EVA-Gaussian: 3D Gaussian-based Real-time Human Novel View Synthesis under Diverse Multi-view Camera Settings

A. More Implementation Details

Network architectures. Our Gaussian position estimation network $\mathcal{D}_{\theta_1}^P$ utilizes a U-Net as the backbone. The architecture incorporates four stages of 2× down-sampling using average pooling to extract essential feature details. Symmetrically, the network features four stages of $2 \times$ upsampling, achieved through transpose convolutional neural networks. The EVA module is incorporated before the $4\times$, $8\times$, $16\times$ down-sampling and up-sampling blocks. The channel dimension starts at 64 prior to the first downsampling block, doubling after each down-sampling block and halving after each up-sampling block, which is facilitated by two residual blocks [2]. The Gaussian attribute estimation network $\mathcal{D}_{ heta_2}^A$ also employs a U-Net backbone, but it does not include the EVA modules and performs only two stages of 2× down-sampling. The architecture of the feature refiner, $\mathcal{D}_{\theta_3}^R$, mirrors that of $\mathcal{D}_{\theta_2}^A$, but operates in a recurrent manner.

More training details. The training hyper-parameters are set as follows: $\lambda_1 = 1, \lambda_2 = 1, \lambda_3 = 10^3, \lambda_{\text{opacity}} = 1, \lambda_{\text{opacity}} = 1, \lambda_{\text{render}} = 0.25,$ and t = 0.05. The number of recurrent loops L for the feature refinement module 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, the stereo reconstruction task in the experiment part utilizes 2 source view images. For novel view image supervision, 3 randomly selected views are chosen between each adjacent pair of source views to compute $\mathcal{L}_{\text{refine}}$ and $\mathcal{L}_{\text{render}}$. The learning rate for deep supervised pre-training and overall network training is initialized to 0.0002 and decreases linearly with the number of training epochs.

B. Real-world Data Assessment

We evaluate our model on HuMMan [1], a real-world dataset captured with RGB cameras at a resolution of 1024. Using images from the front two cameras (ID: 1 and ID: 9) as inputs, we infer the 3D Gaussians through EVA-Gaussian and render novel views from the viewpoint of ID: 0. The visualization results, as illustrated in Fig. 1, demonstrate that EVA-Gaussian produces high-quality novel view im-



Figure 1. Visualization results on real-world data. Minor artifacts on the human boundary mainly arise from the noisy human mask. Notably, GPS-Gaussian cannot generate reasonable outcome.

ages in real-world settings. Notably, GPS-Gaussian fails to produce reasonable results due to the sparse input view angles of only 90° , which further underscores the robustness of our approach.

C. More Visualization Results

In this section, we present additional visualization results in Fig. 2 to compare our method with SOTA approaches GPS-Gaussian and ENeRF on the Thuman2.0 [5] and Thumansit [6] 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 contrast, 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.

D. 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 $o := [o_1, \cdots, o_i, \cdots, o_N] \in \mathbb{R}^N$ of all considered 3D Gaussians and the collection of all 3D Gaussian scaling

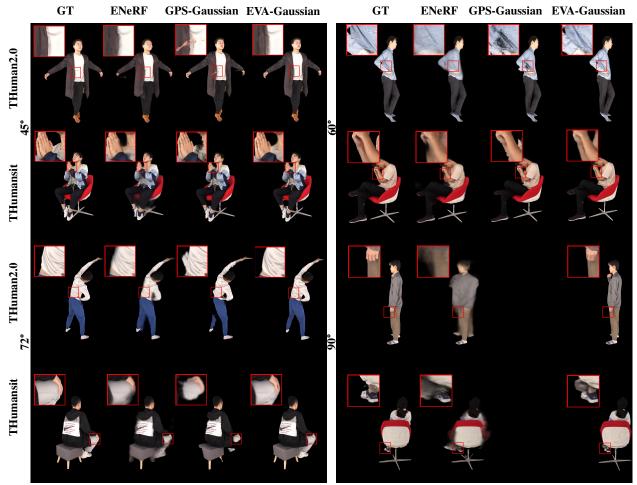


Figure 2. Qualitative comparison on THuman2.0 and THumansit. EVA-Gaussian demonstrates superior novel view rendering quality under diverse camera settings.

factors as:

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

where $s^1:=[s^1_1,s^1_2,\cdots,s^1_N]\in\mathbb{R}^N,\ s^2:=[s^2_1,s^2_2,\cdots,s^2_N]\in\mathbb{R}^N$ and $s^3:=[s^3_1,s^3_2,\cdots,s^3_N]\in\mathbb{R}^N.$ For the 3D Gaussian with the i-th greatest depth, the associated scaling matrix is constructed from the corresponding scaling factors as:

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

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

$$D(\mathbf{x}) = \sum_{i=1}^{N} z_i o_i G_i'(\mathbf{x}) \prod_{j=1}^{i-1} (1 - o_j G_j'(\mathbf{x}))), \quad (3)$$

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.

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

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

Theorem D.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:

$$\lim_{\substack{\boldsymbol{o} \to \mathbf{1} \\ \hat{\mathbf{S}} \to \mathbf{0}^+}} D(\boldsymbol{x}') = Z(\boldsymbol{x}'). \tag{4}$$

Theorem D.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 D.1, we introduce the following lemma from the well-known Moore-Osgood Theorem in [3]:

Lemma D.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 $p\in\mathbb{N}$ 10 and $\gamma_{k,\infty}:=\lim_{p\to\infty}\gamma_{k,p}$ exists for every $k\in\mathbb{N}$. 11 If (i) $\lim_{p\to\infty}\sup_{k\in\mathbb{N}}d_{\Gamma}(\gamma_{k,p},\gamma_{k,\infty})=0$ and (ii) 12 $\lim_{k\to\infty}d_{\Gamma}(\gamma_{k,p},\gamma_{\infty,p})=0, \forall p\in\mathbb{N}$, then the joint 13 limit $\lim_{k,p\to\infty}\gamma_{k,p}$ exists. In particular, it holds that 14 $\lim_{k,p\to\infty}\gamma_{k,p}=\lim_{p\to\infty}\gamma_{\infty,p}=\lim_{k\to\infty}\gamma_{k,\infty}$. 15

Lemma D.1 can be regarded as a special case of Theorem 16 7.11 from [4]. This lemma states that for a doubly-indexed 17

sequence, if the sequence converges uniformly with respect to one 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 D.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:

$$\lim_{\substack{o \to 1 \\ \hat{S} \to 0^{+}}} D(\mathbf{x}') = \lim_{\substack{o \to 1 \\ \hat{S} \to 0^{+}}} \sum_{i=1}^{N} (z_{i} o_{i} G_{i}'(\mathbf{x}') \prod_{j=1}^{i-1} (1 - o_{j} G_{j}'(\mathbf{x}')))$$

$$= \lim_{\substack{o \to 1 \\ \hat{S} \to 0^{+}}} \sum_{i=1}^{N} (z_{i} o_{i} e^{-\frac{1}{2} (\mathbf{x}' - \boldsymbol{\mu}_{i})^{T} (\mathbf{J} \mathbf{W} \mathbf{R}_{i} \mathbf{S}_{i} \mathbf{S}_{i}^{T} \mathbf{R}_{i}^{T} \mathbf{W}^{T} \mathbf{J}^{T})^{-1} (\mathbf{x}' - \boldsymbol{\mu}_{i})}$$

$$= \lim_{\substack{o \to 1 \\ \hat{S} \to 0^{+}}} \sum_{i=1}^{N} (z_{i} o_{i} e^{-\frac{1}{2} (\mathbf{x}' - \boldsymbol{\mu}_{j})^{T} (\mathbf{J} \mathbf{W} \mathbf{R}_{j} \mathbf{S}_{j} \mathbf{S}_{j}^{T} \mathbf{R}_{j}^{T} \mathbf{W}^{T} \mathbf{J}^{T})^{-1} (\mathbf{x}' - \boldsymbol{\mu}_{j})}$$

$$= \lim_{\substack{i=1 \\ \hat{S} \to 0^{+}}} \lim_{\substack{o \to 1 \\ o \to 1}} \sum_{i=1}^{N} (z_{i} o_{i} e^{-\frac{1}{2} (\mathbf{x}' - \boldsymbol{\mu}_{i})^{T} (\mathbf{J} \mathbf{W} \mathbf{R}_{i} \mathbf{S}_{i} \mathbf{S}_{i}^{T} \mathbf{R}_{i}^{T} \mathbf{W}^{T} \mathbf{J}^{T})^{-1} (\mathbf{x}' - \boldsymbol{\mu}_{i})}$$

$$= \lim_{\substack{i=1 \\ j=1}} (1 - o_{j} e^{-\frac{1}{2} (\mathbf{x}' - \boldsymbol{\mu}_{j})^{T} (\mathbf{J} \mathbf{W} \mathbf{R}_{j} \mathbf{S}_{j} \mathbf{S}_{j}^{T} \mathbf{R}_{j}^{T} \mathbf{W}^{T} \mathbf{J}^{T})^{-1} (\mathbf{x}' - \boldsymbol{\mu}_{j})}),$$

$$(7)$$

where (a) is from Lemma D.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 $\hat{\boldsymbol{S}} \to \mathbf{0}^+$ and as $\boldsymbol{o} \to \mathbf{1}$. This implies that the joint limit of \boldsymbol{o}

and s can be decomposed into the separate limits of o and s. Thus, we have:

$$\lim_{\hat{S} \to 0^{+}} \lim_{o \to 1} \sum_{i=1}^{N} (z_{i} o_{i} e^{-\frac{1}{2} (x' - \mu_{i})^{T} (JW R_{i} S_{i} S_{i}^{T} R_{i}^{T} W^{T} J^{T})^{-1} (x' - \mu_{i})}$$

$$\prod_{j=1}^{i-1} (1 - o_{j} e^{-\frac{1}{2} (x' - \mu_{j})^{T} (JW R_{j} S_{j} S_{j}^{T} R_{j}^{T} W^{T} J^{T})^{-1} (x' - \mu_{j})}))$$

$$= \lim_{\hat{S} \to 0^{+}} \sum_{i=1}^{N} (\lim_{o_{i} \to 1} z_{i} o_{i} e^{-\frac{1}{2} (x' - \mu_{i})^{T} (JW R_{i} S_{i} S_{i}^{T} R_{i}^{T} W^{T} J^{T})^{-1} (x' - \mu_{i})}$$

$$\lim_{(o_{j}, \dots, o_{i-1}) \to 1} \prod_{j=1}^{i-1} (1 - o_{j} e^{-\frac{1}{2} (x' - \mu_{j})^{T} (JW R_{j} S_{j} S_{j}^{T} R_{j}^{T} W^{T} J^{T})^{-1} (x' - \mu_{j})}))$$

$$= \lim_{\hat{S} \to 0^{+}} \sum_{i=1}^{N} (z_{i} e^{-\frac{1}{2} (x' - \mu_{i})^{T} (JW R_{i} S_{i} S_{i}^{T} R_{i}^{T} W^{T} J^{T})^{-1} (x' - \mu_{i})}$$

$$\prod_{j=1}^{i-1} (1 - e^{-\frac{1}{2} (x' - \mu_{j})^{T} (JW R_{j} S_{j} S_{j}^{T} R_{j}^{T} W^{T} J^{T})^{-1} (x' - \mu_{j})})).$$

$$(9)$$

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:

$$\mathbf{S}_{i}^{'} := \begin{pmatrix} s_{i}^{1} & 0 & 0\\ 0 & s_{i}^{1} & 0\\ 0 & 0 & s_{i}^{1} \end{pmatrix}. \tag{10}$$

Consequently, we have:

$$\lim_{\hat{S} \to 0^{+}} \sum_{i=1}^{N} (z_{i}e^{-\frac{1}{2}(x'-\mu_{i})^{T}}(JWR_{i}S_{i}S_{i}^{T}R_{i}^{T}W^{T}J^{T})^{-1}(x'-\mu_{i}))$$

$$\prod_{j=1}^{i-1} (1 - e^{-\frac{1}{2}(x'-\mu_{j})^{T}}(JWR_{j}S_{j}S_{j}^{T}R_{j}^{T}W^{T}J^{T})^{-1}(x'-\mu_{j})))$$

$$= \lim_{\hat{S} \to 0^{+}} \sum_{i=1}^{N} (z_{i}e^{-\frac{1}{2}(x'-\mu_{i})^{T}}(JWR_{i}S_{i}S_{i}^{T}R_{i}^{T}W^{T}J^{T})^{-1}(x'-\mu_{i}))$$

$$\prod_{j=1}^{i-1} (1 - e^{-\frac{1}{2}(x'-\mu_{j})^{T}}(JWR_{j}S_{j}S_{j}^{T}R_{j}^{T}W^{T}J^{T})^{-1}(x'-\mu_{j})))$$

$$= \lim_{s^{1} \to 0^{+}} \sum_{i=1}^{N} (z_{i}e^{-\frac{1}{2}(x'-\mu_{i})^{T}}(JWR_{i}S_{i}'S_{i}'^{T}R_{i}^{T}W^{T}J^{T})^{-1}(x'-\mu_{i}))$$

$$\prod_{j=1}^{i-1} (1 - e^{-\frac{1}{2}(x'-\mu_{j})^{T}}(JWR_{j}S_{j}'S_{j}'^{T}R_{j}^{T}W^{T}J^{T})^{-1}(x'-\mu_{j}))).$$
(12)

From (12), we see that when $x' = \mu_i$, it gives

$$e^{-\frac{1}{2}(\mathbf{x}' - \mu_i)^T (\mathbf{J} \mathbf{W} \mathbf{R}_i \mathbf{S}_i' \mathbf{S}_i'^T \mathbf{R}_i^T \mathbf{W}^T \mathbf{J}^T)^{-1} (\mathbf{x}' - \mu_i)} = 1.$$
(13)

Otherwise, if $x' \neq \mu_i$, we have

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

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

$$= 0. \tag{14}$$

By combining Eq. 5 - Eq. 14, we have

$$\lim_{\substack{o \to 1 \\ \hat{S} \to 0^{+}}} D(\mathbf{x}') = \sum_{i=1}^{N} \left(\lim_{s_{i}^{1} \to 0^{+}} z_{i} e^{-\frac{1}{2} (\mathbf{x}' - \boldsymbol{\mu}_{i})^{T} (\mathbf{J} \mathbf{W} \mathbf{R}_{i} \mathbf{S}'_{i} \mathbf{S}'_{i}^{T} \mathbf{R}_{i}^{T} \mathbf{W}^{T} \mathbf{J}^{T})^{-1} (\mathbf{x}' - \boldsymbol{\mu}_{i})} \right)$$

$$\prod_{j=1}^{i-1} \lim_{s_{j}^{1} \to 0^{+}} \left(1 - e^{-\frac{1}{2} (\mathbf{x}' - \boldsymbol{\mu}_{j})^{T} (\mathbf{J} \mathbf{W} \mathbf{R}_{j} \mathbf{S}'_{j} \mathbf{S}'_{j}^{T} \mathbf{R}_{j}^{T} \mathbf{W}^{T} \mathbf{J}^{T})^{-1} (\mathbf{x}' - \boldsymbol{\mu}_{j})} \right)$$

$$= Z(\mathbf{x}'),$$
(15)

which completes the proof.

18 References

37

38

39

40

41

- [1] Zhongang Cai, Daxuan Ren, Ailing Zeng, Zhengyu Lin, Tao 19 Yu, Wenjia Wang, Xiangyu Fan, Yang Gao, Yifan Yu, Liang 20 Pan, Fangzhou Hong, Mingyuan Zhang, Chen Change Loy, 21 22 Lei Yang, and Ziwei Liu. HuMMan: Multi-modal 4d human dataset for versatile sensing and modeling. In 17th European 23 Conference on Computer Vision, Tel Aviv, Israel, October 23-24 27, 2022, Proceedings, Part VII, pages 557-577. Springer, 25 2022. 1 26
- [2] 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, pages 770–778, 2016.
- [3] 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.
- [4] Walter Rudin et al. Principles of mathematical analysis.
 McGraw-hill New York, 1964. 3
 - [5] 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), 2021. 1
- [6] 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. 1