcal{M}^F(\mathcal{V}^F)$, $\mathcal{M}^B(\mathcal{V}^B)$, $\mathcal{M}^L(\mathcal{V}^L)$, $\mathcal{M}^R(\mathcal{V}^R)$, each representing different postures and having color only on one side.

Canonical Point Cloud Fusion. These meshes are deformed to a canonical pose using inverse Linear Blend Skinning [35].

$$\mathcal{V}_{can}^k = \left(\sum_{i=1}^N w_i^k G_i^k \right)^{-1} \mathcal{V}^k, \text{ for } k \in \{F, B, L, R\} \quad (4)$$

$$\mathcal{P}_{can}(\mathcal{V}_{can}) = \text{Fuse}(\mathcal{M}_{can}^F, \mathcal{M}_{can}^B, \mathcal{M}_{can}^L, \mathcal{M}_{can}^R) \quad (5)$$

As a result, each canonical pose point cloud part is colored on one side only. We then fuse the four canonical point clouds to form an initial point cloud in canonical space, comprising approximately 50,000 points. Following [24], we convert the initial point cloud into canonical Gaussians. Additionally, our model can also initialize using the bare SMPL’s canonical pose mesh (6890 vertices with white colors). While this approach requires a slightly longer convergence time, it still ensures a visually appealing result of comparable quality.

3.3. Human-centric Gaussian Forward Animation

Advantages of Gaussian Approach. Traditional NeRF-based methods [15, 22, 43, 45], involving inverse LBS for deformation and ray sampling in observation space, face efficiency issues at high resolutions (e.g., $512 \times 512 \times 64 \approx 16,777k$), as depicted in Fig. 2. These methods struggle with real-time rendering due to the vast number of sampling points and the slow inverse LBS deformation process.

Our approach deviates from these conventional methods by optimizing explicit Gaussians in canonical space and

animating them into observation space. This method results in more efficient rendering, especially suitable for dynamic scenes.

Human-centric Gaussian Translation. Given the *original Gaussian position* x , the *transformation matrix* of the i -th bone $G_i(\theta, \beta)$, and the *blend weight* of the i -th bone w_i . The *deformed Gaussian position* x' is represented as:

$$x' = \sum_{i=1}^N w_i G_i(\theta, \beta) x, \quad (6)$$

where $G_i(\theta, \beta)$ is computed by SMPL vertex which is the nearest to the i -th Gaussian.

Human-centric Gaussian Rotation. Gaussians inherently display **anisotropic** characteristics, meaning they exhibit different properties in various directions. Therefore, accurate rotation adjustments are crucial when transitioning between canonical and observation poses, ensuring the model’s adaptability across new observation viewpoints.

Unlike positions, rotations cannot be directly derived using Linear Blend Skinning (LBS). To define the rotation of Gaussians, we anchor each Gaussian to the SMPL mesh by identifying the closest triangular facet based on Euclidean distance. The position of a triangular facet’s centroid, denoted as f_j^p , is calculated as the mean of its vertices f_j^1 , f_j^2 , and f_j^3 :

$$f_j^p = \frac{f_j^1 + f_j^2 + f_j^3}{3}. \quad (7)$$

For the i -th Gaussian, the nearest triangle facet is determined by finding the minimum distance to x_i :

$$j^* = \operatorname{argmin}_j \|x_i - f_j^p\|. \quad (8)$$

We then adopt the rotation $R_{f_{j^*}}$ of this facet as the Gaussian’s rotation transformation matrix, which we called Triangular Face Rotation:

$$R_i = R_{f_{j^*}} = e_{\text{can}_{j^*}} e_{\text{ob}_{j^*}}^\top, \quad (9)$$

$$r'_i = \operatorname{Rot}(R_i, r_i), \quad (10)$$

where r'_i is rotated Gaussian rotation, $e_{\text{can}_{j^*}}$ and $e_{\text{ob}_{j^*}}$ represent the orthonormal bases of the facet in the canonical and observed poses, respectively. These bases are computed from the edge vectors’ normalized cross products. To enhance computational efficiency during training, we precompute R_f leveraging the known distribution of SMPL poses.

Rotation of Spherical Harmonics. To ensure that the rotation of spherical harmonics aligns with the human body’s posture, we rotate these functions during rendering. As a part of Gaussians’ attributes, spherical harmonics should rotate together with Gaussians to precisely represent the view-dependent colors. Given the rotation of i -th Gaussian

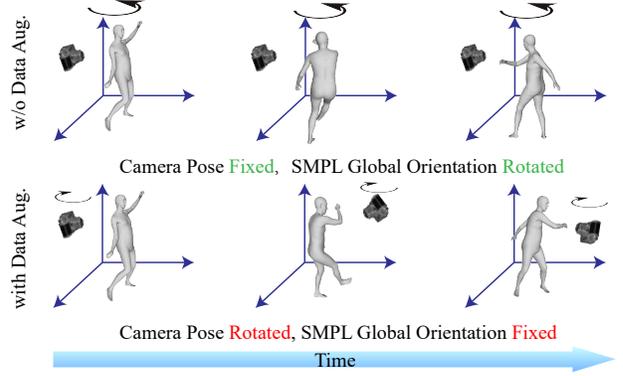


Figure 4. **Data Augmentation Technique** (§ 3.4). Through this method, we simulate the camera rotating around the human to make up multiple camera poses.

R_i , we can easily get rotated spherical harmonics by applying R_i into its direction d_i , that is:

$$d'_i = \text{SH_Rot}(R_i, d_i) = R_i^\top d_i \quad (11)$$

3.4. Human-centric Gaussian Refinement

Adaptive Gaussian Refinement. To capture the uniqueness of each frame, we use the *frame index* t in *Positional Encoding (PE)* to get *per-frame feature* $\gamma(t)$. The *Gaussian positions*, represented as x , are fed through a *Multilayer Perceptron (MLP)* to calculate position, rotation, and scale residuals. These are then used to adjust the Gaussian parameters for each frame, ensuring accuracy and consistency. The adaptive adjustments are represented as follows:

$$\Delta x, \Delta r, \Delta s = F_\Theta(\gamma(x), \gamma(t), \theta, \beta), \quad (12)$$

where θ, β are the parameters of SMPL *pose* and *shape*. This results in deformed parameters: *position* $x'' = x + \Delta x$, *rotation* $r'' = \operatorname{Rot}(\Delta r, r')$, and *scale* $s'' = s + \Delta s$.

Spherical Harmonics Refinement. Similar to the refinement process of Gaussians, given Δr_i of the i -th Gaussian, we can easily get refined Spherical harmonics by this formula:

$$d''_i = \text{SH_Rot}(\Delta r_i, d'_i) = \text{quat_to_rotmat}(\Delta r_i)^\top d'_i \quad (13)$$

Thus, we get the refined Gaussians $\mathcal{G}(x'', r'', s'', d'')$, and send them to the fast rasterization rendering module of [24] to get high-fidelity results.

3.5. Data Augmentation Technique

During our experimental process, we noted that spherical harmonics coefficients tended to overfit when limited to inputs from a single, fixed camera perspective. This resulted in pronounced color biases, as depicted in Fig. 6. To mitigate this issue, we adopted a data augmentation approach that enhances the diversity of camera perspectives. With

Method	Publication	Res.	Train	ZJU-MoCap [44]				MonoCap [15]			
				PSNR↑	SSIM↑	LPIPS*↓	FPS↑	PSNR↑	SSIM↑	LPIPS*↓	FPS↑
<i>Static scene reconstruction method</i>											
3D GS + SMPL Init. [24]	SIGGRAPH 23	512	~ 5min	26.57	0.935	71.70	156	28.47	0.972	29.57	156
3D GS + SMPL Init. [24]	SIGGRAPH 23	1024	~ 5min	26.53	0.944	58.23	51.3	27.47	0.970	34.37	51.3
<i>Time-consuming human reconstruction method</i>											
NeuralBody [44]	CVPR 21	512	~ 10h	29.03	0.964	42.47	3.5	32.36	0.986	16.70	3.5
AnimNeRF [43]	ICCV 21	512	~ 10h	29.77	0.965	46.89	2.1	31.07	0.985	19.47	2.1
AnimSDF [45]	Arxiv 22	512	~ 10h	30.38	0.975	37.23	1.3	32.48	0.988	13.18	1.3
HumanNeRF [62]	CVPR 22	512	~ 10h	30.66	0.969	33.38	0.36	32.68	0.987	15.52	0.36
<i>Time-efficient human reconstruction method</i>											
InstantAvatar [22]	CVPR 23	512	~ 5min	29.21	0.936	82.42	8.75	32.18	0.977	24.98	8.75
InstantNvr [15]	CVPR 23	512	~ 5min	30.87	0.971	40.11	1.53	32.61	0.988	16.68	1.53
Ours		512	~ 100s	31.29	0.964	39.50	104	33.20	0.983	16.55	104
Ours		512	~ 5min	31.79	0.965	35.75	104	32.63	0.982	16.51	104
InstantAvatar [22]	CVPR 23	1024	~ 5min	27.79	0.912	97.33	3.83	32.10	0.978	24.95	3.83
InstantNvr [15]	CVPR 23	1024	~ 5min	30.89	0.974	41.70	0.54	<i>OOM</i>	<i>OOM</i>	<i>OOM</i>	<i>OOM</i>
Ours		1024	~ 100s	31.00	0.968	40.71	68	32.20	0.983	18.31	68
Ours		1024	~ 5min	30.93	0.967	39.87	68	32.13	0.983	17.01	68

Table 1. **Comparison with SOTA (§ 4.2).** We compare Human101 with several baseline methods. (1) Static scene reconstruction method: 3D GS [24]. (2) Time-consuming human reconstruction methods: HumanNeRF [62], AnimSDF [45], NeuralBody [44] and AnimNeRF [43]. (3) Time-efficient human reconstruction methods: InstantNvr [15], InstantAvatar [22]. LPIPS* = LPIPS $\times 10^3$, and “OOM” means out of GPU memory when training. For the FPS metric, we evaluate by calculating the inference time provided by the official pre-trained models. We have marked out **best** and **second best** metrics of time-efficient human reconstruction methods.

SMPL global orientation *rotation matrix* R_s , *global translation matrix* T_s and *camera pose* R_c, T_c . We can get the following equation, describing how SMPL coordinates x_s are transformed into camera coordinates x_c :

$$x_c = R_c(R_s x_s + T_s) + T_c. \quad (14)$$

If we assume that the SMPL coordinates align with the world coordinates, (i.e. $R'_s = E, T'_s = O$), we can get:

$$x_c = R'_c x_s + T'_c. \quad (15)$$

With Eq. (14) and Eq. (15), we can easily get:

$$R'_c = R_c R_s, T'_c = R_c T_s + T_c. \quad (16)$$

This method effectively simulates the camera encircling the subject, as shown in Fig. 4, leading to a more varied orientation in the spherical harmonics. It successfully diversifies the orientations represented by the spherical harmonics, preventing overfitting and the associated color distortions.

3.6. Training

Setting. Our model takes single-view camera parameters, SMPL parameters, frame indices, and images as inputs. During training, frames are randomly sampled from the video. The predicted image, denoted as \hat{I} , is constrained using the L1 loss and S3IM [63] loss. In contrast to the method in [15], we chose not to use the VGG loss due to its high computational demands, which could slow down our training process. For now, we have also decided against incorporating regularization terms. The loss function is formulated as:

$$\mathcal{L}_{rgb} = \lambda_1 \|\hat{I} - I_{gt}\| + \lambda_2 \text{S3IM}(\hat{I}, I_{gt}) \quad (17)$$

where λ_1 and λ_2 are the weighting factors for L1 and S3IM loss respectively, and I_{gt} is the ground truth image.

3.7. Implementation Details

Following [15, 24], we employ the Adam optimizer [25], setting distinct learning rates for different parameters. Our model is trained on an RTX 3090 GPU, and it seamlessly reaches a level of performance comparable to previous methods in just about 100 seconds. Following this, the model requires roughly 10,000 iterations, culminating in convergence within about 5 minutes. Notably, while our model excels at single-view reconstruction, it further enhances accuracy and results when applied to multi-view reconstructions, all within the same 5-minute timeframe. For an in-depth comparison of multi-view reconstructions, we direct interested readers to the supplementary materials. Additionally, a comprehensive overview of our network structure and hyper-parameters can also be found in the supplementary section.

4. Experiments

4.1. Datasets

ZJU-MoCap Dataset. ZJU-Mocap [44] is a prominent benchmark in human modeling from videos, supplying foreground human masks and SMPL parameters. Similar to [15], our experiments engage 6 human subjects (377, 386, 387, 392, 393, 394) from the dataset. Training utilizes one camera, while the remaining cameras are designated for evaluation. Each subject contributes 100 frames for training.

MonoCap Dataset. The MonoCap Dataset combines four multi-view videos from the DeepCap [16] and DynaCap [17] datasets, collected by [43]. This dataset provides essential details like camera parameters, SMPL parameters and human masks. We choose the same 4 subjects as [15] for better comparison. Further details about all the sequences in our study can be found in the supplementary material.

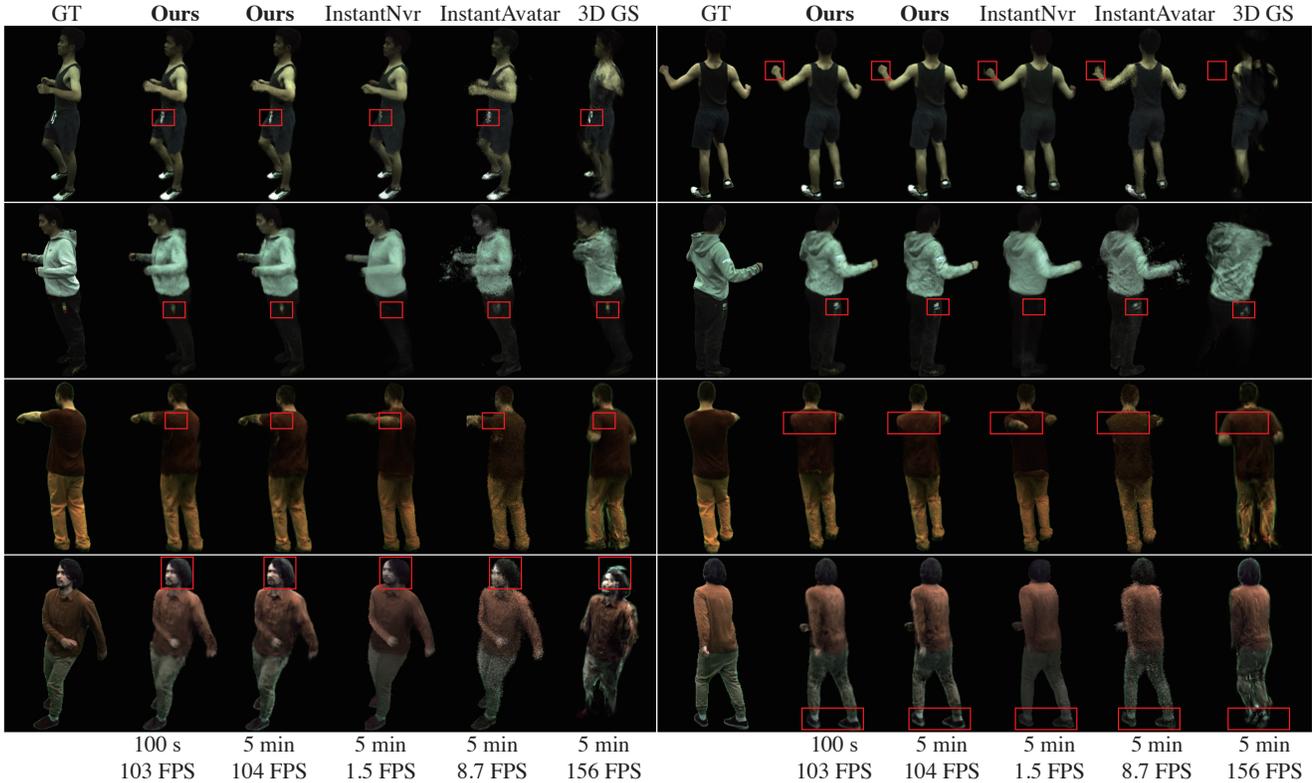


Figure 5. Compare with the state-of-the-art works (§ 4.2). For a fair comparison, we show the results of 512×512 resolutions. Our model delivers results with superior visual quality and richer details, achieving a $67 \times$ increase in FPS compared with the state-of-the-art time-efficient methods. Please **Q** zoom in for a more detailed observation.

4.2. Comparison with the state-of-the-art methods

Baselines. We compare our method with some previous human reconstruction methods [15, 22, 43–45] and baseline method 3D GS [24]. The methods for human reconstruction can be categorized into 3 groups:

- **Static scene reconstruction method.** 3D GS [24] is the backbone of our model. While its original point cloud initialization method failed to extract a point cloud from a single fixed-view camera, we use the canonical SMPL vertices together with white colors as its initial point cloud.
- **Time-consuming human reconstruction methods.** NeuralBody [44] utilizes structured SMPL data together with per-frame latent codes to optimize neural human radiance fields. AnimatableNeRF(AnimNeRF) [43] and AnimatableSDF(AnimSDF) [45] use SMPL deformation and pose-dependent neural blend weight field to model dynamic humans. HumanNeRF [62] further optimizes volumetric human representations, and improves detail quality of rendered image. However, due to the slow optimization of MLPs, these methods usually converge very slow. For instance, [62] takes about 72 hours on 4 RTX 2080Ti GPUs to totally converge.
- **Time-efficient human reconstruction methods.** Utilizing [39]’s voxel grid representation, Instantnvr [15] man-

age to shorten the convergence time into 5 minutes. InstantAvatar [22] combines [39] with a rapid deformation method [8], achieving a fast reconstruction in minutes.

Metrics. We choose PSNR, SSIM, LPIPS [69] as visual quality evaluation metrics, and frame per second (FPS) as rendering speed evaluation metrics. For better comparison, we show $LPIPS^* = LPIPS \times 10^3$ instead. Tab. 1 shows our results compared with others. Here we list only the average metric values of all selected characters on a dataset due to the size limit while more detailed qualitative and quantitative comparisons are in the supplementary material.

Discussion on quantitative results. Tab. 1 presents a comprehensive quantitative comparison between our method and other prominent techniques like InstantNvr [15], InstantAvatar [22], 3D GS [24], HumanNeRF [62], AnimSDF [45], NeuralBody [44], and AnimNeRF [43]. Remarkably, our approach achieves optimization within approximately 100 seconds, yielding results that are comparable with or surpass [15, 22, 24] in terms of PSNR and LPIPS. For 5 minutes’ results, we achieved the highest PSNR and LPIPS scores among rapid reconstruction methods. While our SSIM scores are quite high, they do not reach the state-of-the-art level. This is partly attributed to the characteristics of spherical harmonics in our model. Spherical harmonics, by their na-

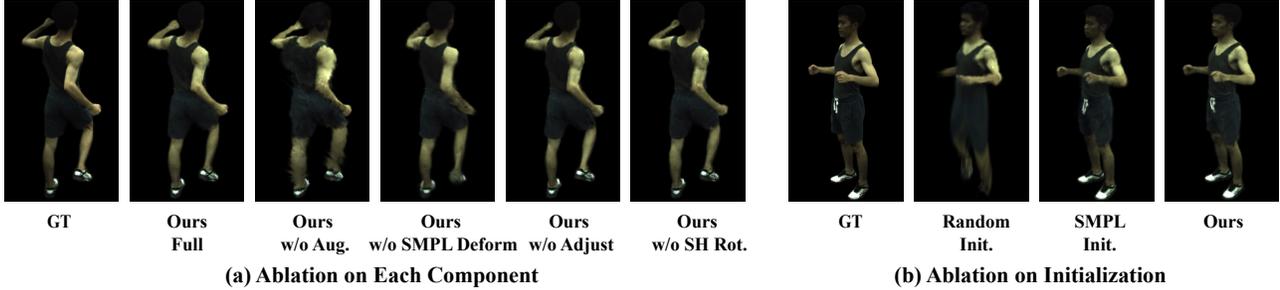


Figure 6. Ablation studies on the Sequence 377 of ZJU-MoCap dataset [44] (§ 4.3). (a) Removing our proposed components leads to and blurry appearance and artifacts. (b) Initialization plays a pivotal role in the geometry quality.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS* \downarrow
w/o SMPL Deformation	28.70	0.964	41.61
w/o Augmentation	29.00	0.961	41.30
w/o Frame Embedding	32.26	0.977	22.51
w/o Gaussian Adjustment	32.55	0.977	23.20
Ours Full	32.18	0.977	21.32

Table 2. Ablation study on component efficacy (§ 4.3). This table demonstrates the impact of individual components in our method. By selectively disabling each part, we validate their effectiveness. The experiments were conducted on Sequence 377 of the ZJU-MoCap dataset [44].

Method	PSNR \uparrow	SSIM \uparrow	LPIPS* \downarrow
Random Init.	26.98	0.951	62.35
SMPL Init.	31.84	0.974	28.60
Ours	32.18	0.977	21.32

Table 3. Ablation on initialization method (§ 4.3). We use Sequence 377 on ZJU-MoCap[44] to test the effectiveness of our canonical human initialization process.

ture, are somewhat limited in capturing high-frequency color information. Moreover, for inference speed, our model is **67 times faster** than InstantNvr [15] and **11 times faster** than InstantAvatar [22] in 512×512 resolution.

Discussion on qualitative results. Fig. 5 showcases a comparison of our model with time-efficient reconstruction works [15, 22, 24]. Our method stands out by providing the most detailed representation and minimal artifacts, as highlighted in Fig. 5 with red boxes around key details. In contrast, [15]’s backside results exhibit unnatural colors on the body’s front due to light penetration and loss of details like missing logos on trousers. InstantAvatar [22] generates noticeable scattered artifacts around the body. Meanwhile, 3D GS [24], lacking a deformation module for dynamic scenes, results in severe limb truncations and facial distortions. Furthermore, in terms of rendering speed, except for [24], our inference speed surpasses all compared methods.

4.3. Ablation Study

Our ablation study results on the ZJU-MoCap [44] 377 sequence are displayed in Tab. 2, Tab. 3, and Fig. 6 (a). Due to the space limitations, more experiment results can be found in the supplementary material.

Human-centric Gaussian rigid deformation. Omitting it, (*i.e.* without SMPL-based rigid translation and rotation, *w/o*

SMPL Deform.), leads to noticeable deficiencies in body parts and a loss of detail, as shown in Fig. 6 (a). This experiment highlights the limitations of simple SMPL deformation in fully capturing complex human motions. The results underscore the necessity of integrating more sophisticated deformation techniques to accurately model human movements.

Spherical rotation. The results, illustrated in Fig. 6 (a), highlight that omitting spherical rotation leads to “spiky” artifacts on the human body, compromising visual quality and causing abnormal lighting effects in the rendered images.

Data augmentation. Our data augmentation technique is crucial, as evidenced by overfitting of spherical harmonics and noticeable skin artifacts in novel view synthesis when it’s absent, as depicted in Fig. 6 (a).

Adaptive Gaussian Refinement. The lack of the Adaptive Gaussian Refinement module results in the model’s inability to capture subtle human deformations (like in fingers), leading to visible artifacts (shown in Fig. 6 (a)).

Initialization method. Our comprehensive study demonstrates the effectiveness of our novel canonical human initialization method. As shown in Fig. 6 (b), this method significantly enhances the quality of human body reconstruction.

5. Conclusion

In this paper, we introduced Human101, a novel framework for single-view human reconstruction using 3D Gaussian Splatting (3D GS) [24]. Our method efficiently reconstructs high-fidelity dynamic human models within just 100 seconds using a fixed-perspective camera. The integration of a novel Canonical Human Initialization, Human-centric Gaussian Forward Animation, and Human-centric Gaussian Refinement, coupled with 3D GS’s explicit representation, significantly improve the rendering speed. Moreover, this enhancement in speed does not sacrifice visual quality. Experiments demonstrate that Human101 outperforms up to 67 times in FPS compared with the state-of-the-art methods and maintain comparable or better visual quality. Human101 sets a new standard in human reconstruction from single-view videos. This breakthrough lays the groundwork for further advances and applications in immersive technologies.

References

- [1] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022. 2
- [2] Alexander W. Bergman, Petr Kellnhofer, Wang Yifan, Eric R. Chan, David B. Lindell, and Gordon Wetzstein. Generative neural articulated radiance fields. In *NeurIPS*, 2022. 3
- [3] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J. Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *Computer Vision – ECCV 2016*. Springer International Publishing, 2016. 3
- [4] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 2
- [5] Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh. Openpose: Realtime multi-person 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 2
- [6] Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In *International Conference on Computer Vision (ICCV)*, 2021. 3
- [7] Xu Chen, Tianjian Jiang, Jie Song, Jinlong Yang, Michael J Black, Andreas Geiger, and Otmar Hilliges. gdna: Towards generative detailed neural avatars. *arXiv*, 2022. 3
- [8] Xu Chen, Tianjian Jiang, Jie Song, Max Rietmann, Andreas Geiger, Michael J. Black, and Otmar Hilliges. Fast-snarf: A fast deformer for articulated neural fields. *Pattern Analysis and Machine Intelligence (PAMI)*, 2023. 3, 7
- [9] Yue Chen, Xuan Wang, Xingyu Chen, Qi Zhang, Xiaoyu Li, Yu Guo, Jue Wang, and Fei Wang. Uv volumes for real-time rendering of editable free-view human performance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16621–16631, 2023. 3
- [10] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. *arXiv preprint arXiv:2208.00277*, 2022. 3
- [11] Enric Corona, Albert Pumarola, Guillem Alenyà, Gerard Pons-Moll, and Francesc Moreno-Noguer. Smplicit: Topology-aware generative model for clothed people. In *CVPR*, 2021. 3
- [12] Zijian Dong, Chen Guo, Jie Song, Xu Chen, Andreas Geiger, and Otmar Hilliges. Pina: Learning a personalized implicit neural avatar from a single rgb-d video sequence. *arXiv*, 2022. 3
- [13] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. In *SIGGRAPH Asia 2022 Conference Papers*, 2022. 3
- [14] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps, 2021. 3
- [15] Chen Geng, Sida Peng, Zhen Xu, Hujun Bao, and Xiaowei Zhou. Learning neural volumetric representations of dynamic humans in minutes. In *CVPR*, 2023. 1, 2, 3, 4, 6, 7, 8, 12, 13, 16, 17, 18
- [16] Marc Habermann, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Deepcap: Monocular human performance capture using weak supervision. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2020. 6, 12
- [17] Marc Habermann, Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Real-time deep dynamic characters. *ACM Transactions on Graphics*, 40(4), 2021. 6, 12
- [18] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. *ICCV*, 2021. 3
- [19] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *ACM Transactions on Graphics (TOG)*, 41(4):1–19, 2022. 3
- [20] Timothy Jeruzalski, David I. W. Levin, Alec Jacobson, Paul Lalonde, Mohammad Norouzi, and Andrea Tagliasacchi. Nilbs: Neural inverse linear blend skinning, 2020. 3
- [21] Boyi Jiang, Yang Hong, Hujun Bao, and Juyong Zhang. Self-recon: Self reconstruction your digital avatar from monocular video. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3
- [22] Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Instan-tavatar: Learning avatars from monocular video in 60 seconds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16922–16932, 2023. 1, 2, 4, 6, 7, 8, 12, 16, 17, 18
- [23] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. Neuman: Neural human radiance field from a single video. In *Proceedings of the European conference on computer vision (ECCV)*, 2022. 2
- [24] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. 2, 3, 4, 5, 6, 7, 8, 13, 16, 17, 19
- [25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. 6
- [26] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *ICCV*, 2019. 2
- [27] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. In *Computer Graphics Forum*, pages 29–43. Wiley Online Library, 2021. 3
- [28] Jonas Kulhanek and Torsten Sattler. Tetra-nerf: Representing neural radiance fields using tetrahedra. *arXiv preprint arXiv:2304.09987*, 2023. 3
- [29] Youngjoong Kwon, Lingjie Liu, Henry Fuchs, Marc Habermann, and Christian Theobalt. Deliffas: Deformable light fields for fast avatar synthesis. *Advances in Neural Information Processing Systems*, 2023. 3
- [30] Christoph Lassner and Michael Zollhofer. Pulsar: Efficient sphere-based neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1440–1449, 2021. 3

- [31] Ruilong Li, Julian Tanke, Minh Vo, Michael Zollhofer, Jürgen Gall, Angjoo Kanazawa, and Christoph Lassner. Tava: Template-free animatable volumetric actors. 2022. [3](#)
- [32] Ruilong Li, Hang Gao, Matthew Tancik, and Angjoo Kanazawa. Nerfacc: Efficient sampling accelerates nerfs. *arXiv preprint arXiv:2305.04966*, 2023. [3](#)
- [33] Siyou Lin, Hongwen Zhang, Zerong Zheng, Ruizhi Shao, and Yebin Liu. Learning implicit templates for point-based clothed human modeling. In *ECCV*, 2022. [3](#)
- [34] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–65:14, 2019. [3](#)
- [35] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, 2015. [3](#), [4](#)
- [36] Marko Mihajlovic, Yan Zhang, Michael J Black, and Siyu Tang. LEAP: Learning articulated occupancy of people. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. [3](#)
- [37] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. [3](#)
- [38] Claus Müller. *Spherical harmonics*. Springer, 2006. [3](#)
- [39] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, 2022. [2](#), [3](#), [7](#), [18](#)
- [40] Atsuhiko Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Unsupervised learning of efficient geometry-aware neural articulated representations. In *European Conference on Computer Vision*, 2022. [3](#)
- [41] Nikolay Patinkin, Dmitry Senushkin, Anna Vorontsova, and Anton Konushin. Neural global illumination for inverse rendering. In *2023 IEEE International Conference on Image Processing (ICIP)*, pages 1580–1584. IEEE, 2023. [3](#)
- [42] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. [3](#)
- [43] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, 2021. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#), [12](#), [13](#), [16](#), [17](#)
- [44] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *CVPR*, 2021. [2](#), [6](#), [7](#), [8](#), [12](#), [13](#), [16](#), [17](#)
- [45] Sida Peng, Shangzhan Zhang, Zhen Xu, Chen Geng, Boyi Jiang, Hujun Bao, and Xiaowei Zhou. Animatable neural implicit surfaces for creating avatars from videos. *arXiv preprint arXiv:2203.08133*, 2022. [1](#), [2](#), [4](#), [6](#), [7](#), [12](#), [13](#)
- [46] Ruslan Rakhimov, Andrei-Timotei Ardelean, Victor Lepitsky, and Evgeny Burnaev. Npbg++: Accelerating neural point-based graphics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15969–15979, 2022. [3](#)
- [47] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *ICCV*, pages 14335–14345, 2021. [3](#)
- [48] Darius Rückert, Linus Franke, and Marc Stamminger. Adop: Approximate differentiable one-pixel point rendering. *ACM Transactions on Graphics (ToG)*, 41(4):1–14, 2022. [3](#)
- [49] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. [2](#)
- [50] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *CVPR*, 2020. [2](#)
- [51] Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J. Black. SCANimate: Weakly supervised learning of skinned clothed avatar networks. In *Proceedings IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. [3](#)
- [52] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. [2](#), [3](#)
- [53] Tomas Simon, Hanbyul Joo, Iain Matthews, and Yaser Sheikh. Hand keypoint detection in single images using multiview bootstrapping. In *CVPR*, 2017. [2](#)
- [54] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proc. CVPR*, 2019. [3](#)
- [55] Noah Snavely, Steven M. Seitz, and Richard Szeliski. Photo tourism: Exploring photo collections in 3d. *ACM Trans. Graph.*, 25(3):835–846, 2006. [4](#)
- [56] Yu Sun, Qian Bao, Wu Liu, Yili Fu, Black Michael J., and Tao Mei. Monocular, One-stage, Regression of Multiple 3D People. In *ICCV*, 2021. [2](#)
- [57] Yu Sun, Wu Liu, Qian Bao, Yili Fu, Tao Mei, and Michael J Black. Putting People in their Place: Monocular Regression of 3D People in Depth. In *CVPR*, 2022.
- [58] Yu Sun, Qian Bao, Wu Liu, Tao Mei, and Michael J. Black. TRACE: 5D Temporal Regression of Avatars with Dynamic Cameras in 3D Environments. In *CVPR*, 2023. [2](#)
- [59] Garvita Tiwari, Nikolaos Sarafianos, Tony Tung, and Gerard Pons-Moll. Neural-gif: Neural generalized implicit functions for animating people in clothing. In *International Conference on Computer Vision (ICCV)*, 2021. [3](#)
- [60] Shaofei Wang, Katja Schwarz, Andreas Geiger, and Siyu Tang. Arah: Animatable volume rendering of articulated human sdf. In *European Conference on Computer Vision*, 2022. [3](#)
- [61] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *CVPR*, 2016. [2](#)

- [62] Chung-Yi Weng, Brian Curless, Pratul P. Srinivasan, Jonathan T. Barron, and Ira Kemelmacher-Shlizerman. HumanNeRF: Free-viewpoint rendering of moving people from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16210–16220, 2022. [2](#), [3](#), [6](#), [7](#)
- [63] Zeke Xie, Xindi Yang, Yujie Yang, Qi Sun, Yixiang Jiang, Haoran Wang, Yunfeng Cai, and Mingming Sun. S3im: Stochastic structural similarity and its unreasonable effectiveness for neural fields. In *International Conference on Computer Vision*, 2023. [6](#)
- [64] Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J. Black. ICON: Implicit Clothed humans Obtained from Normals. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13296–13306, 2022. [2](#), [4](#)
- [65] Yuliang Xiu, Jinlong Yang, Xu Cao, Dimitrios Tzionas, and Michael J. Black. ECON: Explicit Clothed humans Optimized via Normal integration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. [2](#), [4](#)
- [66] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields, 2021. [3](#)
- [67] Jianfeng Zhang, Zihang Jiang, Dingdong Yang, Hongyi Xu, Yichun Shi, Guoxian Song, Zhongcong Xu, Xinchao Wang, and Jiashi Feng. Avatargen: A 3d generative model for animatable human avatars. In *Arxiv*, 2022. [3](#)
- [68] Qiang Zhang, Seung-Hwan Baek, Szymon Rusinkiewicz, and Felix Heide. Differentiable point-based radiance fields for efficient view synthesis. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–12, 2022. [3](#)
- [69] 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. [7](#)
- [70] Ruiqi Zhang, Jie Chen, and Qiang Wang. Explicifying neural implicit fields for efficient dynamic human avatar modeling via a neural explicit surface. *arXiv preprint arXiv:2308.05112*, 2023. [1](#), [3](#)
- [71] Zechuan Zhang, Li Sun, Zongxin Yang, Ling Chen, and Yi Yang. Global-correlated 3d-decoupling transformer for clothed avatar reconstruction. *arXiv preprint arXiv:2309.13524*, 2023. [2](#), [4](#)
- [72] Zerong Zheng, Xiaochen Zhao, Hongwen Zhang, Boning Liu, and Yebin Liu. Avatarrex: Real-time expressive full-body avatars. *ACM Transactions on Graphics (TOG)*, 42(4), 2023. [1](#), [3](#)

Human101: Training 100+FPS Human Gaussians in 100s from 1 View

Supplementary Material

A. Overview

Overview of the Supplementary Material:

- Implementation Details § B:
 - Conventions in Symbolic Operations. § B.1
 - Dataset. § B.2
 - Baseline Implementation Details. § B.3
 - Hyperparameters. § B.4
 - Network Structure. § B.5
 - Canonical Human Initialization. § B.6
 - Details of Triangular Face Rotation Matrices. § B.7
- Additional Quantitative Results § C:
 - Novel View Results. § C.1
 - Multi-view Results. § C.2
- Additional Qualitative Results § D:
 - Depth Visualization Results. § D.1
 - Novel Pose Results. § D.2
- More Experiments § E:
 - Memory Efficiency Comparison. § E.1
 - Ablation Study. § E.2
 - Failure Cases. § E.3
- Downstream Applications § F:
 - Composite Scene Rendering Results. § F.1
- More Discussions § G:
 - Data Preprocessing Technique. § G.1
 - Limitations. § G.2
 - Ethics Considerations. § G.3

B. Implementation Details

B.1. Conventions in Symbolic Operations

In our work, the rotation operations involve various types of rotational quantities (such as rotation matrices and quaternions). For simplicity, we represent these rotation operations in the format of “multiplication” in the main text. Here, we detail this representation more concretely:

For Gaussian rotation r_i , when optimizing r_i , it is considered a quaternion. While rotating it by the triangular face rotation matrix R_i , we first convert R_i into a unit quaternion and express this process using quaternion multiplication. Thus, this operation is denoted as:

$$r'_i = \text{quat_multi}(\text{rotmat_to_quat}(R_i), r_i) \quad (18)$$

When applying rotation in the refinement module, the predicted Δr_i is a quaternion. Therefore, this rotation is expanded as:

$$r''_i = \text{quat_multi}(\Delta r_i, r'_i) \quad (19)$$

For spherical harmonics, modifying spherical coefficients directly is not efficient. A more effective approach is to inversely rotate the view direction d_i . Specifically, we first calculate the direction from the camera center P_c to the final Gaussian position x''_i as the view direction. Then, we apply the inverse rotation transformation to view directions as the input for SH evaluation. Specifically, we have:

$$d_i = x''_i - P_c \quad (20)$$

$$d'_i = \text{SH_Rot}(R_i, d_i) \quad (21)$$

$$d''_i = \text{SH_Rot}(\text{quat_to_rotmat}(\Delta r_i), d'_i) \quad (22)$$

The function SH_Rot takes a rotation matrix and a view direction as input, returning a rotated view direction:

$$\text{SH_Rot}(R, d) = R^{-1}d = R^T d \quad (23)$$

B.2. Dataset

ZJU-MoCap. For ZJU-MoCap Dataset [44], we choose 6 subjects (377, 386, 387, 392, 393, 394) for evaluation. Because other subjects tend to not appear on the full side in a single fixed view. And following [15], we use camera 04 for training and other views for testing. Due to the low quality of the images in camera 03 for Subject 377 and Subject 392, we filter out these two views.

MonoCap. MonoCap is re-collected by [45], with Lan & Marc 1024 × 1024 resolution, selected from DeepCap dataset [16] and olek & vlad 1295 × 940 resolution selected from [17]. For better comparison, we show the FPS results of Lan and Marc. The DeepCap dataset [16] and DynaCap dataset [17] are only granted for non-commercial academic purposes. They prohibit the redistribution of that data. The users should also sign a license. More frame-selecting details are illustrated in Tab. 4.

B.3. Baseline Implementation Details

For Neural Body [44], Animatable NeRF [43], and AnimatableSDF [45], we utilized the results released in [15]. We also tested their rendering speeds by inferring with their pre-trained models on the same device using a single RTX 3090 GPU.

The work presented in [22] did not have an implementation for the ZJU-MoCap and MonoCap Datasets due to their slightly varied SMPL definition. Consequently, we adjusted the deformer in [22] to match the SMPL vertices. It’s important to note that [22] is designed specifically for monocular datasets, and it refines the SMPL parameters before metric evaluation. For a fair comparison, we adhered

Dataset	subject	Training view	Testing view index	Start Frame	End Frame	Frame Interval
ZJU-MoCap [44]	386, 387, 393, 394	4	Remaining	0	500	5
	377, 392	4	Remaining except 3	0	500	5
	Lan	0	Remaining	620	1120	5
MonoCap [43]	Marc	0	Remaining	35000	35500	5
	Olek	44	0, 5, 10, 15, 20, 25, 30, 35, 40, 45, 49	12300	12800	5
	Vlad	66	0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100	15275	15775	5

Table 4. **Dataset settings** (§ B.2).

to the same SMPL and camera parameters provided by the ZJU-MoCap Dataset and MonoCap Dataset, as with other baseline methods [15, 43–45], and chose not to refine the SMPL parameters before evaluation.

Regarding the 3D Gaussian Splatting [24], COLMAP could not determine valid camera parameters due to the input of monocular fixed-view video frames. As a solution, we opted to use the SMPL vertices from the initial frame as the input point cloud positions and designated the point cloud colors as white. Given that [24] is primarily a static multi-view 3D reconstruction method, achieving convergence in our setup proved challenging. Hence, we present the outcomes at 30k iterations, consistent with its original settings.

B.4. Hyperparameters

We experimentally fine-tuned our model employing a set of hyperparameters tailored for optimal performance. Regarding the spherical harmonics, we employed third-degree spherical harmonics for their balance of computational efficiency and representational fidelity. Uniquely, we increment the degree of spherical harmonics every 500 iterations, culminating at a maximum degree of three. For the learnable MLP component, we set the learning rate at 2×10^{-3} . During the optimization of Gaussians, we implemented an opacity reset at every 1500 iterations to refine transparency values.

B.5. Network Structure

To compensate for the rigid position and rotation using the learnable MLP, we employ straightforward linear layers, featuring a total of 5 layers with $n_{\text{hidden.dim}} = 64$. ReLU serves as the activation function between these layers, while no activation function is applied to the output. For SMPL parameters, we use a simple linear layer to compress its feature dimension. We use Positional Encoding with a frequency of 10. Fig. 7 demonstrates the structure of the linear networks.

B.6. Canonical Human Initialization

In the initialization process, we use an algorithmic approach instead of manually selecting four photos. Our objective is to select four images where the person’s angles on each are approximately 90 degrees apart. Additionally, it is preferable that the person’s pose in these images closely resembles the canonical pose. This ensures minimal accuracy loss when deforming the point cloud estimated by econ into the canonical

pose. To achieve this, we undertake the following steps:

1. Identify suitable image pairs. We traverse the dataset’s frames and for each frame index in frame index T , we maintain a set C_i , C_i records all frame indices whose angle δ with frame index i is between 80-100 degrees. The formula is as follows:

$$C_i = \{j \mid 80 \leq \delta_{ij} \leq 100, \forall j \neq i \text{ and } j > i\}, \forall i \in T \quad (24)$$

The angle δ_{ij} between frames i and j is derived by calculating the difference in angles of the global rotation matrices R_{global} from the SMPL parameters of the two frames. The formula is as follows:

$$R_{\text{diff } ij} = R_{\text{global } i}^{-1} \cdot R_{\text{global } j} \quad (25)$$

$$\delta_{ij} = \text{as_euler}(R_{\text{diff } ij}) \quad (26)$$

2. Select a suitable group of frames. The second part involves identifying a set of four images that meet the criteria, executed through a four-level nested loop. Initially, frame i is selected, followed by choosing j from the set C_i . Subsequently, k is selected from j ’s set C_j , and l from C_k . For each selected group of frames (i, j, k, l) , the algorithm first checks if the angular difference between every two frames exceeds 80 degrees. Then, it computes the distance between the pose’s joint positions in these images and the joint positions of the canonical pose. Finally, the group of frames with the smallest distance is selected. The process is shown in Algorithm 1.

B.7. Details of Triangular Face Rotation Matrices

The process of computing rotation matrices involves two main steps: first, determining the orthonormal basis vectors (e_{can} and e_{ob}) that describe the orientation of each triangular facet of the SMPL model in the canonical and target poses, respectively; second, constructing the rotation matrix from these basis vectors.

For each triangular facet f constituted by vertices A , B , and C , and edges AB , AC , and BC , we define the first unit direction vector as:

$$\vec{a} = \frac{\vec{AB}}{\|\vec{AB}\|} \quad (27)$$

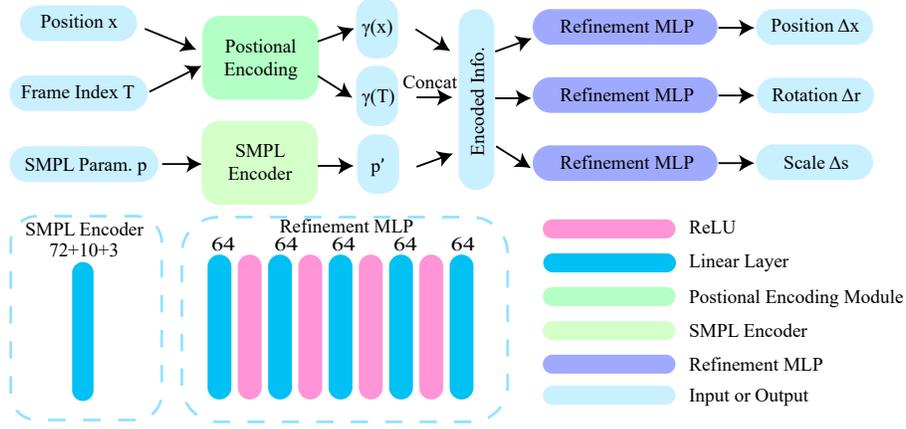


Figure 7. **Network structure** (§ B.5). This diagram presents the network architecture for refining the attributes of Gaussians. The input position x , frame index T , and SMPL parameters p are first processed through positional encoding and an SMPL encoder, respectively. The encoded information $\gamma(x)$, $\gamma(T)$, and p' are then concatenated and passed through a series of refinement MLPs to produce adjustments in position Δx , rotation Δr , and scale Δs . Each refinement MLP is composed of linear layers and employs ReLU activation functions.

Algorithm 1: Frame Selection (§ B.6)

Data: Sets T, C
Result: Best frame indices I_{best} and minimum distance d_{min}

```

 $d_{min} \leftarrow \infty$ 
 $I_{best} \leftarrow \emptyset$ 
for  $i \in T$  do
  for  $j \in C_i$  do
    for  $k \in C_j$  do
      for  $l \in C_k$  do
        if  $\delta_{ik} > 80^\circ$  and  $\delta_{il} > 80^\circ$  and  $\delta_{jl} > 80^\circ$  then
           $d \leftarrow$ 
            distance(pose( $i, j, k, l$ ), canonical pose)
          if  $d < d_{min}$  then
             $d_{min} \leftarrow d$ 
             $I_{best} \leftarrow \{i, j, k, l\}$ 
          end if
        end if
      end for
    end for
  end for
end for

```

Then, we use the normal of the triangular plane as the second unit direction vector:

$$\vec{b} = \frac{\vec{AB} \times \vec{AC}}{\|\vec{AB} \times \vec{AC}\|} \quad (28)$$

Subsequently, the third direction vector is derived from the

cross-product of the first two unit vectors:

$$\vec{c} = \vec{a} \times \vec{b} \quad (29)$$

Combining these vectors, we obtain the orthonormal basis for the triangular facet:

$$e = (\vec{a}, \vec{b}, \vec{c}) \quad (30)$$

Having acquired the orthonormal bases in both canonical and observation spaces, the triangular face rotation matrix is computed as:

$$R_f = e_{can} e_{ob}^\top \quad (31)$$

C. Additional Quantitative Results

C.1. Novel View Results

For the novel view setting, Tab. 5 and Tab. 6 show our results separately for resolution 512×512 and 1024×1024 .

C.2. Multi-view Results

While our model was not specifically designed for multi-view training data, we have conducted tests on the ZJU-MoCap Dataset to assess its performance in such scenarios. The results, as depicted in Fig. 9, demonstrate the model's capability to handle multi-view inputs.

D. Additional Qualitative Results

D.1. Depth Visualization

As shown in Fig. 8, our method, with its explicit representation, achieves a superior depth representation. This illustrates the advantages of our approach in terms of geometric accuracy.

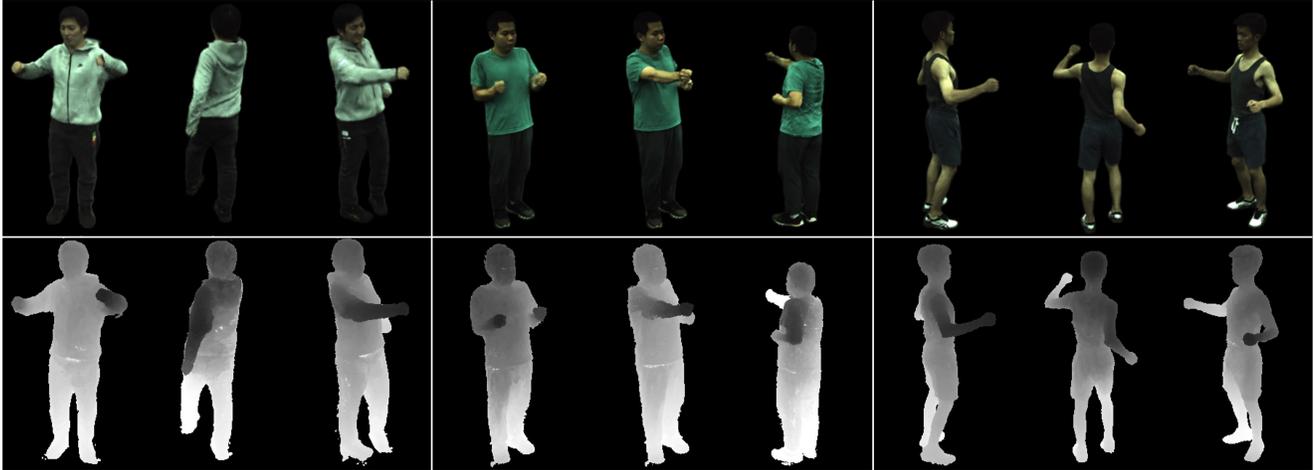


Figure 8. Depth visualization results (§ D.1).

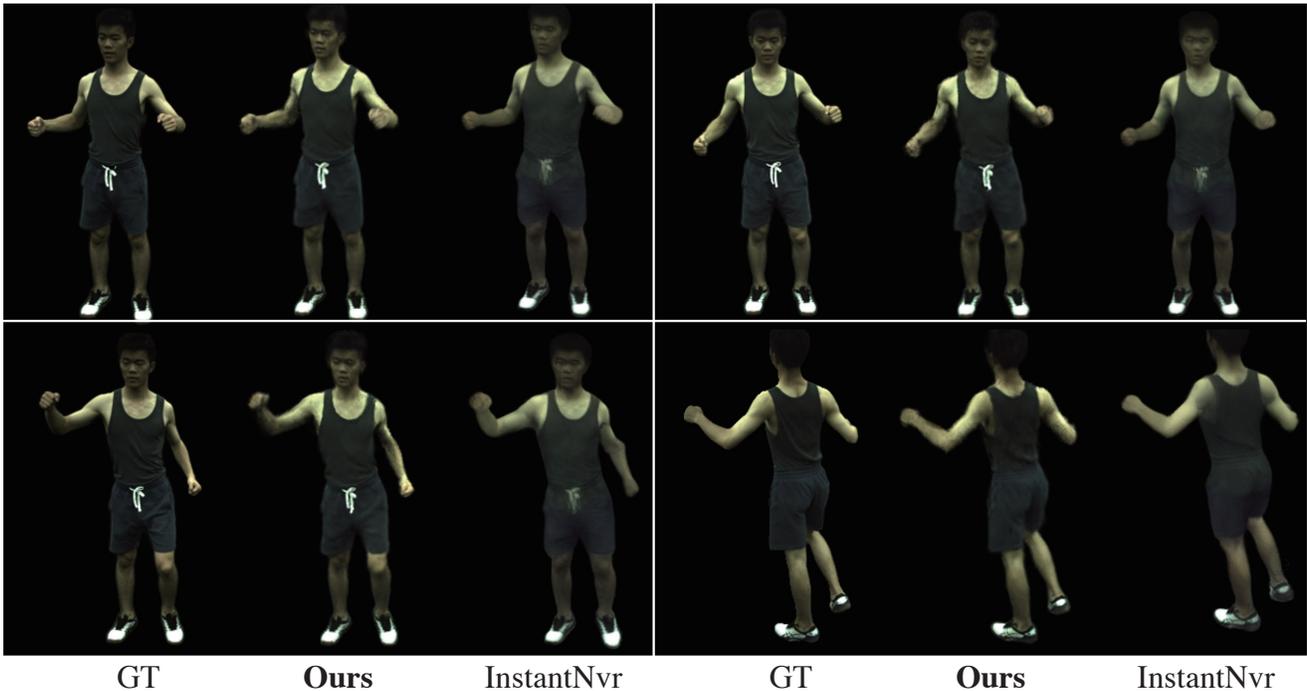


Figure 9. Multi view results. Qualitative results of methods trained with 4 views on the Sequence 377 of the ZJU-MoCap dataset.

D.2. Novel Pose Results

The results of our model trained on Subject 377 for unseen poses are shown in Fig. 10. Compared to the outcomes from InstantNvr, our results are less prone to artifacts and unnatural limb distortions. Simultaneously, our color reproduction is closer to the ground truth, with more preserved details in image brightness.

E. Additional Experiments

E.1. Memory Efficiency Comparison

To assess the efficiency of our model, we compared its resource consumption during the inference process with recent works in the field. In our comparison, we focused on three key metrics: training time GPU memory consumption (“Train Memory”), inference GPU memory consumption (“Infer Memory”), and disk space required for storing the model checkpoints (“Model Size”). In our work, the model

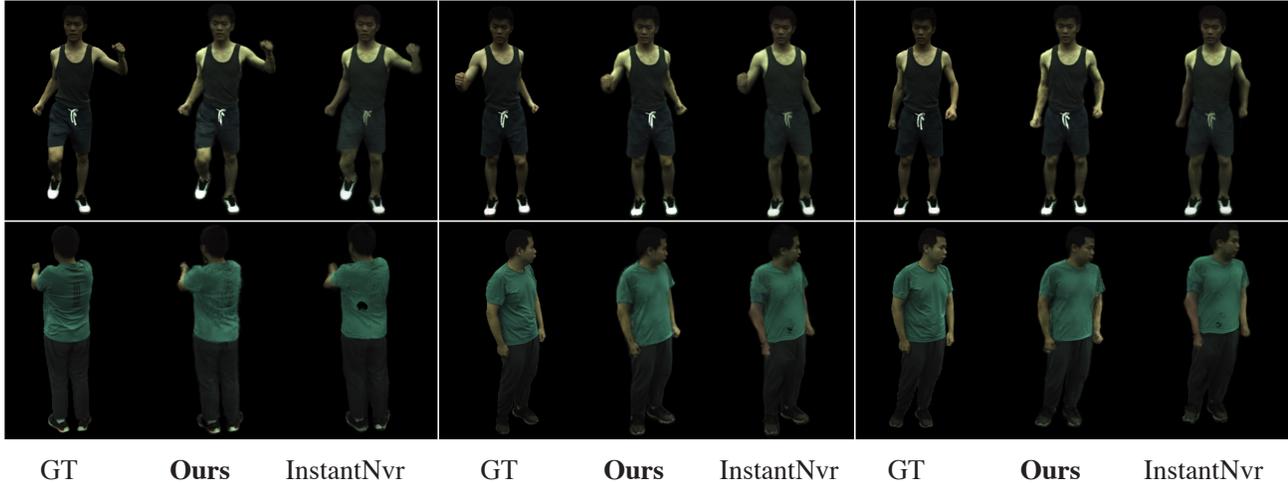


Figure 10. **Novel pose results.** (§ D.2) We show the results of unseen poses of our model and [15]. Results show that our model is less likely to produce artifacts or holes in unseen pose synthesis.

ZJU-MoCap [44]									
Method	Training Time	PSNR	SSIM	LPIPS*	FPS	PSNR	SSIM	LPIPS*	FPS
Subject		377				386			
3D GS [24]	5min	26.17	0.949	60.96	156	30.17	0.951	51.81	156
InstantNvr [15]	5min	31.69	0.981	32.04	1.53	33.16	0.979	38.67	1.53
InstartAvatar [22]	5min	29.90	0.961	49.00	8.75	30.67	0.917	111.5	8.75
Ours	100s	32.18	0.977	24.65	104	33.94	0.972	36.03	104
Ours	5min	32.02	0.976	21.35	104	33.78	0.969	33.73	104
Subject		387				392			
3D GS [24]	5min	24.56	0.922	80.61	156	26.72	0.932	79.61	156
InstantNvr [15]	5min	27.73	0.961	55.90	1.53	31.81	0.973	39.25	1.53
InstartAvatar [22]	5min	27.49	0.928	86.30	8.75	29.39	0.934	96.90	8.75
Ours	100s	28.32	0.956	47.76	104	32.22	0.966	41.89	104
Ours	5min	28.26	0.956	44.57	104	32.11	0.967	39.23	104
Subject		393				394			
3D GS [24]	5min	25.01	0.923	85.80	156	26.79	0.932	71.38	156
InstantNvr [15]	5min	29.46	0.964	46.68	1.53	31.26	0.969	39.89	1.53
InstartAvatar [22]	5min	28.17	0.931	86.60	8.75	29.64	0.943	64.20	8.75
Ours	100s	29.69	0.957	46.52	104	31.37	0.967	40.16	104
Ours	5min	29.52	0.956	44.15	104	31.25	0.968	36.86	104
MonoCap [43]									
Method	Training Time	PSNR	SSIM	LPIPS*	FPS	PSNR	SSIM	LPIPS*	FPS
Subject		Lan				Marc			
3D GS [24]	5min	28.76	0.970	30.19	156	30.16	0.972	30.76	156
InstantNvr [15]	5min	32.78	0.987	17.13	1.53	33.84	0.989	16.92	1.53
InstartAvatar [22]	5min	32.43	0.978	20.90	8.75	33.88	0.979	24.40	8.75
Ours	100s	32.63	0.982	14.21	104	34.84	0.983	19.21	104
Ours	5min	32.56	0.982	13.20	104	35.02	0.983	17.25	104
Subject		Olek				Vlad			
3D GS [24]	5min	28.32	0.961	45.24	147	23.13	0.961	51.16	147
InstantNvr [15]	5min	34.95	0.991	13.93	1.48	28.88	0.984	18.72	1.48
InstartAvatar [22]	5min	34.21	0.980	20.60	8.43	28.20	0.972	34.00	8.43
Ours	100s	34.31	0.982	15.07	101	28.96	0.977	23.56	101
Ours	5min	34.09	0.983	14.09	101	28.84	0.977	21.49	101

Table 5. 512×512 results of each subject on ZJU-MoCap dataset and Monocap dataset for **novel view synthesis** (§ C.1).

ZJU-MoCap [44]									
Method	Training Time	PSNR	SSIM	LPIPS*	FPS	PSNR	SSIM	LPIPS*	FPS
Subject		377				386			
3D GS [24]	5min	26.03	0.957	50.22	51.3	30.17	0.958	47.97	51.3
InstantNvr [15]	5min	31.69	0.981	32.04	0.5	33.16	0.979	38.67	0.5
InstartAvatar [22]	5min	27.74	0.933	87.91	3.83	28.81	0.916	97.72	3.83
Ours	100s	31.76	0.977	30.27	68	33.66	0.973	37.30	68
Ours	5min	31.64	0.976	27.99	68	33.42	0.973	36.03	68
Subject		387				392			
3D GS [24]	5min	24.57	0.931	64.75	51.3	26.69	0.9432	60.72	51.3
InstantNvr [15]	5min	27.93	0.968	49.11	0.5	31.89	0.977	42.49	0.5
InstartAvatar [22]	5min	26.15	0.890	107.7	3.83	27.98	0.9052	106.9	3.83
Ours	100s	27.95	0.959	47.56	68	31.97	0.970	41.65	68
Ours	5min	28.02	0.960	46.03	68	31.86	0.969	40.83	68
Subject		393				394			
3D GS [24]	5min	24.97	0.932	67.65	51.3	26.72	0.941	58.07	51.3
InstantNvr [15]	5min	29.32	0.969	48.36	0.5	31.36	0.968	39.58	0.5
InstartAvatar [22]	5min	27.43	0.899	102.6	3.83	28.62	0.926	81.20	3.83
Ours	100s	29.52	0.961	46.08	68	31.10	0.964	41.39	68
Ours	5min	29.42	0.960	44.64	68	31.04	0.963	40.07	68
MonoCap [43]									
Method	Training Time	PSNR	SSIM	LPIPS*	FPS	PSNR	SSIM	LPIPS*	FPS
Subject		Lan				Marc			
3D GS [24]	5min	28.44	0.974	25.95	51.3	30.13	0.9762	26.66	51.3
InstantNvr [15]	5min	32.61	0.988	12.73	0.5	33.76	0.989	17.01	0.5
InstartAvatar [22]	5min	32.89	0.982	17.30	3.83	33.72	0.982	21.81	3.83
Ours	100s	31.77	0.982	16.38	68	34.43	0.984	20.29	68
Ours	5min	31.72	0.982	15.55	68	34.56	0.985	18.96	68
Subject		Olek				Vlad			
3D GS [24]	5min	28.34	0.966	33.12	49.6	23.14	0.962	51.73	49.6
InstantNvr [15]	5min	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM
InstartAvatar [22]	5min	34.10	0.983	18.10	3.42	28.27	0.967	42.60	3.42
Ours	100s	34.04	0.984	16.19	63	28.53	0.979	20.37	63
Ours	5min	33.85	0.983	15.32	63	28.40	0.980	19.11	63

Table 6. 1024×1024 results of each subject on ZJU-MoCap dataset and Monocap dataset for **novel view synthesis** (§ C.1).

Method	Training Time	PSNR	SSIM	LPIPS*	FPS
NeuralBody	~10hours	32.99	0.983	26.8	3.5
HumanNeRF	~10hours	32.28	0.982	19.6	0.36
AnimatableNeRF	~10hours	32.31	0.980	32.2	2.1
AnimatableSDF	~10hours	32.63	0.983	32.0	1.3
InstantNvr	~13mins	32.55	0.981	26.5	1.5
Ours	~ 5mins	33.90	0.981	24.92	104

Table 7. **Multi-view results comparison** (§ C.2). Though our model is not designed for multi-view settings, we do experiments on 4 views of Sequence 377. Our model produces remarkable results using much less time while achieving good visual quality and evaluation metrics and much higher FPS.

size is computed by the sum of point cloud size and MLP checkpoint size.

As illustrated in Tab. 8, Human101 demonstrates notable memory efficiency compared to prior methods [15, 22]. During training, our model employs a strategy aligned with downstream applications, opting for direct run-time querying of the neural network for rendering. This decision not only conserves space but also facilitates real-time rendering capabilities, as opposed to pre-storing query results which would increase storage requirements and impede real-time

performance.

E.2. Ablation Study

Sparse Input Frames. Our model consistently delivers impressive results even with fewer input video frames. For Subject 377, as detailed in Tab. 9, we showcase our performance metrics for varying frame counts, specifically at 250, 100, 50, and 25 frames.

Positional Encoding. In our experiments, we explored different positional encoding strategies for Gaussian positions,

Method	Resolution	Train Memory	Infer Memory	Model Size
InstantAvatar [22]	512×512	4542M	3964M	151M
	1024×1024	4542M	4020M	151M
	642×470	4516M	3966M	151M
	1285×940	4654M	4038M	151M
InstantNvr [15]	512×512	19132M	4816M	3.2G
	1024×1024	23320M	4816M	3.2G
	642×470	21868M	7660M	3.2G
	1285×940	OOM	-	-
Ours	512×512	1878M	956M	12M+292K
	1024×1024	4146M	1842M	12M+292K
	642×470	1932M	1008M	12M+292K
	1285×940	4726M	2038M	12M+292K

Table 8. **Memory efficiency comparison (§ E.1).** For all resolutions in the dataset, we test the memory efficiency by Training GPU memory consumption (“Train Memory”), Inference GPU memory consumption (“Infer Memory”), and the size of the checkpoints (“Model Size”). Results demonstrate that our model utilizes much less GPU memory and disk usage than [15] while maintaining comparable or better visual quality. Note: when inferring, we don’t precompute and save Gaussians in target space while we choose to query the network for each frame. This methodological choice significantly reduces the storage requirements and makes it possible for Human101 to apply for more flexible use cases.

Frame Num	PSNR	SSIM	LPIPS*
25	31.66	0.974	24.78
50	32.00	0.975	22.26
100	32.18	0.977	21.32
250	32.17	0.977	19.17

Table 9. **Ablation study on frame number (§ E.2).** Our model still maintains good visual quality using sparse frame inputs even with only 25 images to train.

Method	PSNR	SSIM	LPIPS*
NoEnc	32.13	0.976	24.47
GridEnc	31.99	0.975	29.47
PE(Ours)	32.18	0.977	21.32

Table 10. **Ablation study on encoding method (§ E.2).** The results demonstrate that the positional encoding method produces better quality than no encoding (“NoEnc”) and grid-encoding (“GridEnc”).

specifically comparing Instant-ngp [39]’s grid encoding against the traditional sine and cosine positional encoding. While grid encoding can experimentally accelerate the fitting process on the training frames, it also tends to make the model more susceptible to overfitting. Consequently, as demonstrated in Tab. 10, this results in suboptimal performance on novel view test frames.

Degree of Spherical Harmonics. We have also performed ablation experiments to determine the optimal degree of spherical harmonics for our reconstruction task. As indicated by Tab. 11, increasing the degree of spherical harmonics leads to improved reconstruction quality. However, higher

Degree	PSNR	SSIM	LPIPS*
0	31.77	0.974	24.05
1	32.04	0.976	22.55
2	32.13	0.976	21.60
3(Ours)	32.18	0.977	21.32

Table 11. **Ablation study on the degree of spherical harmonics (§ E.2).** We evaluate the impact of the harmonics’ degree on the quality of reconstruction, with the degree of 3 (our chosen configuration) offering a trade-off between reconstruction detail and computational efficiency.

degrees bring a greater computational load. Consequently, we have chosen to adopt third-degree spherical harmonics for fitting in our final model, balancing accuracy with computational efficiency.

Converge Speed on different Initialization. Fig. 11 demonstrates that the choice of initialization method significantly impacts the model’s convergence speed. Furthermore, Fig. 12 shows different initialization strategies result in varying numbers of Gaussians at convergence. Generally, for the same scene, a larger number of Gaussians at convergence corresponds to richer reconstructed details. Since querying the MLP network is the more time-consuming factor during the inference phase, an increase in the number of Gaussians does not substantially affect the rendering FPS.

E.3. Failure Cases

Our model adeptly processes both monocular and multi-view video inputs, achieving high-fidelity reconstructions from sparse view inputs within a brief training duration. However, it is important to acknowledge the model’s limitations. In in-

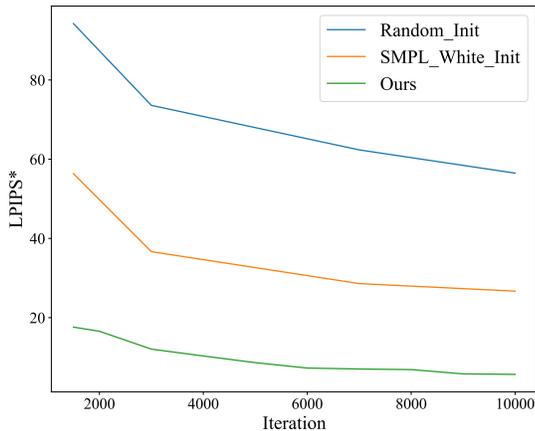


Figure 11. **Ablation study on convergence speed.** (§ E.2) We compare training view LPIPS results with the initialization method to be random initialization (“Random_Init”), bare SMPL with white color initialization (“SMPL_White_Init.”) and our Canonical Human Initialization method (“Ours”) separately.

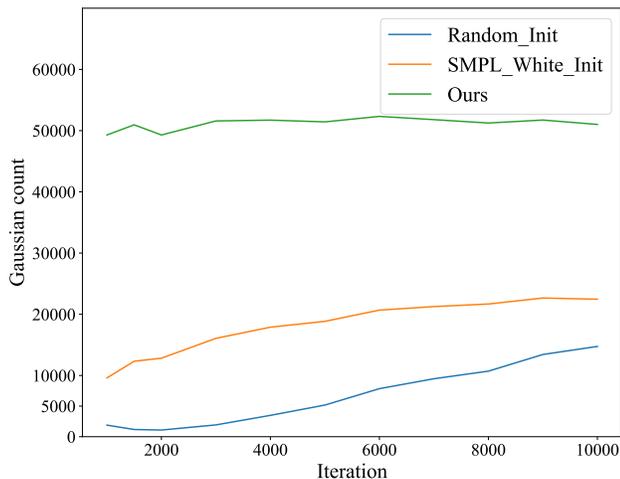


Figure 12. **Ablation study on Gaussian count** (§ E.2). We compare the number of Gaussians at different stages using various initialization methods. A superior initialization approach necessitates a greater number of Gaussians to represent the geometry more precisely, which generally yields better results.

stances where the input video fails to provide precise masks — for example, during intense movement where flowing hair carries unmasked background elements — this can result in visual artifacts, as depicted in Fig. 13.



Figure 13. **Failure case** (§ E.3). When dealing with intense movement where flowing hair carries unmasked background elements, our model may produce artifacts due to the complex human motion.



Figure 14. **Composite scene rendering** (§ F.1). We render the avatar integrated with the scene.

F. Application

F.1. Composite Scene Rendering

Rendering a human figure against a plain color background alone is not ideal for further downstream applications. Thanks to the explicit representation capability of 3D Gaussian Splatting (3D GS) [24], we can effortlessly segregate dynamic human figures from static scenes by explicitly splicing the Gaussians. This splicing process is natural and allows for the easy separation of static backgrounds and dynamic human elements.

As demonstrated in Fig. 14, this functionality facilitates downstream applications. In the example, the background and the human subject are trained separately and then composited during the rendering process. See supplementary videos for better results.

G. More Discussions

G.1. Discussions on Data Preprocessing Technique

Given that our task operates within a single-camera setting, we empirically observed during our experiments that, within fixed-view monocular videos, spherical harmonic coefficients tend to overfit to a singular direction. This leads to subpar generalization for free-view videos, resulting in numerous artifacts. To address this, we employed a data augmentation strategy that mimics a multi-camera environment.

With access to the SMPL parameters detailing the global rotation of the human subject, it's intuitive to keep the human orientation static while allowing the camera to orbit around the figure. This mimics a nearly equivalent process. Using this technique, we simulate varying camera viewpoints to render the dynamic human across different frames, markedly boosting the generalizability of the spherical harmonic functions.

However, this trick isn't devoid of limitations. In real-world scenarios, due to the diffuse reflection of light, we often perceive varying colors for the same object from different viewpoints. Our strategy overlooks this variance, providing an approximation that might not always align perfectly with real-world lighting conditions.

G.2. Limitations

While Human101 marks a significant advancement in dynamic human reconstruction, it is not without its limitations:

- Dependency on SMPL parameter accuracy. Human101 is significantly affected by the accuracy of SMPL parameter estimation. Inaccurate parameters can introduce substantial noise, complicating the reconstruction process.
- Requirement for complete body visibility in training data. Our model achieves the best results when training data includes all body parts relevant to the task. Partial visibility, where some body parts are not fully captured, may lead to artifacts in the reconstructed model.

Addressing these limitations could involve integrating more comprehensive human body priors, providing a pathway for future enhancements to our framework.

G.3. Ethics Considerations

Ethical considerations, particularly around privacy, consent, and data security, are critical in the development and application of Human101. Ensuring informed consent for all participants and transparent communication about the project's capabilities and limitations is essential to respect privacy and avoid misrepresentation. Secure handling and storage of sensitive human data are paramount to prevent unauthorized access and misuse. Additionally, acknowledging the potential for misuse of this advanced technology, we emphasize the need for ethical guidelines to govern its responsible use. Our commitment is to uphold high ethical standards in all aspects of Human101, safeguarding the respectful and secure use of human data.

G.4. Broader Impact

The development of Human101 has significant implications across various domains. Its ability to rapidly reconstruct high-quality, realistic human figures from single-view videos holds immense potential in fields such as virtual reality, animation, and telepresence. This technology can enhance user

experiences in gaming, film production, and virtual meetings, offering more immersive and interactive environments. However, its potential misuse in creating deepfakes or violating privacy cannot be ignored. It's crucial to balance innovation with responsible use, ensuring that Human101 serves to benefit society while minimizing negative impacts. Ongoing dialogue and regulation are necessary to navigate the ethical challenges posed by such advanced technology. Overall, Human101 stands to make a substantial impact in advancing digital human modeling while prompting necessary discussions on technology's ethical use.