EVA-GAUSSIAN: 3D GAUSSIAN-BASED REAL-TIME HUMAN NOVEL VIEW SYNTHESIS UNDER DIVERSE CAMERA SETTINGS

Anonymous authors

000 001

002 003

004

006

017 018 019

021

023

027

031 032

034

039

040

041

042

043

044

045

046

Paper under double-blind review



024 Figure 1: Qualitative comparison of novel view synthesis on the THuman2.0 dataset, with the angle 025 between the stereo views being 72 degree and GT representing the ground truth. We compare our proposed EVA-Gaussian against the state-of-the-art approaches GPS-Gaussian (Zheng et al., 2024) 026 and ENeRF (Lin et al., 2022). The quantitative metrics of PSNR \uparrow , SSIM \uparrow , LPIPS \downarrow , and inference time demonstrate that our proposed method achieves superior reconstruction quality, while 028 enabling real-time reconstruction under sparse viewing conditions and high resolution settings. 029

ABSTRACT

The feed-forward based 3D Gaussian Splatting method has demonstrated exceptional capability in real-time human novel view synthesis. However, existing approaches are restricted to dense viewpoint settings, where camera view angles are less than 60 degrees. This limitation constrains their flexibility in free-viewpoint rendering across a wide range of camera view angle discrepancies. To address this limitation, we propose a real-time pipeline named EVA-Gaussian for 3D human novel view synthesis across diverse multi-view camera settings. Specifically, we first introduce an Efficient cross-View Attention (EVA) module to accurately estimate the position of each 3D Gaussian from the source images. Then, we integrate the source images with the estimated Gaussian position map to predict the attributes and feature embeddings of the 3D Gaussians. Moreover, we employ a recurrent feature refiner to correct artifacts caused by geometric errors in position estimation and enhance visual fidelity. To further improve synthesis quality, we incorporate a powerful anchor loss function for both 3D Gaussian attributes and human face landmarks. Experimental results on the THuman2.0 and THumansit datasets showcase the superiority of our EVA-Gaussian approach in rendering quality across diverse camera settings. Project page: https: //anonymousiclr2025.github.io/iclr2025/EVA-Gaussian.

048 051

052

INTRODUCTION 1

3D reconstruction and novel view synthesis have long been fundamental yet complex tasks in visual data representation and computer vision. Recent advancements in fast 3D reconstruction and novel view synthesis for humans have shown immense potential in applications such as holographic
communication, real-time teaching, and augmented/virtual reality (AR/VR), where time efficiency
is critical for user experience and downstream processing. Nonetheless, existing methods either rely
on dense input views and precise templates as prior knowledge (Qian et al., 2024; Lei et al., 2024;
Hu et al., 2024; Wen et al., 2024; Kocabas et al., 2024; Kwon et al., 2024) or are restricted to specific
camera poses (Zheng et al., 2024; Tu et al., 2024). None of these approaches have fully developed
a pipeline for real-time human reconstruction under diverse, especially sparse, camera viewpoints,
which remains a significant challenge.

062 In recent years, Neural Radiance Fields (NeRFs) (Mildenhall et al., 2021) have emerged as a promis-063 ing technique for 3D reconstruction. These models employ neural networks to predict the color and 064 density of sampled 3D points along camera rays and aggregate these predictions to synthesize novel images with high fidelity. Despite their effectiveness, NeRFs suffer from substantial time con-065 sumption during both the training and rendering phases. Although various advancements, such as 066 multi-resolution hash encoding (Müller et al., 2022) and feed-forward neural scene prediction (Yu 067 et al., 2021a; Xu et al., 2024), have been made to reduce the time for reconstruction and rendering, 068 the achievable speeds remain insufficient for real-time applications. 069

More recently, 3D Gaussian Splatting (3DGS) (Kerbl et al., 2023) has been introduced as a solution 071 to this rendering bottleneck. 3DGS utilizes a set of discrete Gaussian representations to model complex 3D scenes and leverages the α -blending technique to enable real-time novel view rendering. 072 In the field of 3D human avatar reconstruction, previous works (Qian et al., 2024; Lei et al., 2024; 073 Hu et al., 2024; Wen et al., 2024; Kocabas et al., 2024; Kwon et al., 2024) have employed 3DGS 074 as a representation for humans to achieve animatable full-body human avatar reconstruction. These 075 methods, however, rely on precise human templates as priors, and their training and reconstruction 076 processes can take from minutes to hours, which impedes their use in real-time applications such 077 as holographic communication. While a feed-forward human reconstruction method (Zheng et al., 2024) has achieved fast reconstruction and real-time rendering with stereo inputs, the stereo settings 079 and small angle change between camera viewpoints restrict the overall reconstruction quality under sparse camera settings and lead to sub-optimal performance with more than two input views.

081 To address these limitations and enable real-time 3D reconstruction of humans using various camera positions and different numbers of cameras, we propose EVA-Gaussian, a novel 3D Gaussian-based 083 pipeline for real-time human novel view synthesis. Our method attaches 3D Gaussians to the sur-084 face of human body through multi-view depth estimation and aligns their positions closely with 085 point cloud locations. A key innovation of our method is the introduction of an Efficient cross-View Attention (EVA) module for multi-view 3D Gaussian position estimation (see Sec. 4.2). Specif-087 ically, we employ a U-Net (Ronneberger et al., 2015) as the backbone and further use dedicated 880 window-embedded cross-view attention to infer multi-view position correspondences. This attention mechanism enables EVA-Gaussian to effectively process multiple inputs from cameras positioned 089 at various viewpoint angles, thereby ensuring robust performance across a wide range of viewing 090 angles, even under extremely sparse camera settings. Besides, we incorporate a Gaussian attribute 091 estimation module that takes the EVA output and the original RGB images as input to estimate the 092 remaining 3D Gaussian attributes (see Sec. 4.3). Furthermore, we embed an additional attribute, referred to as feature, into each Gaussian for further feature splatting and image quality refinement, 094 thereby mitigating the position estimation error introduced by the EVA module (see Sec. 4.4). In 095 addition, we employ an anchor loss to penalize the inconsistency between multi-view face land-096 marks, which achieves better supervision for human faces (see Sec. 4.5). We conduct extensive 097 experiments on the THuman2.0 (Yu et al., 2021b) and THumanSit (Zhang et al., 2023) datasets. The 098 results, as exemplified in Fig. 1, demonstrate that our proposed EVA-Gaussian outperforms existing feed-forward synthesis approaches in rendering quality, while enabling real-time reconstruction and rendering. Moreover, our approach generalizes well to settings with diverse numbers of cameras and 100 significant changes in camera viewpoint angles. In summary, our main contributions are as follows: 101

- 102
- 103 104 105

107

• We propose a novel pipeline for fast feed-forward 3D human reconstruction, called *EVA-Gaussian*, that comprises three main stages: 1) a multi-view 3D Gaussian position estimation stage, 2) a 3D Gaussian attributes estimation stage, and 3) a feature refinement stage.

• We introduce an EVA module to enhance multi-view correspondence retrieval, leading to improved 3D Gaussian position estimation and enhanced novel view synthesis under diverse view numbers and sparse camera settings.

- We employ a recurrent feature refiner that fuses splatted RGB images and feature maps to mitigate geometric artifacts caused by position estimation errors. Moreover, we incorporate an anchor loss that utilizes facial landmarks as anchor points to better supervise Gaussian position estimation, thereby enhancing the quality of synthesized novel view images.
 - Extensive experiments on THuman2.0 and THumansit demonstrate the effectiveness and superiority of our proposed pipeline over existing methods in terms of rendered novel view quality and inference speed, especially under sparse camera settings.
- 118 2 RELATED WORKS
- 119

109

110

111

112

113

114

115 116 117

120 **3DGS-based Human Reconstruction.** 3D Gaussians Splatting has recently emerged as an effective 121 technique for 3D human reconstruction. However, most previous works (Lei et al., 2024; Hu et al., 122 2024; Wen et al., 2024; Qian et al., 2024; Kocabas et al., 2024; Pan et al., 2024) bind 3D Gaussians 123 to a predefined human mesh model, such as SMPL (Loper et al., 2023) or SMPL-X (Pavlakos et al., 124 2019). This approach generates 3D Gaussians and human models in a canonical space and then 125 transforms them to match the target human pose using the predefined weights of the human model. This iterative binding process, however, is extremely time-consuming. Moreover, these methods 126 require human templates as inputs at each frame, which incurs extra computational cost and poten-127 tially misleads the reconstruction procedure due to the errors in pose estimation. These limitations 128 significantly hinder their applicability in real-world scenarios. 129

- 130 **Fast Generalizable 3D Reconstruction.** In the field of NeRF rendering, pixelNeRF (Yu et al., 131 2021a) pioneers the approach of predicting features per pixel from a single image in a feed-forward manner for 3D reconstruction. While subsequent works (Chen et al., 2021; Wang et al., 2021; Lin 132 et al., 2022) have followed this feed-forward NeRF pipeline, they still suffer from the extensive time 133 consumption of the NeRF rendering process. Besides, their reconstruction results are often unsatis-134 factory in sparse camera settings. The introduction of 3DGS has helped mitigate the rendering speed 135 issue of high-quality novel view synthesis methods. Notably, pixelSplat (Charatan et al., 2024) and 136 Splatter Image (Szymanowicz et al., 2024) are the first to combine the feed-forward inference and 137 3DGS, which predict 3D Gaussian attributes for each pixel and project them back to the 3D space for 138 novel view synthesis in a real-time manner. Nevertheless, these methods still struggle with inaccu-139 rate estimation of 3D Gaussian positions. MVSplat (Chen et al., 2024) and MVGaussian (Liu et al., 140 2024a) address this issue by leveraging cost-volume modules, thereby achieving better novel view 141 image quality. Moreover, latentSplat (Wewer et al., 2024) attaches a latent vector to each 3D Gaus-142 sian and uses this vector for novel view refinement through a diffusion decoder and generative loss, which significantly improves image quality on extrapolation views. Despite these advancements, 143 these approaches do not fully exploit prior knowledge about human images and camera settings, 144 which limits their performance on human reconstruction and novel view synthesis. In the field of 145 3D generation, LGM (Tang et al., 2024) and GSLRM (Zhang et al., 2024) employ a transformer-146 based feed-forward pipeline for predicting 3D Gaussian attributes from multi-view images in real-147 world scenarios, but these networks are often too complex for real-time reconstruction. Our work 148 utilizes a carefully designed memory-efficient attention mechanism, thereby enabling generalizable 149 3D reconstruction at high resolutions with a much lower temporal and computational cost.
- 150

The work most closely related to ours is GPS-Gaussian (Zheng et al., 2024), which proposes a stereo 151 matching network for 3D Gaussian position estimation and employs two 3-layer U-Nets to predict 152 3D Gaussian scales, rotations, and opacities. Although GPS-Gaussian has demonstrated the capa-153 bility for real-time human reconstruction and novel view synthesis, it suffers from severe distortions 154 under sparse camera settings and mismatch across multiple viewpoints. Subsequent works have 155 attempted to alleviate these issues. For instance, Tele-Aloha (Tu et al., 2024) introduces an image 156 blending and cascaded disparity estimation method for human reconstruction with four input views. 157 However, this approach is tailored to a specific system and struggles to generalize to sparser camera 158 settings. On the other hand, GHG (Kwon et al., 2024) achieves real-time 3D Gaussian-based human novel view synthesis in a feed-forward manner, but it requires additional human template priors, 159 thus inheriting the limitations associated with template-based human reconstruction methods. In 160 contrast, our proposed method eliminates the need for human templates and is specifically designed 161 to generalize effectively across various sparse camera settings.

162 3 PRELIMINARY

3D Gaussian Splatting (3DGS) uses a set of 3D Gaussian distributions to represent a 3D sence based on a sequence of multi-view RGB images (Kerbl et al., 2023). Each 3D Gaussian can be mathematically expressed as:

$$G(\boldsymbol{x}) = e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})},$$
(1)

where $\mu \in \mathbb{R}^3$ and $\Sigma \in \mathbb{R}^{3\times 3}$ denote the position vector and the covariance matrix, respectively. The covariance matrix Σ is constrained to be semi-positive definite, which is guaranteed by decomposing it into a rotation matrix $R \in \mathbb{R}^{3\times 3}$ and a scaling matrix $S \in \mathbb{R}^{3\times 3}$ as:

$$\Sigma = RSS^T R^T, \tag{2}$$

where R is further represented by a quaternion $q \in \mathbb{R}^4$, and S is further represented by a scaling vector $s \in \mathbb{R}^3$. In addition, to facilitate rendering from a specific camera perspective, 3DGS leverages a view transformation W to transfer the 3D Gaussians from the world space to the camera space, and employs a projective transformation J to approximate the projection of these 3D Gaussians onto the 2D image plane.

The final rendered color for each pixel is computed using an α -blending function that accumulates the contributions of all the projected 3D Gaussians on the pixel:

$$c_{\text{pixel}} = \sum_{i=1}^{N} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j),$$
 (3)

where N denotes the number of Gaussians, α_i is the product of the *i*-th projected 2D Gaussian and its optimized opacity value $o_i \in [0, 1]$, and $c_i \in \mathbb{R}^3$ denotes the color attribute of the Gaussian, which is represented using spherical harmonics coefficients to enable view-dependent color modeling.

191

192

182

183

185

186

167 168

173

4 METHODOLOGY

4.1 OVERVIEW

In this paper, we focus on fast human 3D reconstruction and novel view synthesis under diverse camera settings. Our objective is to reconstruct a 3D scene from a set of *n* sparse-view RGB images $\{I_i\}_{i=1}^n, I_i \in \mathbb{R}^{H \times W \times 3}$, captured from different viewpoints surrounding a human subject, where the angle between any two adjacent camera views is denoted by Δ , and synthesize arbitrary novel view images at any camera position in real time. To achieve this, we propose *EVA-Gaussian*, a method that utilizes deep neural networks and 3D Gaussian Splatting to enhance novel image quality while achieving real-time reconstruction.

200 Specifically, we employ 3DGS to represent each source image I_i as a set of 3D Gaussians. Each 201 pixel in the foreground corresponds to a unique 3D Gaussian. We use U_i to denote the number of Gaussians for source image *i*. The proposed EVA-Gaussian predicts the positions and attributes of 202 3D Gaussians in the form of attribute maps $\{M_i\}_{i=1}^n = \{P_i, O_i, S_i, Q_i, F_i\}_{i=1}^n$ from the image set 203 $\{I_i\}_{i=1}^n$, where P_i, O_i, S_i, Q_i , and F_i denote the attribute maps for Gaussian positions, opacities, 204 scales, quaternions, and features of source image i, respectively. Notably, in the feature map F_i = 205 $\{f_i^u\}_{u=1}^{U_i}$, each element $f_i^u \in \mathbb{R}^{32}$ serves as a new attribute associated with each 3D Gaussian, 206 which will be used later in Sec. 4.4 to remove artifacts caused by geometric errors in $\{P_i\}_{i=1}^n$. 207 Mathematically, the procedure of EVA-Gaussian is expressed as: 208

209

$$\{M_i\}_{i=1}^n = \mathcal{D}_{\theta}(\{I_i\}_{i=1}^n), \tag{4}$$

210 where θ denotes the parameters of the learnable neural network.

The framework of EVA-Gaussian is depicted in Fig. 2. EVA-Gaussian splits the process of predicting Gaussian maps into three stages. In the first stage, it employs a U-Net architecture with an Efficient cross-View Attention module (EVA) to obtain enhanced multi-view predictions of the 3D Gaussian position maps $\{P_i\}_{i=1}^n$, as elaborated in Sec. 4.2. In the second stage, a Gaussian attribute

prediction network, detailed in Sec. 4.3, takes the predicted 3D Gaussian position maps $\{P_i\}_{i=1}^n$



Figure 2: **Framework of EVA-Gaussian**. EVA-Gaussian takes sparse-view images captured around a human subject as input and performs three key stages: (1) estimating the positions of 3D Gaussians, (2) inferring the remaining attributes (i.e., opacities, scales, quaternions, and features) of these Gaussians, and (3) refining the output image in a recurrent manner.

and the original RGB images $\{I_i\}_{i=1}^n$ as input to estimate the remaining attributes of 3D Gaussians. The predicted 3D Gaussians from all source images are then aggregated to render target views using differential rasterization (Kerbl et al., 2023). In the final stage, the rendered RGB image \hat{I}^0 and its corresponding feature map F_{novel} are fused for further refinement using the network described in Sec. 4.4. In addition, an anchor loss is introduced during the training stage to enhance the overall reconstruction quality, as depicted in Sec. 4.5.

4.2 GAUSSIAN POSITION ESTIMATION244

253

The variations in depth across the surface of human body may appear minimal. However, these nuances are critically important, particularly in regions such as the face and hands that contain a wealth of semantic information. Even slight inaccuracies in depth estimation within these areas can lead to significant degradation in visual quality and fidelity. This underscores the necessity for precise estimation of 3D Gaussian positions to enable effective and high-fidelity human reconstruction.

To tackle this challenge, we employ a U-Net based architecture $\mathcal{D}_{\theta_1}^P$ to estimate the 3D Gaussian position maps $\{P_i\}_{i=1}^n$ from multi-view images $\{I_i\}_{i=1}^n$, which is expressed as:

$$[\mathbf{P}_{i}]_{i=1}^{n} = \mathcal{D}_{\theta_{1}}^{P}(\{\mathbf{I}_{i}\}_{i=1}^{n}).$$
(5)

To ensure accurate depth estimation across diverse camera angles or arbitrary input views, we pro-254 pose an EVA module, as illustrated in Fig. 3. This module is integrated into the three lowest resolu-255 tion layers of the U-Net basebone $\mathcal{D}^{P}_{\theta_1}$ to facilitate the multi-view correspondence retrieval and infor-256 mation exchange. We use j to denote the index of each of these three layers, with j = -1, j = -2, 257 and j = -3 representing the lowest, the second-lowest, and the third-lowest resolution layers, re-258 spectively. EVA takes multiple intermediate image features $E_i^j \in \mathbb{R}^{R^j \times C^j}, \forall i \in \{1, \dots, n\}, \forall j \in \{1, \dots, n\}, \forall j$ 259 $\{-1, -2, -3\}$, as input and outputs the corresponding enhanced image features \widetilde{E}_i^j , where R^j and 260 C^{j} denote the total number of pixels and the channel dimension of each pixel at layer j, respectively. 261 Before the execution of attention mechanisms, a learnable positional embedding γ is added to the 262 intermediate feature E_i^j to improve the understanding of image coordinates.

In contrast to other feed-forward scene reconstruction methods (Chen et al., 2024; Charatan et al., 2024; Liu et al., 2024a) that apply to a low resolution of 256×256, our approach aims for highquality human reconstruction at 1024 resolution. Given that the corresponding pixels from the reference views are located in adjacent locations only under human-centric camera settings, calculating attention scores across a concatenated multi-view image sequence or feature map for the entire image, as adopted by Chen et al. (2024), is highly inefficient. To improve efficiency, EVA computes cross-attention only within a local window, which is shifted by half the window size at each attention Position Embedding γ E_1^{j} E_1^{j} E_1^{j

Figure 3: Efficient cross-View Attention (EVA) module for 3D Gaussian position estimation. EVA takes multi-view image features as input, embeds them into window patches using a shifted algorithm, and performs cross-view attention between the features from different views.

iteration. This shifted window embedding allows for low computation complexity and better GPUmemory utilization, while maintaining high performance.

In the cross-view attention stage, the intermediate feature E_i^j is linearly transformed into query, key, and value matrices, denoted by Q_i^j , K_i^j , and V_i^j , respectively. For each E_i^j , we calculate cross-view attention with the key matrices \overline{K}^j fused from all the reference image features, excluding the target feature itself. That is, $\overline{K}_i^j = \mathcal{T}_i^j(K_1^j, \cdots, K_{i-1}^j, K_{i+1}^j, \cdots, K_n^j)$, where the fusion \mathcal{T}_i^j is achieved by using a fully connected neural network. Consequently, each attention map is derived from its associated query, the fused key, and its corresponding value as:

$$\boldsymbol{A}_{i}^{j} = \operatorname{softmax}(\boldsymbol{Q}_{i}^{j} \overline{\boldsymbol{K}}_{i}^{jT} / \sqrt{C^{j}}) \boldsymbol{V}_{i}^{j}, \tag{6}$$

where $A_i^j, \forall i \in \{1, 2, \dots, n\}$, denotes the resultant attention output.

Notably, when the scale of each Gaussian is sufficiently small, the 3D Gaussian position of a pixel aligns precisely with its corresponding value on the depth map. A detailed proof of this property is provided in Appendix C. Based on this observation, we train the position estimation network \mathcal{D}^{P}_{θ} to obtain the position maps $\{P_i\}_{i=1}^{n}$ with the mean squared error (MSE) loss function:

$$\mathcal{L}_{depth} = ||\boldsymbol{P}_i - \boldsymbol{P}_i^{gt}||_2, \tag{7}$$

where P_i^{gt} denotes the ground truth depth map.

4.3 GAUSSIAN ATTRIBUTE ESTIMATION

To complete the estimation of 3D Gaussian maps $\{M_i\}_{i=1}^n$, we employ a shallow U-Net $\mathcal{D}_{\theta_2}^A$ to estimate the remaining attributes O_i , S_i , Q_i , F_i . This network takes the estimated 3D Gaussian position maps $\{P_i\}_{i=1}^n$ from the first stage in Sec. 4.2 and the original RGB images $\{I_i\}_{i=1}^n$ as input, and outputs the 3D Gaussian attributes O_i , S_i , Q_i , F_i , which is expressed as:

$$[\boldsymbol{O}_i, \boldsymbol{S}_i, \boldsymbol{Q}_i, \boldsymbol{F}_i]_{i=1}^n = \mathcal{D}_{\boldsymbol{\theta}_2}^A(\{\boldsymbol{I}_i\}_{i=1}^n \oplus \{\boldsymbol{P}_i\}_{i=1}^n).$$
(8)

The resulting estimated 3D Gaussian maps $\{M_i\}_{i=1}^n = \{P_i, O_i, S_i, Q_i, F_i\}_{i=1}^n$ are then utilized to render novel views using the process described in Sec. 3. The network $\mathcal{D}_{\theta_2}^A$ is trained by using a combination of MSE loss and structural similarity index measure (SSIM) (Wang et al., 2004) loss between the rendered novel view image \hat{I}^0 and the ground truth I^{gt} as follows:

 $\mathcal{L}_{\text{render}} = ||\hat{\boldsymbol{I}}^0 - \boldsymbol{I}^{\text{gt}}||_2 + \lambda_{\text{render}} (1 - \text{SSIM}(\hat{\boldsymbol{I}}^0, \boldsymbol{I}^{\text{gt}})),$ (9)

319 where λ_{render} denotes the weighting factor for the SSIM loss.

{

321 4.4 FEATURE SPLATTING AND REFINEMENT

The 3D Gaussian position maps P_i estimated in Sec. 4.2 inevitably contain some degree of error, which may lead to distortions and artifacts in the rendered RGB images. To mitigate these issues,



270

271

272

273 274

275

276

277 278

279

284

295

297

298

299

300

301 302 303

304 305

306

311

312

317

318

320



Figure 4: Attribute regularization. We regularize the opacities and scales of Gaussians, as well as 333 the position mismatch among the Gaussians in the landmark collection. The optimization of position mismatch terminates when it falls below a specific tolerance. 334

we propose a post-splatting refinement method to correct the position estimates. Recent studies 336 (Berriel Martins & Civera, 2024) have demonstrated that feature vector representations can capture 337 scene information more effectively than spherical harmonics, resulting in significant improvements 338 in novel view synthesis, particularly in scenarios with limited overlapping views. Inspired by this 339 finding, we incorporate a feature vector, i.e., $f_i^u \in \mathbb{R}^{32}$ mentioned in Sec. 4.1, as an additional 340 attribute for each Gaussian to more precisely capture its spatial characteristics. 341

During the splatting process, we first aggregate the 3D Gaussians from all source views. Then, the 342 color values of these 3D Gaussians are rendered using Eq. 3. Concurrently, the feature values of the 343 3D Gaussians are splatted onto the image plane using a modified α -blending function as follows: 344

$$f_{\text{pixel}} = \sum_{j=1}^{N} f_j \alpha_j \prod_{l=1}^{j-1} (1 - \alpha_l),$$
(10)

where f_{pixel} is the feature vector for the corresponding pixel on the feature map of the novel view image \vec{F}_{novel} , f_j denotes the feature vector for the 3D Gaussian with the *j*-th greatest depth, and $N = \sum_{i=1}^{n} U_i$ is the total number of 3D Gaussians from all source views.

Moreover, we employ a carefully designed recurrent U-Net \mathcal{D}_{a}^{R} that takes both the RGB and feature images as input and project them onto the RGB space for the final output through L recurrent loops. This recurrent procedure is expressed as follows:

$$\hat{\boldsymbol{I}}^{l} = \mathcal{D}_{\boldsymbol{\theta}_{3}}^{R}(\hat{\boldsymbol{I}}^{l-1} \oplus \boldsymbol{F}_{\text{novel}}), \hat{\boldsymbol{I}}^{l} \in \mathbb{R}^{H \times W \times 3}, \boldsymbol{F}_{\text{novel}} \in \mathbb{R}^{H \times W \times 32}, l \in \{1 \cdots L\}.$$
(11)

Similar to the Gaussian attribute estimation, the loss function for supervising the final output is a combination of MSE loss and SSIM loss between refined image \hat{I}^L and ground truth I^{gt} as follows:

$$\mathcal{L}_{\text{refine}} = ||\hat{\boldsymbol{I}}^{L} - \boldsymbol{I}^{\text{gt}}||_{2} + \lambda_{\text{refine}} (1 - \text{SSIM}(\hat{\boldsymbol{I}}^{L}, \boldsymbol{I}^{\text{gt}})), \qquad (12)$$

where λ_{refine} denotes the weighting factor for the SSIM loss.

4.5 ATTRIBUTE REGULARIZATION

362 Since human faces are critical for identification and expression understanding, improving the reconstruction of human faces is much more important than that of other body parts. Previous works like 364 GPS-Gaussian (Zheng et al., 2024) treat the entire human body equally and ignore expressive infor-365 mation contained in human faces. Moreover, they fail to ensure the consistency between depth maps 366 and 3D Gaussian locations, resulting in sub-optimal reconstruction quality in the facial regions.

367 To address this issue, we propose a regularization term to enhance the overall reconstruction quality. 368 Specifically, our proposed anchor loss regularizes the scales and opacities of Gaussians to ensure 369 consistency between the geometry of predicted depth maps and the 3D Gaussian positions. It also 370 aligns the Gaussians from different views to force their locations to the same landmark. We adopt 371 MediaPipe (Lugaresi et al., 2019) to annotate human facial landmarks and compute the anchor loss 372 to regularize the 3D landmark Gaussian scales, opacities, and positions as follows:

$$\mathcal{L}_{\text{anchor}} = \sum_{i,j \in \mathbb{V}} \sum_{m_i \in \mathbb{M}_i, m_j \in \mathbb{M}_j} \max\left\{ ||\Pi^{-1}(\boldsymbol{m}_i, \boldsymbol{P}_i(\boldsymbol{m}_i)) - \Pi^{-1}(\boldsymbol{m}_j, \boldsymbol{P}_j(\boldsymbol{m}_j))||_2, t \right\}$$
$$+ \lambda_{\text{opacity}} \sum_{i=1}^{N} ||\boldsymbol{O}_i \log(\boldsymbol{O}_i)||_1 + \lambda_{\text{scale}} \sum_{i=1}^{N} ||\boldsymbol{S}_i||_2, \tag{13}$$

335

345 346 347

348

349

350 351

352

353

354 355

356

357 358

359 360

361

373 374 375

376
377
$$+ \lambda_{\text{opacity}} \sum_{i=1}^{N} ||\boldsymbol{O}_i \log(\boldsymbol{O}_i)||_1 + \lambda_{\text{scale}} \sum_{i=1}^{N} ||\boldsymbol{S}_i||_2,$$

PIPS↓ 0.0954	$\frac{\text{THumansit}}{0.8880}$	PSNR↑) LPIPS↓	THuman2.0 SSIM↑		$\Delta = 45^{\circ}$
PIPS↓).0954	SSIM↑ 1 0.8880	PSNR↑	LPIPS↓	SSIM↑	DCNDA	<u>_</u> +J
).0954	0.8880	02.21		oom	PSINK	
		23.31	0.0824	0.9156	25.19	pixelSplat
).0532	0.9223	24.97	0.0346	0.9515	28.05	MVSplat
).0297	0.9641	25.20	0.0283	0.9706	26.44	MVSGaussian
).0334	0.9567	27.06	0.0238	0.9696	29.62	ENeRF
).0251	0.9671	28.02	0.0224	0.9762	30.30	GPS-Gaussian
).0249	0.9696	29.16	0.0198	0.9782	31.11	EVA-Gaussian
).0334).0251).0249	0.9567 0.9671 0.9696	27.06 28.02 29.16	0.0238 0.0224 0.0198	0.9696 0.9762 0.9782	29.62 30.30 31.11	GPS-Gaussian EVA-Gaussian

Table 1: Comparison with feed-forward 3D reconstruction methods at a resolution of 256×256 . Better results are marked in a deeper color.

where $\{\mathbb{M}_i\}_{i=1}^n$ denotes the collection of all landmarks on the 2D image plane, \mathbb{V} denotes the collection of source views, and Π^{-1} represents the process of reprojection from 2D image to 3D space.

Since the MediaPipe landmark estimation is not perfectly accurate, we introduce a factor t to control the tolerance for mismatch errors. This tolerance facilitates the optimization by activating the loss only when the landmark distance exceeds t. Therefore, this approach optimizes the facial reconstruction loss to a sufficiently low level and avoids being misguided by potential errors in the MediaPipe estimation. This procedure is illustrated in Fig. 4.

By integrating the loss functions in the three stages, i.e., \mathcal{L}_{depth} , \mathcal{L}_{refine} , and the proposed regularization term \mathcal{L}_{anchor} , the overall training loss of the proposed EVA-Gaussian is given by:

$$\mathcal{L}_{\text{EVA-Gaussian}} = \mathcal{L}_{\text{depth}} + \lambda_1 \mathcal{L}_{\text{render}} + \lambda_2 \mathcal{L}_{\text{refine}} + \lambda_3 \mathcal{L}_{\text{anchor}}, \tag{14}$$

where λ_1 , λ_2 , and λ_3 are weights used to balance the different loss terms.

Since the 3D Gaussian position and attribute estimation stages can be executed within tens of milliseconds, and feature refinement is lightweight, taking less than ten milliseconds, EVA-Gaussian is capable of rapidly reconstructing 3D human subjects from a collection of RGB images and rendering novel views in a real-time manner.

403 404 405

413

394

395

396 397

380 381 382

384 385 386

5 EXPERIMENTS

406 407 5.1 Experiment Setup

Implementation details. Our EVA-Gaussian is trained on 1024×1024 pixel images across multiple training views using a single NVIDIA A800 GPU for 100K iterations with the AdamW (Loshchilov & Hutter, 2017) optimizer, unless otherwise specified. For the 3D Gaussian position estimation stage, it is first pretrained under the supervision of ground truth depth maps. Baselines are trained using their publicly available code. More implementation details are provided in Appendix A.

Datasets. We conduct experiments on two open-source human body datasets: THuman2.0 (Yu et al., 2021b) and THumanSit (Zhang et al., 2023). THuman2.0 contains 526 unique human models with their corresponding SMPL parameters, among which 100 individuals are randomly selected for our evaluation. The THumanSit dataset has a similar structure, containing 72 human models with around 60 poses for each, and we randomly choose 5 individuals with all poses for our evaluation.

Metrics. We report results on commonly used metrics: PSNR, SSIM (Wang et al., 2004), and LPIPS (Zhang et al., 2018), computed over the entire image, as some methods may produce artifacts outside the human bounding box (Lin et al., 2022; Zheng et al., 2024). We also include the inference time to demonstrate the real-time reconstruction capability of our method.

423 424

5.2 STEREO RECONSTRUCTION

Comparison with state-of-the-art feed-forward reconstruction methods. We first compare our approach against state-of-the-art (SOTA) feed-forward reconstruction methods, including ENeRF (Lin et al., 2022), pixelSplat (Charatan et al., 2024), MVSplat (Chen et al., 2024), MVSGaussian (Liu et al., 2024b), and GPS-Gaussian (Zheng et al., 2024). All experiments are conducted in a stereo-view setting, where the angle between the two camera views $\Delta = 45^{\circ}$. The attention modules in the scene reconstruction methods (Charatan et al., 2024; Liu et al., 2024b; Chen et al., 2024) are inefficient in their utilization of GPU memory, limiting their ability to train effectively at a high resolution of 1024×1024 . Therefore, we also conduct a fair comparison of all methods at a

Table 2: Comparison of feed-forward human reconstruction methods under different camera	angle
settings, at a resolution of 1024×1024 . Better results are marked in a deeper color. Notably,	GPS-
Gaussian fails to work effectively when $\Delta = 90^\circ$, as it is unable to meet its rectification require	ment.

THuman2.0		$\Delta=45^\circ$			$\Delta = 60^{\circ}$			$\Delta=72^\circ$			$\Delta=90^\circ$	
1024×1024	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
ENeRF	27.94	0.9573	0.0367	26.16	0.9452	0.0516	24.61	0.9309	0.0705	22.85	0.8990	0.1147
GPS-Gaussian	29.63	0.9703	0.0174	27.36	0.9630	0.0249	24.25	0.9519	0.0480	/	/	/
EVA-Gaussian	30.46	0.9730	0.0178	28.29	0.9654	0.0248	27.54	0.9614	0.0297	26.31	0.9555	0.0391
THumansit		$\Delta = 45^{\circ}$			$\Delta = 60^{\circ}$			$\Delta = 72^{\circ}$			$\Delta = 90^{\circ}$	
$\begin{array}{c} \text{THumansit} \\ 1024 \times 1024 \end{array}$	PSNR↑	$\Delta = 45^{\circ}$ SSIM \uparrow	LPIPS↓	PSNR↑	$\Delta = 60^{\circ}$ SSIM \uparrow	LPIPS↓	PSNR↑	$\Delta = 72^{\circ}$ SSIM \uparrow	LPIPS↓	PSNR↑	$\Delta = 90^{\circ}$ SSIM \uparrow	LPIPS↓
$\frac{\text{THumansit}}{1024 \times 1024}$ ENeRF	PSNR↑ 25.61	$\Delta = 45^{\circ}$ SSIM \uparrow 0.9397	LPIPS↓ 0.0494	PSNR↑ 23.80	$\Delta = 60^{\circ}$ SSIM \uparrow 0.9168	LPIPS↓ 0.0745	PSNR↑ 22.48	$\Delta = 72^{\circ}$ SSIM \uparrow 0.8956	LPIPS↓ 0.0985	PSNR↑ 21.20	$\begin{array}{c} \Delta = 90^{\circ} \\ \text{SSIM} \uparrow \\ 0.8571 \end{array}$	LPIPS↓ 0.1406
$\begin{tabular}{c} THumansit\\ 1024 \times 1024\\ \hline ENeRF\\ GPS-Gaussian \end{tabular}$	PSNR↑ 25.61 27.05	$\Delta = 45^{\circ}$ SSIM 0.9397 0.9584	LPIPS↓ 0.0494 0.0227	PSNR↑ 23.80 25.19	$\Delta = 60^{\circ}$ SSIM 0.9168 0.9480	LPIPS↓ 0.0745 0.0351	PSNR↑ 22.48 21.48	$\Delta = 72^{\circ}$ SSIM 0.8956 0.9276	LPIPS↓ 0.0985 0.0713	PSNR↑ 21.20 /	$\begin{array}{c} \Delta = 90^{\circ} \\ \text{SSIM} \uparrow \\ \hline 0.8571 \\ / \end{array}$	LPIPS↓ 0.1406 /



Figure 5: Qualitative comparison on THuman2.0 and THumansit. EVA-Gaussian achieves superior novel view rendering quality under diverse camera settings. Additional visualization results are provided in Appendix B.

resolution of 256×256. The quantitative results presented in Table 1 demonstrate that EVA-Gaussian achieves the best novel view quality in terms of PSNR, SSIM, and LPIPS, while maintaining the second-fastest inference speed.

Comparison under diverse angle changes between camera views. We further evaluate the per-formance of our method across four different angles between the two camera views, i.e., $\Delta =$ $45^{\circ}, 60^{\circ}, 72^{\circ}$, and 90° , at a high resolution of 1024×1024 . As shown in Table 2, our EVA-Gaussian outperforms all baseline methods on all metrics, achieving a maximum PSNR advantage of 5.12 dB. Notably, thanks to our EVA module, EVA-Gaussian remains effective even under extremely sparse camera settings, e.g., $\Delta = 90^{\circ}$. In contrast, GPS-Gaussian fails to work effectively due to its re-liance on stereo rectification. Fig. 5 presents the qualitative results of novel view rendering, where EVA-Gaussian outperforms previous SOTA methods in rendering quality, especially in scenarios with large viewpoint discrepancies.

5.3 MULTI-VIEW RECONSTRUCTION

We conduct experiments to compare our method against GPS-Gaussian under multi-view settings.
Table 3 presents the quantitative results. Our method demonstrates a significant advantage over the
baseline, with a more than 1.5 dB improvement. Notably, the rendering quality of GPS-Gaussian
drops significantly due to the mismatch among multiple inferences. In contrast, our method maintains high performance, thanks to the cross-view consistency ensured by our proposed EVA module.

5.4 ABLATION STUDY

We conduct a detailed ablation study on THuman2.0 in a stereo-view setting, where the angle between the two views $\Delta = 45^{\circ}$, as shown in Table 4 and Fig. 6. We gradually incorporate the EVA module, feature refinement module, and anchor loss to evaluate their individual contributions. The absence of the EVA module results in significant degradation across all metrics, and the network

Table 3: Comparison with GPS-Gaussian under different camera number settings. The results in bold represent the best performance. Our EVA-Gaussian achieves SOTA performance across various metrics, primarily due to the multi-view consistency enabled by our proposed EVA module.

1024 × 1024	THuman2.0 ($\Delta = 45^{\circ}$)				THumansit ($\Delta = 45^{\circ}$)							
1024×1024	3 views			4 views		3 views			4 views			
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
GPS-Gaussian	28.74	0.9655	0.0200	28.51	0.9636	0.0218	26.87	0.9523	0.0243	26.50	0.9498	0.0267
EVA-Gaussian	30.76	0.9722	0.0175	30.35	0.9707	0.0189	28.64	0.9596	0.0255	28.32	0.9582	0.0260

Table 4: Quantitative results of the ablation study on THuman2.0 in a stereo-view setting, where the angle between the two views $\Delta = 45^{\circ}$, at a resolution of 1024×1024 .

1024×1024	THuman2.0 ($\Delta = 45^{\circ}$)						
1024×1024	w/o EVA module	w/o feature refinement	w/o anchor loss	Full model			
PSNR↑	23.41	29.31	30.34	30.46			
SSIM↑	0.9380	0.9676	0.9724	0.9730			
LPIPS↓	0.0659	0.0191	0.0186	0.0178			



Figure 6: Qualitative visualization results of the ablation study on THuman2.0. Each module shows its effectiveness for a better visual output. The feature refinement (FR) module corrects geometric errors in the initial estimations, and the anchor loss further refines critical areas, such as the face, for generating novel view images with higher fidelity.

> struggles to perform multi-view 3D Gaussian geometry prediction. When feature refinement is excluded, visualizations reveal artifacts in critical areas, such as the hands and feet. Moreover, the lack of anchor loss leads to unreliable geometry predictions, particularly in the facial region, which in turn degrades the performance across all metrics, with a notable impact on LPIPS.

CONCLUSION

In this paper, we introduce EVA-Gaussian, a novel real-time 3D human reconstruction pipeline that employs multi-view attention-based 3D Gaussian position estimation and comprehensive fea-ture refinement. To ensure robust performance, the method is trained using both photometric loss and anchor loss. Quantitative and qualitative evaluations on benchmark datasets demonstrate that our EVA-Gaussian achieves state-of-the-art performance while maintaining a competitive inference speed, particularly under sparse camera settings.

While EVA-Gaussian synthesizes high-fidelity novel views, there remain several areas for improve-ment. For instance, the attention module can consume substantial GPU memory when processing a large number of input views or high-resolution images. In addition, the naive reprojection of pixels into 3D space may introduce conflicts in overlapping areas, leading to redundancy in the 3D rep-resentation. These limitations can be effectively addressed by incorporating RGBD information or developing overlap area detection techniques.

540	REFERENCES
541	

560

561

- T Berriel Martins and Javier Civera. Feature splatting for better novel view synthesis with low 542 overlap. arXiv e-prints, pp. arXiv-2405, 2024. 543
- 544 Zhongang Cai, Daxuan Ren, Ailing Zeng, Zhengyu Lin, Tao Yu, Wenjia Wang, Xiangyu Fan, Yang Gao, Yifan Yu, Liang Pan, Fangzhou Hong, Mingyuan Zhang, Chen Change Loy, Lei Yang, and 546 Ziwei Liu. HuMMan: Multi-modal 4d human dataset for versatile sensing and modeling. In 17th European Conference on Computer Vision, Tel Aviv, Israel, October 23–27, 2022, Proceedings, 547 Part VII, pp. 557–577. Springer, 2022. 548
- 549 David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaus-550 sian splats from image pairs for scalable generalizable 3d reconstruction. In *Proceedings of the* 551 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 19457–19467, 2024. 552
- Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. 553 Mysnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In Proceedings 554 of the IEEE/CVF international conference on computer vision, pp. 14124–14133, 2021. 555
- 556 Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. Mysplat: Efficient 3d gaussian splatting from sparse multi-view 558 images. arXiv preprint arXiv:2403.14627, 2024.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778, 2016. 562
- Yihui He, Rui Yan, Katerina Fragkiadaki, and Shoou-I Yu. Epipolar transformers. In Proceedings 563 of the ieee/cvf conference on computer vision and pattern recognition, pp. 7779–7788, 2020.
- 565 Liangxiao Hu, Hongwen Zhang, Yuxiang Zhang, Boyao Zhou, Boning Liu, Shengping Zhang, and 566 Liqiang Nie. Gaussianavatar: Towards realistic human avatar modeling from a single video via 567 animatable 3d gaussians. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 634-644, 2024. 568
- 569 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splat-570 ting for real-time radiance field rendering. ACM Trans. Graph., 42(4):139–1, 2023. 571
- Muhammed Kocabas, Jen-Hao Rick Chang, James Gabriel, Oncel Tuzel, and Anurag Ranjan. Hugs: 572 Human gaussian splats. In Proceedings of the IEEE/CVF conference on computer vision and 573 pattern recognition, pp. 505–515, 2024. 574
- 575 Youngjoong Kwon, Baole Fang, Yixing Lu, Haoye Dong, Cheng Zhang, Francisco Vicente Car-576 rasco, Albert Mosella-Montoro, Jianjin Xu, Shingo Takagi, Daeil Kim, et al. Generalizable human 577 gaussians for sparse view synthesis. arXiv preprint arXiv:2407.12777, 2024.
- 578 Jiahui Lei, Yufu Wang, Georgios Pavlakos, Lingjie Liu, and Kostas Daniilidis. Gart: Gaussian 579 articulated template models. In Proceedings of the IEEE/CVF Conference on Computer Vision 580 and Pattern Recognition, pp. 19876–19887, 2024.
- Haotong Lin, Sida Peng, Zhen Xu, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Efficient 582 neural radiance fields for interactive free-viewpoint video. In SIGGRAPH Asia 2022 Conference 583 Papers, pp. 1–9, 2022. 584
- 585 Tianqi Liu, Guangcong Wang, Shoukang Hu, Liao Shen, Xinyi Ye, Yuhang Zang, Zhiguo Cao, Wei 586 Li, and Ziwei Liu. Fast generalizable gaussian splatting reconstruction from multi-view stereo. arXiv preprint arXiv:2405.12218, 2024a.
- 588 Tianqi Liu, Guangcong Wang, Shoukang Hu, Liao Shen, Xinyi Ye, Yuhang Zang, Zhiguo Cao, 589 Wei Li, and Ziwei Liu. Mysgaussian: Fast generalizable gaussian splatting reconstruction from 590 multi-view stereo. arXiv preprint arXiv:2405.12218, 2024b.
- Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: 592 A skinned multi-person linear model. In Seminal Graphics Papers: Pushing the Boundaries, Volume 2, pp. 851-866. 2023.

627

631

632

- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
- Camillo Lugaresi, Jiuqiang Tang, Hadon Nash, Chris McClanahan, Esha Uboweja, Michael Hays,
 Fan Zhang, Chuo-Ling Chang, Ming Guang Yong, Juhyun Lee, et al. Mediapipe: A framework
 for building perception pipelines. *arXiv preprint arXiv:1906.08172*, 2019.
- 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. *Communications* of the ACM, 65(1):99–106, 2021.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–15, 2022.
- Panwang Pan, Zhuo Su, Chenguo Lin, Zhen Fan, Yongjie Zhang, Zeming Li, Tingting Shen, Yadong
 Mu, and Yebin Liu. Humansplat: Generalizable single-image human gaussian splatting with
 structure priors. *arXiv preprint arXiv:2406.12459*, 2024.
- Antonis Papapantoleon, Dylan Possamaï, and Alexandros Saplaouras. Stability of backward stochastic differential equations: the general lipschitz case. *Electronic Journal of Probability*, 28:1–56, 2023.
- Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios
 Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single
 image. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
 pp. 10975–10985, 2019.
- ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶²⁰
 ⁶²⁰
 ⁶²¹
 ⁶²¹
 ⁶²¹
 ⁶²¹
 ⁶²¹
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁰
 ⁶¹⁰
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁰
 ⁶¹¹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹¹
 ⁶¹¹
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed ical image segmentation. In *Medical image computing and computer-assisted intervention– MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceed- ings, part III 18*, pp. 234–241. Springer, 2015.
- ⁶²⁶ Walter Rudin et al. *Principles of mathematical analysis*, volume 3. McGraw-hill New York, 1964.
- Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Efficient attention:
 Attention with linear complexities. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 3531–3539, 2021.
 - Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8922–8931, 2021.
- Stanislaw Szymanowicz, Chrisitian Rupprecht, and Andrea Vedaldi. Splatter image: Ultra-fast
 single-view 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10208–10217, 2024.
- Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm:
 Large multi-view gaussian model for high-resolution 3d content creation. *arXiv preprint arXiv:2402.05054*, 2024.
- Hanzhang Tu, Ruizhi Shao, Xue Dong, Shunyuan Zheng, Hao Zhang, Lili Chen, Meili Wang,
 Wenyu Li, Siyan Ma, Shengping Zhang, et al. Tele-aloha: A low-budget and high-authenticity
 telepresence system using sparse rgb cameras. *arXiv preprint arXiv:2405.14866*, 2024.
- Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T
 Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi view image-based rendering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4690–4699, 2021.

669

696 697

699 700

- ⁶⁴⁸
 ⁶⁴⁹
 ⁶⁴⁹ Thou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- Jing Wen, Xiaoming Zhao, Zhongzheng Ren, Alexander G Schwing, and Shenlong Wang. Go mavatar: Efficient animatable human modeling from monocular video using gaussians-on-mesh. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2059–2069, 2024.
- Christopher Wewer, Kevin Raj, Eddy Ilg, Bernt Schiele, and Jan Eric Lenssen. latentsplat:
 Autoencoding variational gaussians for fast generalizable 3d reconstruction. *arXiv preprint arXiv:2403.16292*, 2024.
- Haofei Xu, Anpei Chen, Yuedong Chen, Christos Sakaridis, Yulun Zhang, Marc Pollefeys, Andreas Geiger, and Fisher Yu. Murf: Multi-baseline radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20041–20050, 2024.
- Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from
 one or few images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4578–4587, 2021a.
- Tao Yu, Zerong Zheng, Kaiwen Guo, Pengpeng Liu, Qionghai Dai, and Yebin Liu. Function4d:
 Real-time human volumetric capture from very sparse consumer rgbd sensors. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR2021)*, June 2021b.
- Jiajun Zhang, Yuxiang Zhang, Hongwen Zhang, Boyao Zhou, Ruizhi Shao, Zonghai Hu, and
 Yebin Liu. Ins-hoi: Instance aware human-object interactions recovery. *arXiv preprint arXiv:2312.09641*, 2023.
- Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang Xu.
 Gs-Irm: Large reconstruction model for 3d gaussian splatting. *arXiv preprint arXiv:2404.19702*, 2024.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.
- Shunyuan Zheng, Boyao Zhou, Ruizhi Shao, Boning Liu, Shengping Zhang, Liqiang Nie, and Yebin
 Liu. Gps-gaussian: Generalizable pixel-wise 3d gaussian splatting for real-time human novel
 view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19680–19690, 2024.



Figure 7: Qualitative comparison on THuman2.0 and THumansit. EVA-Gaussian demonstrates superior novel view rendering quality under diverse camera settings.

A MORE IMPLEMENTATION DETAILS

732 Network architectures. Our Gaussian position estimation network $\mathcal{D}_{\theta_1}^P$ utilizes a U-Net as the 733 backbone. The architecture incorporates four stages of $2 \times$ down-sampling using average pooling to 734 extract essential feature details. Symmetrically, the network features four stages of $2 \times$ up-sampling, 735 achieved through transpose convolutional neural networks. The EVA module is incorporated before 736 the $4\times$, $8\times$, $16\times$ down-sampling and up-sampling blocks. The channel dimension starts at 64 prior 737 to the first down-sampling block, doubling after each down-sampling block and halving after each 738 up-sampling block, which is facilitated by two residual blocks (He et al., 2016). The Gaussian attribute estimation network $\mathcal{D}_{\theta_2}^A$ also employs a U-Net backbone, but it does not include the EVA 739 modules and performs only two stages of $2 \times$ down-sampling. The architecture of the feature refiner, 740 $\mathcal{D}^{R}_{\theta_{3}}$, mirrors that of $\mathcal{D}^{A}_{\theta_{2}}$, but operates in a recurrent manner. 741

742 More training details. The training hyper-parameters are set as follows: $\lambda_1 = 1, \lambda_2 = 1, \lambda_3 =$ 743 10^3 , $\lambda_{\text{opacity}} = 1$, $\lambda_{\text{opacity}} = 1$, $\lambda_{\text{render}} = 0.25$, and t = 0.05. The number of recurrent loops L 744 mentioned in Sec. 4.4 is empirically set to L = 1 to enhance temporal efficiency. Each training batch contains 2 to 4 source view images, depending on the specific reconstruction task. For instance, 745 the stereo reconstruction task in Sec. 5.2 utilizes 2 source view images. For novel view image 746 supervision, 3 randomly selected views are chosen between each adjacent pair of source views to 747 compute \mathcal{L}_{refine} and \mathcal{L}_{render} . The learning rate for deep supervised pre-training and overall network 748 training is initialized to 0.0002 and decreases linearly with the number of training epochs. 749

750 751

752

728

729 730

731

B MORE VISUALIZATION RESULTS

In this section, we present additional visualization results in Fig. 7 to compare our method with
 SOTA approaches GPS-Gaussian and ENeRF on the Thuman2.0 (Yu et al., 2021b) and Thumansit
 (Zhang et al., 2023) datasets at a resolution of 1024×1024. The results demonstrate that EVA-Gaussian achieves the highest novel view fidelity across various camera viewpoint settings. In con-

trast, GPS-Gaussian struggles to handle the artifacts produced by errors in geometric predictions, while ENeRF generates much more blurry and low-fidelity results compared to both GPS-Gaussian and EVA-Gaussian. Notably, under settings of large viewpoint discrepancy, e.g., $\Delta = 90^{\circ}$, EVA-Gaussian maintains robust performance, while GPS-Gaussian fails to function effectively in these scenarios.

C PROOF OF DEPTH EQUALITY

In this section, we prove that for each pixel on the 3D Gaussian maps $\{M_i\}_{i=1}^n$, the rendered depth equals to the predicted 3D Gaussian depth.

We begin by defining the collection for opacity parameters as $\boldsymbol{o} := [o_1, \dots, o_i, \dots, o_N] \in \mathbb{R}^N$ of all considered 3D Gaussians and the collection of all 3D Gaussian scaling factors as:

$$\tilde{\boldsymbol{S}} = \begin{bmatrix} \boldsymbol{s}^1, \boldsymbol{s}^2, \boldsymbol{s}^3 \end{bmatrix}^T = \begin{pmatrix} s_1^1 & s_2^1 & \cdots & s_N^1 \\ s_1^2 & s_2^2 & \cdots & s_N^2 \\ s_1^3 & s_2^3 & \cdots & s_N^3 \end{pmatrix},$$
(15)

where $s^1 := [s_1^1, s_2^1, \dots, s_N^1] \in \mathbb{R}^N$, $s^2 := [s_1^2, s_2^2, \dots, s_N^2] \in \mathbb{R}^N$ and $s^3 := [s_1^3, s_2^3, \dots, s_N^3] \in \mathbb{R}^N$. For the 3D Gaussian with the *i*-th greatest depth, the associated scaling matrix is constructed from the corresponding scaling factors as:

$$\mathbf{S}_{i} = \begin{pmatrix} s_{i}^{1} & 0 & 0\\ 0 & s_{i}^{2} & 0\\ 0 & 0 & s_{i}^{3} \end{pmatrix},\tag{16}$$

and its z-value is denoted by z_i . The rendered depth map is expressed as:

$$D(\boldsymbol{x}) = \sum_{i=1}^{N} z_i o_i G'_i(\boldsymbol{x}) \prod_{j=1}^{i-1} (1 - o_j G'_j(\boldsymbol{x}))),$$
(17)

where $x \in \mathbb{R}^2$ is a variable on the coordinate system of the image plane and $G'_i(x)$ is the 2D Gaussian that corresponds to the 3D Gaussian with the *i*-th greatest depth after splatting.

⁷⁸⁸ In the camera's coordinate system, we define a 3D Gaussian as on the reprojected ray of a pixel x', ⁷⁸⁹ in condition that the center of this 3D Gaussian lies along the ray originating from the camera center ⁷⁹⁰ and pointing toward the point [x', 1]. We use Z(x') to denote the z-value of the first 3D Gaussian ⁷⁹¹ that appears on this reprojected ray.

Based on the above definitions, we have the following theorem:

Theorem C.1. When the opacity o approaches 1 and each value in \tilde{S} is sufficiently small, it holds for each pixel x' on the image plane that:

$$\lim_{\substack{\boldsymbol{o} \to \mathbf{1} \\ \hat{\boldsymbol{S}} \to \mathbf{0}^+}} D(\boldsymbol{x}') = Z(\boldsymbol{x}').$$
(18)

Theorem C.1 implies that the z-value of the 3D Gaussian at pixel x is equal to the corresponding value on the depth map when the scale of Gaussian is sufficiently small and the opacity approaches 1. To prove Theorem C.1, we introduce the following lemma from the well-known Moore-Osgood Theorem in (Papapantoleon et al., 2023):

Lemma C.1. (*Moore-Osgood Theorem*) Let (Γ, d_{Γ}) be a metric space and $(\gamma_{k,p})_{k,p\in\mathbb{N}}$ be a sequence such that $\gamma_{\infty,p} := \lim_{k\to\infty} \gamma_{k,p}$ exists for every $p \in \mathbb{N}$ and $\gamma_{k,\infty} := \lim_{p\to\infty} \gamma_{k,p}$ exists for every $k \in \mathbb{N}$. If (i) $\lim_{p\to\infty} \sup_{k\in\mathbb{N}} d_{\Gamma}(\gamma_{k,p},\gamma_{k,\infty}) = 0$ and (ii) $\lim_{k\to\infty} d_{\Gamma}(\gamma_{k,p},\gamma_{\infty,p}) = 0, \forall p \in \mathbb{N}$, then the joint limit $\lim_{k,p\to\infty} \gamma_{n,k}$ exists. In particular, it holds that $\lim_{k,p\to\infty} \gamma_{k,p} = \lim_{p\to\infty} \gamma_{\infty,p} = \lim_{k\to\infty} \gamma_{k,\infty}$.

Lemma C.1 can be regarded as a special case of Theorem 7.11 from (Rudin et al., 1964). This lemma states that for a doubly-indexed sequence, if the sequence converges uniformly with respect to one

761 762 763

764

765

766

777 778 779

780

794

795 796 797

820

829

857

858

index while converging pointwise with respect to the other index, then the limit of the sequence
exists. Moreover, this limit is equivalent to the individual limits obtained by separately considering
each index, regardless of the order in which the limiting processes are performed. This result can be
extended to continuous multi-variable functions. Specifically, if a continuous function demonstrates
uniform convergence with respect to one variable and pointwise convergence with respect to another
variable, then the joint limit of the function with respect to both variables can be decomposed into
the separate limits with respect to each variable considered independently.

Based on this theoretical foundation, we are now ready to proceed with the proof of Theorem C.1.

Proof. When the opacity value $o_i \in \mathbb{R}$ approaches 1 and the scale factor $s_i = [s_i^1, s_i^2, s_i^3] \in \mathbb{R}^3$ is sufficiently small for each Gaussian, the depth value is given by:

$$\lim_{\substack{\boldsymbol{o} \to \mathbf{1} \\ \hat{\boldsymbol{S}} \to \mathbf{0}^+}} D(\boldsymbol{x}') = \lim_{\substack{\boldsymbol{o} \to \mathbf{1} \\ \hat{\boldsymbol{S}} \to \mathbf{0}^+}} \sum_{i=1}^N (z_i o_i G'_i(\boldsymbol{x}') \prod_{j=1}^{i-1} (1 - o_j G'_j(\boldsymbol{x}')))$$

$$= \lim_{\boldsymbol{o} \to \mathbf{1}} \sum_{i=1}^N (z_i o_i e^{-\frac{1}{2} (\boldsymbol{x}' - \boldsymbol{\mu}_i)^T (\boldsymbol{J} \boldsymbol{W} \boldsymbol{R}_i \boldsymbol{S}_i \boldsymbol{S}_i^T \boldsymbol{R}_i^T \boldsymbol{W}^T \boldsymbol{J}^T)^{-1} (\boldsymbol{x}' - \boldsymbol{\mu}_i)}$$
(19)

(20)

)

(24)

 $\prod_{j=1}^{i} (1 - o_j e^{-\frac{1}{2} (\boldsymbol{x}' - \boldsymbol{\mu}_j)^T (\boldsymbol{J} \boldsymbol{W} \boldsymbol{R}_j \boldsymbol{S}_j \boldsymbol{S}_j^T \boldsymbol{R}_j^T \boldsymbol{W}^T \boldsymbol{J}^T)^{-1} (\boldsymbol{x}' - \boldsymbol{\mu}_j)}))$

$$\begin{array}{l} \textbf{830} \\ \textbf{831} \\ \textbf{832} \\ \textbf{832} \\ \textbf{833} \\ \textbf{834} \\ \textbf{835} \\ \textbf{835} \\ \textbf{835} \\ \textbf{836} \\ \textbf{836} \\ \textbf{836} \\ \textbf{836} \\ \textbf{836} \\ \textbf{837} \\ \textbf{83$$

where (a) is from Lemma C.1. Specifically, the function $\sum_{i=1}^{N} (z_i o_i G'_i(\boldsymbol{x}) \prod_{j=1}^{i-1} (1 - o_j G'_j(\boldsymbol{x})))$ is continuous with respect to the two variables \boldsymbol{o} and \boldsymbol{s} . Besides, it converges uniformly as $\hat{\boldsymbol{S}} \to \boldsymbol{0}^+$ and as $\boldsymbol{o} \to \boldsymbol{1}$. This implies that the joint limit of \boldsymbol{o} and \boldsymbol{s} can be decomposed into the separate limits of \boldsymbol{o} and \boldsymbol{s} . Thus, we have:

$$\lim_{\hat{S}\to0^{+}} \lim_{o\to1} \sum_{i=1}^{N} (z_{i}o_{i}e^{-\frac{1}{2}(x'-\mu_{i})^{T}(JWR_{i}S_{i}S_{i}^{T}R_{i}^{T}W^{T}J^{T})^{-1}(x'-\mu_{i})} \\ \prod_{j=1}^{i-1} (1-o_{j}e^{-\frac{1}{2}(x'-\mu_{j})^{T}(JWR_{j}S_{j}S_{j}^{T}R_{j}^{T}W^{T}J^{T})^{-1}(x'-\mu_{j})})) \\ = \lim_{\hat{S}\to0^{+}} \sum_{i=1}^{N} (\lim_{o_{i}\to1} z_{i}o_{i}e^{-\frac{1}{2}(x'-\mu_{i})^{T}(JWR_{i}S_{i}S_{i}^{T}R_{i}^{T}W^{T}J^{T})^{-1}(x'-\mu_{i})} \\ \lim_{(o_{j},\cdots,o_{i-1})\to1} \prod_{j=1}^{i-1} (1-o_{j}e^{-\frac{1}{2}(x'-\mu_{j})^{T}(JWR_{j}S_{j}S_{j}^{T}R_{j}^{T}W^{T}J^{T})^{-1}(x'-\mu_{j})})) \qquad (22) \\ = \lim_{\hat{S}\to0^{+}} \sum_{i=1}^{N} (z_{i}e^{-\frac{1}{2}(x'-\mu_{i})^{T}(JWR_{i}S_{i}S_{i}^{T}R_{i}^{T}W^{T}J^{T})^{-1}(x'-\mu_{i})} \\ \prod_{j=1}^{i-1} (1-e^{-\frac{1}{2}(x'-\mu_{j})^{T}(JWR_{j}S_{j}S_{j}^{T}R_{j}^{T}W^{T}J^{T})^{-1}(x'-\mu_{j})})). \qquad (23)$$

The 3D Gaussians typically assume an ellipsoidal geometric shape. However, when the scaling factors are sufficiently small, the ellipsoid can be approximated as a sphere, such that $s^1 = s^2 = s^3$. As a result, the scaling matrix for the 3D Gaussian with the *i*-th greatest depth becomes:

862
863

$$\boldsymbol{S}_{i}^{'} := \begin{pmatrix} s_{i}^{1} & 0 & 0\\ 0 & s_{i}^{1} & 0\\ 0 & 0 & s_{i}^{1} \end{pmatrix}.$$

864 Consequently, we have: 865 866 $\lim_{\hat{\boldsymbol{S}} \to \boldsymbol{0}^+} \sum_{i=1}^{N} (z_i e^{-\frac{1}{2} (\boldsymbol{x}' - \boldsymbol{\mu}_i)^T (\boldsymbol{J} \boldsymbol{W} \boldsymbol{R}_i \boldsymbol{S}_i \boldsymbol{S}_i^T \boldsymbol{R}_i^T \boldsymbol{W}^T \boldsymbol{J}^T)^{-1} (\boldsymbol{x}' - \boldsymbol{\mu}_i)}$ 867 868 $\prod_{j=1}^{i-1} \left(1 - e^{-\frac{1}{2}(\boldsymbol{x}' - \boldsymbol{\mu}_j)^T (\boldsymbol{J} \boldsymbol{W} \boldsymbol{R}_j \boldsymbol{S}_j \boldsymbol{S}_j^T \boldsymbol{R}_j^T \boldsymbol{W}^T \boldsymbol{J}^T)^{-1} (\boldsymbol{x}' - \boldsymbol{\mu}_j)}\right)\right)$ 870 871 $= \lim_{\substack{i \to 0^+ \\ 1 \to -2^-, 3}} \sum_{i=1}^{N} (z_i e^{-\frac{1}{2} (\mathbf{x}' - \boldsymbol{\mu}_i)^T (\mathbf{J} \mathbf{W} \mathbf{R}_i \mathbf{S}_i \mathbf{S}_i^T \mathbf{R}_i^T \mathbf{W}^T \mathbf{J}^T)^{-1} (\mathbf{x}' - \boldsymbol{\mu}_i)}$ 872 873 874 875 $\prod_{j=1}^{i-1} \left(1 - e^{-\frac{1}{2}(\boldsymbol{x}' - \boldsymbol{\mu}_j)^T (\boldsymbol{J} \boldsymbol{W} \boldsymbol{R}_j \boldsymbol{S}_j \boldsymbol{S}_j^T \boldsymbol{R}_j^T \boldsymbol{W}^T \boldsymbol{J}^T)^{-1} (\boldsymbol{x}' - \boldsymbol{\mu}_j)}\right)$ (25)877 $= \lim_{\boldsymbol{s}^1 \to \boldsymbol{0}^+} \sum_{i=1}^N (z_i e^{-\frac{1}{2} (\boldsymbol{x}' - \boldsymbol{\mu}_i)^T (\boldsymbol{J} \boldsymbol{W} \boldsymbol{R}_i \boldsymbol{S}_i' \boldsymbol{S}_i'^T \boldsymbol{R}_i^T \boldsymbol{W}^T \boldsymbol{J}^T)^{-1} (\boldsymbol{x}' - \boldsymbol{\mu}_i)}$ 878 879 880 $\prod_{i=1}^{i-1} (1 - e^{-\frac{1}{2} (\boldsymbol{x}' - \boldsymbol{\mu}_j)^T (\boldsymbol{J} \boldsymbol{W} \boldsymbol{R}_j \boldsymbol{S}_j' \boldsymbol{S}_j'^T \boldsymbol{R}_j^T \boldsymbol{W}^T \boldsymbol{J}^T)^{-1} (\boldsymbol{x}' - \boldsymbol{\mu}_j)})).$ (26)882 883 884 From (26), we see that when $x' = \mu_i$, it gives that 885 $\rho^{-\frac{1}{2}(\boldsymbol{x}'-\boldsymbol{\mu}_i)^T}(\boldsymbol{J}\boldsymbol{W}\boldsymbol{R}_i\boldsymbol{S}_i'\boldsymbol{S}_i'^T\boldsymbol{R}_i^T\boldsymbol{W}^T\boldsymbol{J}^T)^{-1}(\boldsymbol{x}'-\boldsymbol{\mu}_i)=1.$ 886 (27)887 888 Otherwise, if $x' \neq \mu_i$, we have 889 1 / / 890

$$\lim_{\substack{s_i^1 \to 0^+ \\ s_i^1 \to 0^+}} e^{-\frac{1}{2(s_i^1)^2} (\mathbf{x}' - \boldsymbol{\mu}_i)^T (\mathbf{J} \mathbf{W} \mathbf{R}_i \mathbf{R}_i^T \mathbf{W}^T \mathbf{J}^T)^{-1} (\mathbf{x}' - \boldsymbol{\mu}_i)}$$

$$= \lim_{\substack{s_i^1 \to 0^+ \\ s_i^1 \to 0^+}} e^{-\frac{1}{2(s_i^1)^2} (\mathbf{x}' - \boldsymbol{\mu}_i)^T (\mathbf{J} \mathbf{W} \mathbf{R}_i \mathbf{R}_i^T \mathbf{W}^T \mathbf{J}^T)^{-1} (\mathbf{x}' - \boldsymbol{\mu}_i)}$$

$$= 0.$$
(28)

By combining Eq. 19 – Eq. 28, we have

896 897

905 906

907 908

$$\lim_{\substack{\boldsymbol{o} \to \mathbf{1} \\ \hat{\boldsymbol{S}} \to \mathbf{0}^+}} D(\boldsymbol{x}') = \sum_{i=1}^N (\lim_{s_i^1 \to 0^+} z_i e^{-\frac{1}{2} (\boldsymbol{x}' - \boldsymbol{\mu}_i)^T (\boldsymbol{J} \boldsymbol{W} \boldsymbol{R}_i \boldsymbol{S}_i' \boldsymbol{S}_i'^T \boldsymbol{R}_i^T \boldsymbol{W}^T \boldsymbol{J}^T)^{-1} (\boldsymbol{x}' - \boldsymbol{\mu}_i)} \\ \prod_{i=1}^{i-1} \lim_{s_i^1 \to 0^+} (1 - e^{-\frac{1}{2} (\boldsymbol{x}' - \boldsymbol{\mu}_j)^T (\boldsymbol{J} \boldsymbol{W} \boldsymbol{R}_j \boldsymbol{S}_j' \boldsymbol{S}_j'^T \boldsymbol{R}_j^T \boldsymbol{W}^T \boldsymbol{J}^T)^{-1} (\boldsymbol{x}' - \boldsymbol{\mu}_j)}))$$
(29)

$$=Z(\boldsymbol{x}'), \tag{30}$$

which completes the proof.

D DETAILS ON EFFICIENT CROSS-VIEW ATTENTION

In this section, we further clarify the motivation of our proposed EVA module, detail its architecture,and analyze its advantages over existing attention mechanisms.

Attention mechanisms that leverage multi-view correspondence across different viewpoints, such as epipolar attention (He et al., 2020), have proven beneficial for downstream 3D tasks like pose estimation. Epipolar attention has demonstrated its effectiveness by performing attention for each pixel only with sampled points along its epipolar line. This approach is based on the principle that a pixel in the source image corresponds to a pixel along the epipolar line in the target image. However, the sampling process and attention calculation in traditional epipolar attention are computationally and temporally intensive, as shown in 5. To mitigate this problem, we propose the EVA module, specifically tailored for the camera settings in feed forward human 3D Gaussian reconstruction, Table 5: Comparison of the GPU memory usage of different feed forward 3D Gaussian reconstruction methods. Compared to previous feed-forward methods, our EVA-Gaussian method maintains highly competitive GPU memory usage.

921	Batch Size=1	GPU Memory Usage						
922	Input Image Resolution	$2\times 3\times 128\times 128$	$2\times 3\times 256\times 256$	$2 \times 3 \times 512 \times 512$	$2\times 3\times 1024\times 1024$			
923	PixelSplat (Epipolar Attention)	6099 MiB	13429 MiB	49598 MiB	Out of Memory			
004	MVSplat	3040 MiB	6584 MiB	27082 MiB	Out of Memory			
924	GPS-Gaussian	1909 MiB	2357 MiB	4035 MiB	11215 MiB			
925	EVA-Gaussian	2289 MiB	3171 MiB	7185 MiB	24121 MiB			



Figure 8: The correspondences between two source view images. The red points in the left view and the green points in the right view are the matching correspondent points. After transferring points from the left view to the right view at the exact position, and connecting them with a line, it is intuitive that the connecting lines are nearly parallel to the x-axis.

where cameras are closely positioned and oriented towards the same point on the human body. In this setting, the corresponding connections between matched pairs align parallel to the x-axis, as depicted in Fig. 8.

In contrast to traditional attention mechanisms, our EVA only computes attention weights for each pixel with nearby pixels along the x-axis. Moreover, we implement this attention mechanism within the deeper layers of the UNet architecture, as shown in Fig. 9, to fuse the dense spatial information from nearby pixels for each pixel. In addition, in order to enable a lager receptive field for each pixel and mitigate the potential loss of multi-view correspondences at the boundaries of local windows at a minimal cost, we perform the attention mechanism twice, with the second iteration using a window shifted by half of its size. Fig. 10 illustrates the key differences between EVA and other attention mechanisms including epipolar attention (He et al., 2020), and the cross attention in LoFTR (Sun et al., 2021). A quantitative comparison between our EVA module and other attention modules in terms of both temporal and computational costs is summarized in Table 6, which demonstrates the efficiency gains achieved through our approach. Our EVA module consumes less than 10% of the time overhead required by MVSplat's attention mechanism.

E VISUALIZATION FOR DIFFERENT NOVEL VIEW CAMERA SETUPS

In this section, we present additional visualization results under different novel view camera settings.

965Novel View Synthesis under Diverse Viewpoints. The 3D Gaussians generated under a uniformly966placed camera setup with $\Delta = 60^{\circ}$ illustrate a strong generalization ability to the random novel967view camera setup with pitch and yaw ranging from -25° to $+25^{\circ}$. To be more specific, as shown968in Fig. 11, given a pair of input images at $\Delta = 60^{\circ}$, EVA-Gaussian can effectively infer a human 3D969Gaussian model and render it under (A) yaw=0^{\circ}, pitch=0^{\circ}, (B) yaw $\in [15^{\circ}, 25^{\circ}]$, pitch $\in [15^{\circ}, 25^{\circ}]$,970(C) yaw=0^{\circ}, pitch $\in [15^{\circ}, 25^{\circ}]$, (D) yaw=0^{\circ}, pitch $\in [-25^{\circ}, -15^{\circ}]$ with promising visual quality.

Novel View Synthesis Under Higher Resolutions. We also infer 3D Gaussians using models trained with 1K resolution and render them at both 1K and 2K resolutions. Fig. 12 shows that the 3D

Table 6: Comparison of the temporal and computational efficiency among different attention modules. Notably, the GPU memory consumption of our EVA module does not scale up with window sizes, as the efficient attention algorithm (Shen et al., 2021) is adopted for implementation.



Figure 9: Detailed model structure for the Gaussian Position estimation network $\mathcal{D}_{\theta_1}^P$ in Sec. 4.2. EVA is implemented in the middle layers of the UNet.

Gaussians generated by GPS-Gaussian and EVA-Gaussian can be rendered at 2K resolution. However, artifacts such as holes and incompleteness exist in the 2K renderings of GPS-Gaussian, likely due to the lack of supervision from 2K novel view images. In contrast, although EVA-Gaussian is also not supervised by 2K novel view images, it exhibits greater robustness across different rendering resolutions, which mainly benefits from the stable performance of our feature refinement module.

1015 1016

1007 1008

F CROSS-DOMAIN EVALUATION

1017 1018

In this section, we evaluate the cross-domain capabilities of EVA-Gaussian. We first perform evaluation on THumansit dataset with a model trained on THuman2.0 dataset, and subsequently perform evaluation on THuman2.0 dataset with models trained on THumansit dataset. The camera view angles are consistently set to $\Delta = 60^{\circ}$.

Table 7 presents a comparison on the cross-domain generalization abilities between EVA-Gaussian and GPS-Gaussian. Given that THumansit contains significantly more human models (over 4,000) compared to THuman2.0 (526 models), the performance of both EVA-Gaussian and GPS-Gaussian trained on THumansit perform robustly when evaluated on the THuman2.0 dataset.

1050 1051 1052

1054

1056

1058

1062 1063



Figure 10: Comparison between different attention mechanisms. (A) is a cross attention mechanism adopted by the area of feature matching, e.g. LoFTR Sun et al. (2021) (B) is epipolar attention from epipolar transformer He et al. (2020), (C) and (D) are the proposed Efficient cross-View Attention at different window embedding stage. EVA only does attention with the most relevant pixels, thus greatly reduces the computational and temporal overhead.



Figure 11: Qualitative results under diverse novel view camera settings demonstrate that the 3D Gaussian model inferred by EVA-Gaussian consistently achieves high-quality novel view renderings across various random camera configurations.

1067 Furthermore, when trained on THumansit dataset and evaluated on THuman2.0 dataset, EVA-1068 Gaussian shows a greater performance improvement (+2.11 dB in PSNR) compared to GPS-1069 Gaussian (+1.97 dB in PSNR). This is attributed to the strong data processing ability of attention 1070 modules in EVA-Gaussian, allowing EVA-Gaussian to maintain consistent robustness when pro-1071 vided with sufficient data. This is further illustrated by the evaluation on THumansit, where models trained on THuman2.0 experience a performance decline due to limited data availability; however, EVA-Gaussian still outperforms GPS-Gaussian, achieving a performance gain of 0.41 dB in PSNR under these conditions. In addition, to demonstrate EVA-Gaussian's strong generalization ability 1074 across datasets, we present visualization results in Fig. 13, which shows that EVA-Gaussian do not 1075 suffer greatly from the out-of-domain problem, since we have explicitly introduced inductive bias to our EVA module. 1077

In conclusion, EVA-Gaussian demonstrates effective generalization ability to cross-domain datasets
 and exhibits exceptional robustness when a sufficient amount of data is available. Moreover, EVA-Gaussian consistently achieves superior performance on the out-of-domain data.



Figure 12: Qualitative rendered results across different resolutions, compared with GPS-Gaussian. Both models are trained with 1K resolution images. There exist incomplete artifacts in the 2K rendering of GPS-Gaussian. In contrast, EVA-Gaussian does not exhibit this issue and produces high-quality rendering results at 2K resolution.







Figure 13: Visualization of cross-domain evaluation results for EVA-Gaussian. The left side displays the rendered results generated by EVA-Gaussian trained on the THuman2.0 dataset and evaluated on the THumansit dataset, while the right side shows the rendered results from EVA-Gaussian trained on the THumansit dataset and evaluated on the THuman2.0 dataset.

Table 7: Quantitative results of cross-domain validations, compared with GPS-Gaussian. It demon-strates that our EVA-Gaussian consistently outperforms GPS-Gaussian when evaluated on various training and evaluation datasets.

Mathod	THumar	$nsit \rightarrow TH$	uman2.0	THuman2.0 \rightarrow THumansit			
Methou	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS	
GPS-Gaussian	29.33	0.9733	0.0325	20.86	0.9243	0.0872	
EVA-Gaussian	30.40	0.9751	0.0321	21.27	0.9275	0.0876	



Figure 14: Visualization of EVA-Gaussian on real-world data. Minor artifacts on the human boundary mainly arise from the noisy human mask. Notably, GPS-Gaussian cannot generate reasonable
outcome under this camera setting.

G REAL-WORLD DATA EVALUATION

In this section, we evaluate our model on on real-world data, the HuMMan Cai et al. (2022) dataset,
 which is a real-world dataset captured with RGB cameras at 1K resolution.

We select images from the front two cameras (ID: 1 and ID: 9) as inputs, infer the 3D Gaussians through EVA-Gaussian model and render novel view on the viewpoint of ID:0. The visualization results, as illustrated in Fig. 14, demonstrate that EVA-Gaussian produces high-quality novel view images in real-world settings.

1164 It is important to note that, due to the sparse input view angles of only 90°, GPS-Gaussian is unable to generate reasonable outcomes.