# GPS-Gaussian+: Generalizable Pixel-wise 3D Gaussian Splatting for Real-Time Human-Scene Rendering from Sparse Views

Boyao Zhou*, Shunyuan Zheng*, Hanzhang Tu, Ruizhi Shao, Boning Liu, Shengping Zhang, Liqiang Nie and Yebin Liu

https://yaourtb.github.io/GPS-Gaussian+

**Abstract**—Differentiable rendering techniques have recently shown promising results for free-viewpoint video synthesis of characters. However, such methods, either Gaussian Splatting or neural implicit rendering, typically necessitate per-subject optimization which does not meet the requirement of real-time rendering in an interactive application. We propose a generalizable Gaussian Splatting approach for high-resolution image rendering under a sparse-view camera setting. To this end, we introduce Gaussian parameter maps defined on the source views and directly regress Gaussian properties for instant novel view synthesis without any fine-tuning or optimization. We train our Gaussian parameter regression module on human-only data or human-scene data, jointly with a depth estimation module to lift 2D parameter maps to 3D space. The proposed framework is fully differentiable with both depth and rendering supervision or with only rendering supervision. We further introduce a regularization term and an epipolar attention mechanism to preserve geometry consistency between two source views, especially when neglecting depth supervision. Experiments on several datasets demonstrate that our method outperforms state-of-the-art methods while achieving an exceeding rendering speed.

**Index Terms**—3D Gaussian Splatting, Novel View Synthesis, Free Viewpoint Video

---✦---

## 1 INTRODUCTION

FREE-Viewpoint Video (FVV) synthesis from sparse input views is a challenging and crucial task in computer vision, which is largely used in sports broadcasting, stage performance and telepresence systems [1], [2]. However, early attempts [3], [4] try to solve this problem through a weighted blending mechanism [5] by using a huge number of cameras, which dramatically increases computational cost and latency. On the other hand, NeRF-like differentiable volumetric rendering techniques [6], [7], [8], [9] can synthesize novel views under sparse camera setting [10], but typically suffer from per-scene optimization [6], [7], [8], [9], slow rendering speed [6], [7] and overfitting to input views [11].

In contrast, point-based rendering [16], [17], [18], [19] has drawn long-lasting attention thanks to its high-speed, and even real-time, rendering performance. Once integrated with neural networks, point-based graphics [20], [21] realize a promising explicit representation with comparable realism and extremely superior efficiency in FVV tasks [20], [21]. Recently, 3D Gaussian Splatting (3D-GS) [14] introduces a new representation that the point clouds are formulated as 3D Gaussians with a series of learnable properties including 3D position, color, opacity and anisotropic covariance. By applying $\alpha$-blending [22], 3D-GS provides not only a more reasonable and accurate mechanism for back-propagating
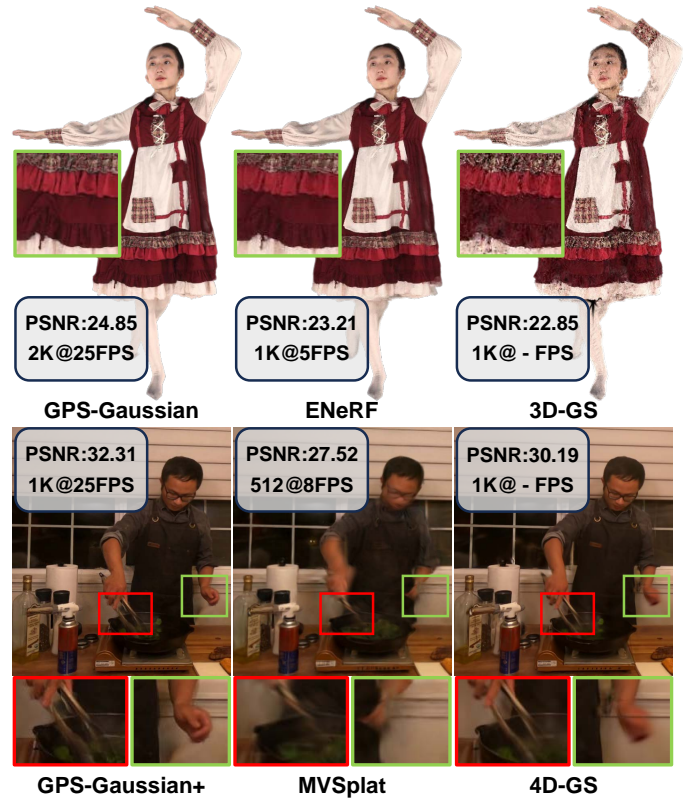
---

* indicates equal contribution.
*Boyao Zhou, Hanzhang Tu, Ruizhi Shao, Boning Liu and Yebin Liu are with Department of Automation, Tsinghua University, Beijing 100084, P.R.China. Shunyuan Zheng and Shengping Zhang are with the School of Computer Science and Technology, Harbin Institute of Technology, Weihai 264209, P.R.China. Liqiang Nie is with the School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen 518055, P.R.China. Corresponding author: Shengping Zhang (s.zhang@hit.edu.cn).*



Fig. 1: **High-fidelity and real-time rendering.** On the top, GPS-Gaussian produces $2K$-resolution rendering of character, while GPS-Gaussian+ renders novel views of human-centered scenes on the bottom. Our methods outperform the state-of-the-art feed-forward implicit rendering method ENeRF [12], explicit rendering method MVSplat [13] and optimization-based methods 3D-GS [14] and 4D-GS [15].

the gradients but also a real-time rendering efficiency for complex scenes. Despite realizing a real-time inference, Gaussian Splatting relies on per-scene [14] or per-frame [23] parameter optimization for several minutes. It is therefore impractical in interactive scenarios as it necessitates the re-optimization of Gaussian parameters once the character or the scene changes.

More recently, some generalizable Gaussian Splatting methods [13], [24], [25], [26] have been proposed to explore novel-view synthesis in a feed-forward way. In general, such methods leverage a learnable geometry prior with a differentiable rendering pipeline to achieve feed-forward inference. For example, Splatter Image [24] regresses directly Gaussians' positions and other properties from a single image. In such an ill-posed setting, 3D consistency is hardly held and image quality is extremely low. PixelSplat [25], MVSplat [13] and MVSGaussian [26] utilize probabilistic depth estimation and multiplane sweeping to represent geometry from multiple source view images. Although such geometry cues allow 3D-GS to generalize to some static scenes, the probabilistic representations can easily make floating artifacts due to uncertain Gaussian positions. More-over, they dramatically increase the computational cost for inference and no longer render novel view image in real-time, even with a low resolution of $256 \times 256$.

In this paper, we propose to integrate binocular stereo-matching [27], [28], [29] as a geometry cue with 3D-GS rendering pipeline to achieve a generalizable Gaussian Splatting. Given a pair of images, stereo-matching can determinately calculate the disparity by searching the cost volume built on the extracted features of two source views. Based on the epipolar geometry, it is straightforward to transform disparity into depth given camera parameters. This determinant geometry representation can be better supervised with rendering loss than probabilistic ones because Gaussian Splatting necessitates a certain position of each primitive.

Specifically, we introduce 2D Gaussian parameter (depth residual, color, scaling, rotation, opacity) maps which are defined on source view image planes, instead of unstructured point clouds. These Gaussian parameter maps allow us to represent a scene with pixel-wise parameters, *i.e.* each pixel corresponding to a specific Gaussian point. Additionally, it enables the application of efficient 2D convolution networks rather than expensive 3D operators. Given the estimated depth map via binocular stereo-matching, it is efficient to lift 2D parameter maps to 3D Gaussian points. Such unprojected Gaussian points from both the two source views constitute the representation of scene and novel view images can be rendered with splatting technique [14]. Since the whole pipeline is fully differentiable, we can jointly train an iterative stereo matching-based depth estimation [29] along with our Gaussian parameter regression with both depth and rendering loss or with only rendering loss.

A preliminary version of this work has been published as a highlight paper [30] in CVPR 2024, in which we propose a real-time framework of human novel view synthesis and train this framework on synthetic human-only data with depth and rendering loss. In the current version, we aim to extend it to human-scene scenarios and no longer require matting technique or depth supervision. Note that wrong matting leads to floating artifacts and it is not triv-

ial to acquire high-quality geometry or depth information from complex human-centered scene data. To achieve high-quality rendering without depth supervision, an epipolar attention is applied in the shared feature extraction module (Sec. 4.2), which can improve stereo-matching accuracy and rendering consistency. In addition, a depth residual map (Sec. 4.4.3) is introduced in order to recover high-frequency details from the predicted depth of stereo-matching when lacking depth supervision. Furthermore, we propose a regularization term in Sec. 4.5.2 to preserve geometry consistency between the two source views and improve the overall stability of the training process when lacking the ground truth of depth. Thanks to these adaptive components, our network can be trained with only rendering loss, making it scalable for more general human-scene scenarios.

In practice, we are able to synthesize high-fidelity free-viewpoint video around 25 FPS on a single modern graphics card. Leveraging the rapid rendering capabilities and broad generalizability inherent in our proposed method, an unseen character with or without background can be instantly rendered without necessitating any fine-tuning or optimization. In summary, our contributions can be summarized as:

- We introduce a generalizable 3D Gaussian Splatting methodology that employs pixel-wise Gaussian parameter maps defined on 2D source image planes to formulate 3D Gaussians in a feed-forward manner.
- We propose a fully differentiable framework composed of an iterative depth estimation module and a Gaussian parameter regression module. The intermediate depth prediction bridges the two components and allows them to benefit from joint training.
- We introduce a regularization term and an epipolar attention mechanism to preserve geometry consistency between the two source views when using only rendering loss. Our method generalizes well to unseen characters even in complicated scenes.
- We develop a real-time FVV system that achieves high-resolution rendering of characters in the scene without any geometry supervision.

## 2 RELATED WORK

**Neural Implicit Representation.** Neural implicit function has recently aroused a surge of interest to represent complicated scenes, in form of occupancy fields [31], [32], [33], [34], radiance fields [6], [7], [35], [36], [37], [38] and signed distance functions [10], [39], [40], [41], [42]. Implicit representation shows the advantage in memory efficiency and topological flexibility for human representation [34], [43], [44] or scene reconstruction [45], [46], especially in a pixel-aligned feature query manner [32], [33]. However, each queried point is processed through the full network, which dramatically increases computational complexity. More recently, numerous methods have extended Neural Radiance Fields (NeRF) [6] to static human modeling [47], [48] and dynamic human modeling from sparse multi-view cameras [7], [10], [36] or a monocular camera [37], [38], [49]. However, these methods typically require a per-subject optimization process and it is non-trivial to generalize these methods to unseen subjects. Previous attempts, *e.g.*, PixelNeRF [11],

IBRNet [50], MVSNeRF [51] and ENeRF [12] resort to image-based features as potent priors for feed-forward scene modeling. Despite the great progress in accelerating the scene-specific NeRF [8], [9], [52], [53], efficient generalizable NeRF for interactive scenarios remains to be further elucidated.

**Deep Image-based Rendering.** Image-based rendering, or IBR in short, synthesizes novel views from a set of multi-view images with a weighted blending mechanism, which is typically computed from a geometry proxy. [54], [55] deploy multi-view stereo from dense input views to produce mesh surfaces as a proxy for image warping. DNR [56] directly produces learnable features on the surface of mesh proxies for neural rendering. Obtaining these proxies is not straightforward since high-quality multi-view stereo and surface reconstruction requires dense input views. Point clouds from SfM [57], [58] or depth sensors [59], [60] can also be engaged as geometry proxies. These methods highly depend on the performance of 3D reconstruction algorithms or the quality of depth sensors. FWD [61] designs a network to refine depth estimations, then explicitly warps pixels from source views to novel views with the refined depth maps. FloRen [62] utilizes a coarse human mesh reconstructed by PIFu [32] to render initialized depth maps for novel views. Arguably FloRen [62] is most related to our preliminary work GPS-Gaussian [30], as it also realizes $360°$ free view human performance rendering in real-time. However, the appearance flow in FloRen merely works in 2D domains, where the rich geometry cues and multi-view geometric constraints only serve as 2D supervisions. The difference is that our approach lifts 2D priors into 3D space and utilizes the point representation to synthesize novel views in a fully differentiable manner.

**Point-based Graphics.** Point-based representation has shown great efficiency and simplicity for various 3D human-centered tasks [63], [64], [65], [66], [67], [68], [69]. Previous attempts integrate point cloud representation with 2D neural rendering [20], [21] or NeRF-like volume rendering [70], [71]. Still, such a hybrid architecture does not exploit the rendering capability of point clouds and takes a long time to optimize on different scenes. Then differentiable point-based [18] and sphere-based [19], [69] rendering have been developed, which demonstrates promising rendering qualities, especially attaching them to a conventional network pipeline [60], [61]. In addition, isotropic points can be substituted by a more reasonable Gaussian point modeling [14], [23], to realize a rapid differentiable rendering framework with a splatting technique. This advanced representation has showcased prominent performance in concurrent 3D human-centered work [72], [73], [74], [75], [76]. However, a per-scene or per-subject optimization strategy limits its real-world application. Although [23], [77] accelerate partly the optimization process by using an on-the-fly strategy, they struggle to handle topology change in dynamic scenes. In this paper, we go further to generalize 3D Gaussians across diverse subjects while maintaining its fast and high-quality rendering properties.

**Free-Viewpoint Video.** Targeting different applications, there are two feasible schemes to produce free-viewpoint videos, one uses a compact 4D representation [10], [15], [78], [79], [80], [81], and the other formulate an individual 3D representation for each discrete timestamp, which can be further subdivided into on-the-fly optimization methods [23], [69], [77], [82] and feed-forward inference methods [26], [30], [83], [84]. The 4D representations [15], [79], [80], [81] cater to volumetric video, which can be played back and viewed from any viewpoint at any time, but the performance degrades when the capturing time goes longer. On the contrary, the on-the-fly optimization method [82] excels at handling long-time sequences. They can realize similar experiences after applying customized compression designs but they typically have higher memory costs than 4D methods. Nevertheless, despite having been accelerated, the essential optimization process is far from real-time. Thus, we orient towards feed-forward methods for interactive scenarios. Among them, MonoFVV [85] and Function4D [86] implement RGBD fusion with depth sensors to attain real-time human rendering. The large variation in pose and clothing makes the feed-forward generalizable free-view rendering a more challenging task, thus recent work [48], [83], [84], [87], [88] simplifies the problem by leveraging human priors [32], [89]. However, an inaccurate prior estimation would mislead the final result. For more general dynamic scenarios, [12] relies on expensive probabilistic geometry estimation, thus they can hardly achieve real-time free-viewpoint rendering, even integrated with Gaussian Splatting [13], [25], [26].

## 3 PRELIMINARY

Since the proposed GPS-Gaussian+ harnesses the power of 3D-GS [14], we give a brief introduction in this section.

3D-GS models a static 3D scene explicitly with point primitives, each of which is parameterized as a scaled Gaussian with 3D covariance matrix $\boldsymbol{\Sigma}$ and mean $\mu$

$$G(\mathcal{X}) = e^{-\frac{1}{2}(\mathcal{X}-\mu)^T \boldsymbol{\Sigma}^{-1}(\mathcal{X}-\mu)} \tag{1}$$

In order to be effectively optimized by gradient descent, the covariance matrix $\boldsymbol{\Sigma}$ can be decomposed into a scaling matrix $\mathbf{S}$ and a rotation matrix $\mathbf{R}$ as

$$\boldsymbol{\Sigma} = \mathbf{R}\mathbf{S}\mathbf{S}^T\mathbf{R}^T \tag{2}$$

Following [90], the projection of Gaussians from 3D space to a 2D image plane is implemented by a view transformation $\mathbf{W}$ and the Jacobian of the affine approximation of the projective transformation $\mathbf{J}$. The covariance matrix $\boldsymbol{\Sigma}'$ in 2D space can be computed as

$$\boldsymbol{\Sigma}' = \mathbf{J}\mathbf{W}\boldsymbol{\Sigma}\mathbf{W}^T\mathbf{J}^T \tag{3}$$

followed by a point-based alpha-blend rendering which bears similarities to that used in NeRF [6], formulated as

$$\mathbf{C}_{color} = \sum_{i \in N} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1}(1 - \alpha_i) \tag{4}$$

where $\mathbf{c}_i$ is the color of each point, and density $\alpha_i$ is reasoned by the multiplication of a 2D Gaussian with covariance $\boldsymbol{\Sigma}'$ and a learned per-point opacity [91]. The color is defined by spherical harmonics (SH) coefficients in [14].

To summarize, the original 3D Gaussians methodology characterizes each Gaussian point by the following attributes: (1) a 3D position of each point $\mathcal{X} \in \mathbb{R}^3$, (2) a color defined by SH $\mathbf{c} \in \mathbb{R}^k$ (where $k$ is the freedom of SH basis), (3) a rotation parameterized by a quaternion $\mathbf{r} \in \mathbb{R}^4$, (4) a scaling factor $\mathbf{s} \in \mathbb{R}^3_+$, and (5) an opacity $\alpha \in [0, 1]$.
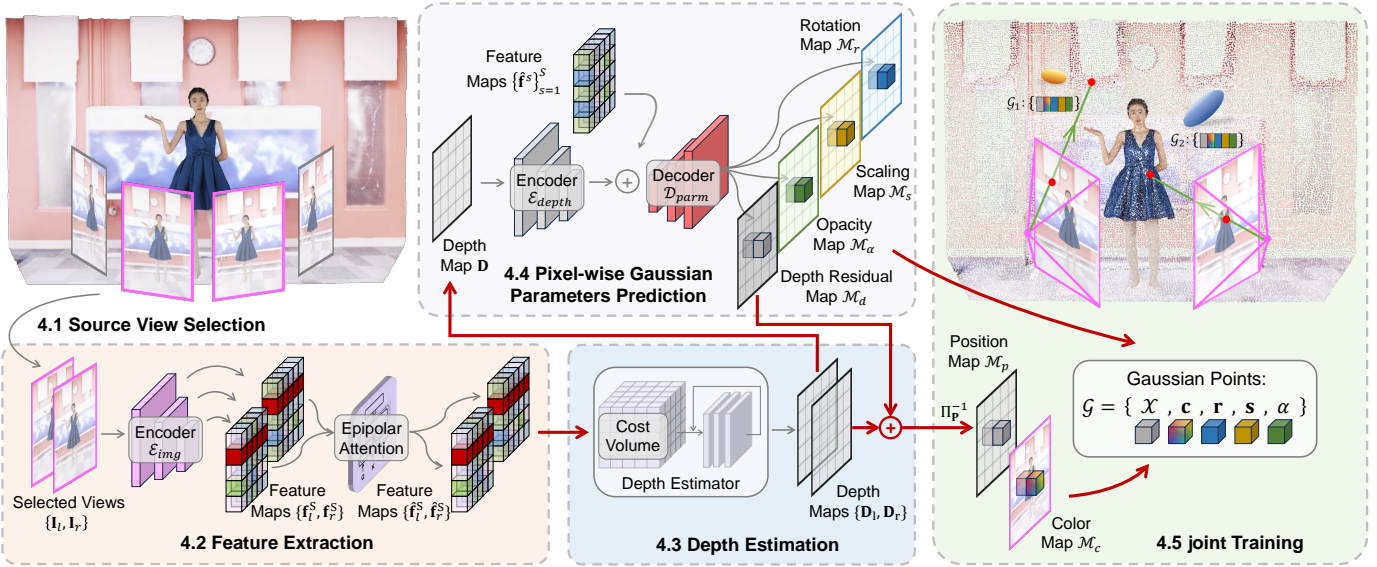
Fig. 2: **Overview of GPS-Gaussian+.** Given RGB images of a human-centered scene with sparse camera views and a target novel viewpoint, we select the adjacent two views on which to formulate our pixel-wise Gaussian representation. We extract the image features by using epipolar attention and then conduct an iterative depth estimation. For each source view, the RGB image serves as a color map, while the other parameters of 3D Gaussians are predicted in a pixel-wise manner. The Gaussian parameter maps defined on 2D image planes of both views are further unprojected to 3D space via refined depth maps and aggregated for novel view rendering. The fully differentiable framework enables a joint training mechanism with only rendering loss and geometry regularization.

## 4 METHOD

The overview of our method is illustrated in Fig. 2. Given RGB images of a human-centered scene with sparse camera views, our method aims to generate high-quality free-viewpoint video of the performer in real-time. Once given a target novel viewpoint, we select the two neighboring views from sparse cameras (Sec. 4.1). Then, image features are extracted from the two input images with a shared image encoder by using epipolar attention (Sec. 4.2), and they are further used to predict the depth maps for both source views with a binocular depth estimator (Sec. 4.3). The colors of 3D Gaussians are directly determined by the corresponding source view pixels, while other parameters of 3D Gaussians are predicted in a pixel-wise manner when feeding the predicted depth values and the former image features into a network (Sec. 4.4). Combined with RGB map of the source view image, these parameter maps formulate the Gaussian representation in 2D image planes and are further unprojected to 3D space with the estimated depth. The unprojected Gaussians from both views are aggregated and rendered to the target viewpoint in a differentiable way, which allows for end-to-end training (Sec. 4.5), even with only rendering loss.

### 4.1 Source View Selection

As a binocular stereo method, we synthesize the target novel view with two adjacent source views. Given $N$ input images $\{\mathbf{I}_n\}_{n=1}^N$, with their camera position $\{C_n\}_{n=1}^N$, source views can be represented by $\mathbf{V_n} = C_n - O$, where $O$ is the center of the scene. Similarly, the target novel view rendering can be defined as $I_{tar}$ with camera position $C_{tar}$ and view $\mathbf{V_{tar}} = C_{tar} - O$. By conducting a dot product of all

source view vectors and the novel view vector, the nearest two views $(v_l, v_r)$ can be selected as the 'working set' of binocular stereo, where $l$ and $r$ stand for 'left' and 'right' view, respectively.

### 4.2 Feature Extraction

The selected images are encoded with a feature extraction module in order to search the corresponding features from one view to another. Once two source view images are rectified, $\mathbf{I}_l, \mathbf{I}_r \in [0,1]^{H \times W \times 3}$ are fed to a shared image encoder $\mathcal{E}_{img}$ with several residual blocks and downsampling layers to extract dense feature maps $\mathbf{f}^s \in \mathbb{R}^{H/2^s \times W/2^s \times D_s}$ where $D_s$ is the dimension at the $s$-th feature scale

$$\langle \{\mathbf{f}_l^s\}_{s=1}^S, \{\mathbf{f}_r^s\}_{s=1}^S \rangle = \mathcal{E}_{img}(\mathbf{I}_l, \mathbf{I}_r) \tag{5}$$

where we set $S = 3$ in our experiments.

Since the image encoder $\mathcal{E}_{img}$ is independent of each other view, it struggles to extract informative features when lacking depth supervision. Thus, we propose to conduct an epipolar attention module on the bottleneck features $\mathbf{f}_{l,r}^S$, in order to exchange useful information from each other view. Note that corresponding pixels from the two rectified images are located on the same horizontal epipolar line. In practice, we rearrange feature map into $H/2^S$ line features $\mathbf{f}^e \in \mathbb{R}^{W/2^S \times D_S}$ and employ multi-head attention [92] $Att$ along each epipolar line

$$\langle \mathbf{Q}, \mathbf{K}, \mathbf{V} \rangle = \langle \mathbf{f}^e \mathbf{W}^Q, \mathbf{f}^e \mathbf{W}^K, \mathbf{f}^e \mathbf{W}^V \rangle$$
$$\hat{\mathbf{f}}_i = \mathbf{f}_i^e + Att(\mathbf{Q}_i, \mathbf{K}_j, \mathbf{V}_j) \tag{6}$$

where $\{i, j\} = \{l, r\}$ or $\{r, l\}$. Processed line features $\{\hat{\mathbf{f}}_k\}_{k=1}^{H/2^S}$ of source view are concatenated into feature map

$\hat{\mathbf{f}}^S \in \mathbb{R}^{H/2^S \times W/2^S \times D_S}$ following some convolution operations. Such an attention-based encoder $\mathcal{E}^{att}$ significantly increases the perceptive field of the feature extractor so that the extracted features can be used to build cost volume in the following section.

## 4.3 Depth Estimation

The depth map is the key component of our framework which lifts the 2D image planes to 3D Gaussian representation. Note that, depth estimation in binocular stereo is equivalent to disparity estimation. For each pixel coordinate $x = (u, v)$ in one view, disparity estimation $\phi_{disp}$ aims to find its corresponding coordinate $(u + \phi_{disp}(u), v)$ in another view, considering the displacement of each pixel is constrained to a horizontal line in rectified stereo. Since there is a one-to-one mapping between disparity maps and depth maps given camera parameters, we do not distinguish them in the following sections. Inspired by [29], we implement this module in an iterative manner mainly because it avoids using prohibitively slow 3D convolutions to filter the cost volume. Given the processed feature maps $\hat{\mathbf{f}}_l^S, \hat{\mathbf{f}}_r^S$, we compute a 3D cost volume $\mathbf{C} \in \mathbb{R}^{H/2^S \times W/2^S \times W/2^S}$ using matrix multiplication

$$\mathbf{C}(\hat{\mathbf{f}}_l^S, \hat{\mathbf{f}}_r^S), \quad C_{ijk} = \sum_h (\hat{\mathbf{f}}_l^S)_{ijh} \cdot (\hat{\mathbf{f}}_r^S)_{ikh} \quad (7)$$

Then, an iterative update mechanism predicts a sequence of depth estimations $\{\mathbf{d}_l^t\}_{t=1}^T$ and $\{\mathbf{d}_r^t\}_{t=1}^T$ by looking up in volume $\mathbf{C}$, where $T$ is the update iterations. For more details about the update operators, please refer to [93]. The outputs of the final iterations $(\mathbf{d}_l^T, \mathbf{d}_r^T)$ are upsampled to full image resolution via a convex upsampling. The depth estimation module $\Phi_{depth}$ can be formulated as

$$\langle \mathbf{D}_l, \mathbf{D}_r \rangle = \Phi_{depth}(\hat{\mathbf{f}}_l^S, \hat{\mathbf{f}}_r^S, K_l, K_r) \quad (8)$$

where $K_l$ and $K_r$ are the camera parameters, $\mathbf{D}_l, \mathbf{D}_r \in \mathbb{R}^{H \times W \times 1}$ are the depth estimations. The classic binocular stereo methods estimate the depth for 'reference view' only, while we pursue depth maps for both input views with a shared-weight network to serve as the position of Gaussian points, which results in a decent efficiency increase.

## 4.4 Pixel-wise Gaussian Parameters Prediction

In 3D-GS [14], each Gaussian point in 3D space is characterized by attributes $\mathcal{G} = \{\mathcal{X}, \mathbf{c}, \mathbf{r}, \mathbf{s}, \alpha\}$, which represent 3D position, color, rotation, scaling and opacity, respectively. In this section, we introduce a pixel-wise manner to formulate 3D Gaussians in 2D image planes. Specifically, the proposed Gaussian maps $\mathbf{G}$ are defined as

$$\mathbf{G}(x) = \{\mathcal{M}_p(x), \mathcal{M}_c(x), \mathcal{M}_r(x), \mathcal{M}_s(x), \mathcal{M}_\alpha(x)\} \quad (9)$$

where $x$ is the coordinate of a valid pixel in a rectified image plane, $\mathcal{M}_p, \mathcal{M}_c, \mathcal{M}_r, \mathcal{M}_s, \mathcal{M}_\alpha$ represents Gaussian parameter maps of position, color, rotation, scaling and opacity, respectively.

### 4.4.1 Color Map

Considering our human-centered scenario is predominantly characterized by diffuse reflection, instead of predicting the sphere harmonic (SH) coefficients, we directly use the source RGB image as the color map

$$\mathcal{M}_c(x) = \mathbf{I}(x) \quad (10)$$

### 4.4.2 Rotation, Scaling and Opacity Map

The remaining Gaussian parameters are related not only to the extracted features $\{\mathbf{f}^s\}_{s=1}^S$ in Sec. 4.2 but also to the spatial cues from estimated depth in Sec. 4.3. The former one provides a global context with image encoder $\mathcal{E}_{img}^{att}$ and the latter one should focus on structural details so that Gaussian parameters can be predicted in a feed-forward manner. Hence, we construct an additional encoder $\mathcal{E}_{depth}$, which takes the depth map $\mathbf{D}$ as input, to complement the coarse geometric awareness for each pixel. The image features and the spatial features are fused by a U-Net like decoder $\mathcal{D}_{parm}$ to regress pixel-wise Gaussian features in full image resolution

$$\mathbf{\Gamma} = \mathcal{D}_{parm}(\mathcal{E}_{img}^{att}(\mathbf{I}) \oplus \mathcal{E}_{depth}(\mathbf{D})) \quad (11)$$

where $\mathbf{\Gamma} \in \mathbb{R}^{H \times W \times D_G}$ is Gaussian features, $\oplus$ stands for concatenations at all feature levels. The prediction heads, each composed of two convolution layers, are adapted to Gaussian features for specific Gaussian parameter map regression. Before being used to formulate Gaussian representations, the rotation map should be normalized since it represents a quaternion

$$\mathcal{M}_r(x) = Norm(h_r(\mathbf{\Gamma}(x))) \quad (12)$$

where $h_r$ is the rotation head. The scaling map and the opacity map need activations to satisfy their range

$$\begin{aligned} \mathcal{M}_s(x) &= Softplus(h_s(\mathbf{\Gamma}(x))) \\ \mathcal{M}_\alpha(x) &= Sigmoid(h_\alpha(\mathbf{\Gamma}(x))) \end{aligned} \quad (13)$$

where $h_s$ and $h_\alpha$ represent the scaling head and opacity head, respectively.

### 4.4.3 Depth Residual Map

We discover that the estimated depth in Sec. 4.3 is still coarse, especially in the case of the absence of depth supervision during training. To add high-frequency details onto coarse geometry, we further design a depth residual map

$$\mathcal{M}_d(x) = \gamma Tanh(h_d(\mathbf{\Gamma}(x))) \quad (14)$$

where $h_d$ represents the depth residual head and we use the $Tanh$ function with a scaling factor $\gamma = 0.5$ to activate the predicted value in a small range. Given the predicted depth map $\mathbf{D}$ in Eq. 8 and the residual value in Eq. 14, a pixel located at $x$ can be immediately unprojected from image planes to 3D space using projection matrix $\mathbf{P} \in \mathbb{R}^{3 \times 4}$ structured with camera parameters $K$

$$\mathcal{M}_p(x) = \Pi_{\mathbf{P}}^{-1}(x, \mathbf{D}(x) + \mathcal{M}_d(x)) \quad (15)$$

Finally, the learnable unprojection in Eq. 15 bridges 2D feature space and 3D Gaussian representation.

## 4.5 Joint Training

The pixel-wise Gaussian parameter maps defined on both source views are then lifted to 3D space and aggregated to render photo-realistic novel view images using the Gaussian Splatting technique in Sec. 3. Our whole framework is fully differentiable so that we jointly train depth estimation (Sec. 4.3) and Gaussian parameters prediction (Sec. 4.4) which typically benefit each other. The full pipeline can be trained with only rendering loss or a combination of depth loss and rendering loss when ground truth depth is available during training. When depth supervision is absent, we should also take into account the geometry consistency between two point clouds unprojected from left-view and right-view depth maps.

### 4.5.1 Training Loss

**Rendering loss.** First, we use rendering loss composed of L1 loss and SSIM loss [94], denoted as $\mathcal{L}_{mae}$ and $\mathcal{L}_{ssim}$ respectively, to measure the difference between the rendered and the ground truth image

$$\mathcal{L}_{render} = \lambda_1 \mathcal{L}_{mae} + \lambda_2 \mathcal{L}_{ssim} \tag{16}$$

where we set $\lambda_1 = 0.8$ and $\lambda_2 = 0.2$ in our experiments.
**Depth loss.** When ground truth depth is available, we minimize the L1 distance between the predicted and ground truth depth over the full sequence of predictions $\{\mathbf{d}^t\}_{t=1}^T$ with exponentially increasing weights, as shown in [29]. Given ground truth depth $\mathbf{d}_{gt}$, the loss is defined as

$$\mathcal{L}_{depth} = \sum_{t=1}^T \mu^{T-t} \|\mathbf{d}_{gt} - \mathbf{d}^t\|_1 \tag{17}$$

where we set $\mu = 0.9$ in our experiments.

### 4.5.2 Geometry Regularization

Ground truth depth is not trivially accessible, especially for complex scene data. Only rendering loss can not ensure geometry consistency between the two input views. Thus we try to minimize Chamfer distance between the unprojected Gaussian points of the two source views as a regularization term to boost the stereo-matching in two directions

$$\mathcal{L}_{CD} = \frac{1}{|\mathcal{P}^l|} \sum_{p_l \in \mathcal{P}^l} \min_{p_r \in \mathcal{P}^r} \|p_l - p_r\|_2 + \\ \frac{1}{|\mathcal{P}^r|} \sum_{p_r \in \mathcal{P}^r} \min_{p_l \in \mathcal{P}^l} \|p_r - p_l\|_2 \tag{18}$$

where $\mathcal{P}^{l,r}$ represents a valid Gaussian point set from the left or right view.

Overall, the final loss function is defined as $\mathcal{L} = \mathcal{L}_{render} + \alpha \mathcal{L}_{CD} + \beta \mathcal{L}_{depth}$. In practice, we set $\beta = 0$ when ground truth depth is not available.

## 5 EXPERIMENTS

### 5.1 Datasets and Metrics

**Human-Scene data.** To train and evaluate our network, we collect character performance data in the scene from DyNeRF [78] and ENeRF-outdoor [12] datasets. We take 4 motion sequences from DyNeRF, each of which contains 300 frames.

The first 220 frames are used as training data, and we evaluate our method on the rest of the data. For ENeRF-outdoor data, we take 4 motion sequences of 300 frames as training data, and 2 motion sequences of unseen characters as test data. We also capture motion sequences of single-character or multiple-character performance in 3 different scenes with a forward-facing camera rig, as shown in Fig. 13(a), to test the robustness of our method across scenes. In particular, 10 cameras are positioned in a line, spanning 1.6 meters. Four cameras with red circles in Fig. 13(a), are used as inputs which compose 3 pairs of source views, and the others serve as novel views during validation. For rendering continuity across source-view pairs, we use all 10 views as supervision during training. For each scene, our captured dataset consists of 3 sequences for training and 2 sequences for test, so there are 15 sequences in total. We train a model on each dataset and the models can generalize to unseen characters in the scene. For our captured data, our model is able to handle all three backgrounds. In terms of human-scene ratio, our data and ENeRF-outdoor capture full-body characters with small-focal cameras, while DyNeRF focuses on upper-body. Due to the original resolution of raw data, we set all images to $1K$ resolution.
**Human-only data.** To learn human priors from a large amount of data, we collect 1700 and 526 human scans from Twindom [95] and THuman2.0 [86], respectively. We randomly select 200 and 100 scans as validation data from Twindom and THuman2.0, respectively. In addition, we uniformly position 8 cameras in a cycle, thus the angle between two neighboring cameras is about $45°$. To test the robustness in real-world scenarios, we capture real data of 4 characters in the same 8-camera setup and prepare 8 additional camera views for evaluation, as shown in Fig. 13(b). For synthetic data, we render images on $2K$ resolution as rendering supervision during training and as ground truth during the test.
**Evaluation metric.** Following ENeRF [12], we evaluate our method and other baselines with PSNR, SSIM [94] and LPIPS [96] as metrics for the rendering results in valid regions of novel views.

### 5.2 Implementation Details

Our method is trained on a single RTX3090 graphics card using AdamW [97] optimizer with an initial learning rate of $2e^{-4}$. For real-captured human-scene data, we train the whole network from scratch for around $100k$ iterations with rendering loss and Chamfer distance. For DyNeRF [78] data, we set $\alpha = 0.005$ for its large range of depth, while we set $\alpha = 0.5$ for ENeRF-outdoor [12] and our captured data. For synthetic human-only data, THuman2.0 [86] and Twindom [95], we have 2 strategies to train networks. We can still train the whole network from scratch for $100k$ iterations with rendering loss and Chamfer distance. When depth information is available during training, the depth estimation module can be firstly trained for $40k$ iterations and we then jointly train depth estimation and Gaussian parameters prediction for $100k$ iterations.

### 5.3 Comparisons

**Baselines.** Considering that our goal is instant novel view synthesis, we compare our GPS-Gaussian+ against general-

TABLE 1: **Quantitative comparison on human-scene datasets.** All methods are evaluated on an RTX 3090 GPU to report the speed of synthesizing one novel view with two $1024 \times 1024$ source images, except MVSplat [13] with two $512 \times 512$ images due to memory cost. Our method uses TensorRT for fast inference. † 4D-GS [15] requires per-sequence optimization, while the other methods perform feed-forward inferences. The best, the second best and the third best are highlighted with different colors.

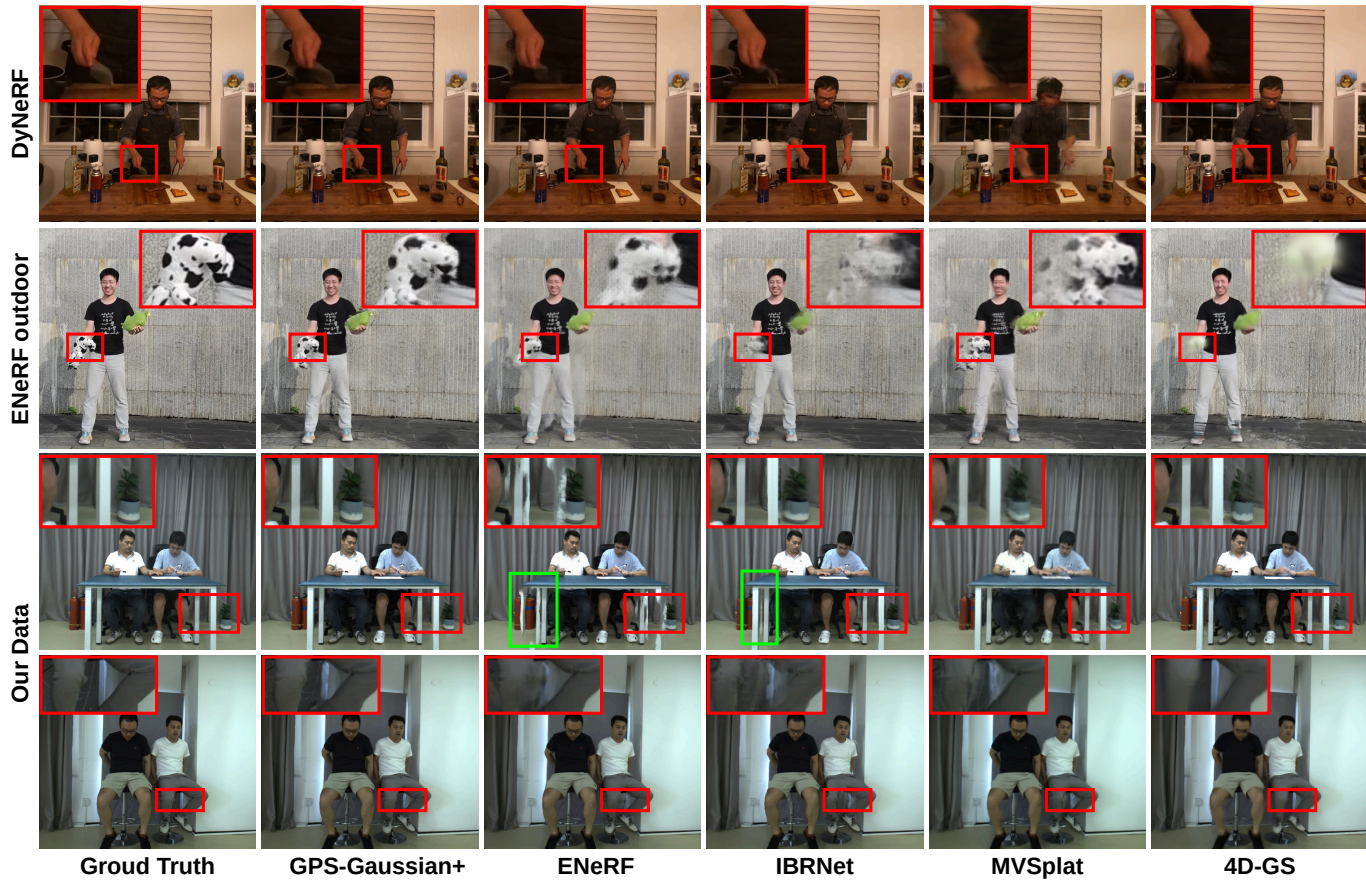| Method | DyNeRF [78] | | | ENeRF-outdoor [12] | | | Our Human-Scene Data | | | FPS |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | |
| 4D-GS [15]† | 31.67 | 0.954 | 0.057 | 21.47 | 0.480 | 0.302 | 28.65 | 0.906 | 0.112 | / |
| MVSplat [13] | 27.89 | 0.895 | 0.159 | 21.50 | 0.543 | 0.324 | 31.56 | 0.947 | 0.111 | 8 |
| IBRNet [50] | 32.10 | 0.944 | 0.067 | 24.19 | 0.626 | 0.269 | 32.09 | 0.948 | 0.077 | 0.25 |
| ENeRF [12] | 32.56 | 0.953 | 0.050 | 23.21 | 0.530 | 0.291 | 32.62 | 0.968 | 0.051 | 5 |
| GPS-Gaussian+ | 33.72 | 0.961 | 0.039 | 23.07 | 0.643 | 0.238 | 33.74 | 0.971 | 0.041 | 25 |



Fig. 3: **Qualitative comparison on human-scene data.** Our method produces high-quality renderings with respect to others.

izable methods including Gaussian Splatting-based method MVSplat [13], implicit method ENeRF [12], image-based rendering method FloRen [62] and hybrid method IBR-Net [50]. In particular, it is difficult to train MVSplat on masked human-only data for its probabilistic modeling and FloRen relies on human prior, so we compare MVSplat only on human-scene data and FloRen only on human-only data. Following our setting, all baselines are trained on the same training set from scratch and take two source views as input for synthesizing the targeted novel view. The preliminary work, GPS-Gaussian, and FloRen use ground truth depths of synthetic human-only data for supervision. We further prepare the comparison with the original 3D-GS [14] for static human-only data and with 4D-GS [15] for sequential human-scene data which are optimized on all input views using the default strategies in the released code.

### 5.3.1 Results on Human-Scene Data

We compare state-of-the-art methods on 3 real captured human-scene datasets. In Table 1, our approach achieves superior or competitive results at the fastest speed with respect to other methods. In particular, our approach makes a great improvement on metric LPIPS which reveals better global rendering quality. We notice that camera parameters are not perfectly calibrated in the ENeRF-outdoor [12] dataset. Although the bad calibration has a tough impact on the back-propagation of rendering loss, our method can still synthesize fine-grained novel view images with more detailed appearances in Fig. 3. Due to the lack of consistent geometry prior, ENeRF and IBRNet can easily make floating artifacts illustrated in Fig. 3. MVSplat and ENeRF rely on multiplane sweeping to infer geometry from

TABLE 2: **Quantitative comparison on human-only datasets.** All methods are evaluated on an RTX 3090 GPU to report the speed of synthesizing one novel view with two $1024 \times 1024$ source images. Our methods and FloRen [62] use TensorRT for fast inference. † 3D-GS [14] requires per-subject optimization, while the other methods perform feed-forward inferences. The <mark>best</mark>, the <mark>second best</mark> and the <mark>third best</mark> are highlighted with different colors. ✓ denotes training with depth supervision.

| Method | Dep.Sup. | THuman2.0 [86] | | | Twindom [95] | | | Human-Only Real Data | | | FPS |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | |
| 3D-GS [14]† | | 24.18 | 0.821 | 0.144 | 22.77 | 0.785 | 0.153 | 22.97 | 0.839 | 0.125 | / |
| FloRen [62] | ✓ | 23.26 | 0.812 | 0.184 | 22.96 | 0.838 | 0.165 | 22.80 | 0.872 | 0.136 | 15 |
| IBRNet [50] | | 23.38 | 0.836 | 0.212 | 22.92 | 0.803 | 0.238 | 22.63 | 0.852 | 0.177 | 0.25 |
| ENeRF [12] | | 24.10 | 0.869 | 0.126 | 23.64 | 0.847 | 0.134 | 23.26 | 0.893 | 0.118 | 5 |
| GPS-Gaussian | ✓ | 25.57 | 0.898 | 0.112 | 24.79 | 0.880 | 0.125 | 24.64 | 0.917 | 0.088 | 25 |
| GPS-Gaussian+ | | 24.72 | 0.894 | 0.129 | 24.23 | 0.871 | 0.141 | 23.45 | 0.904 | 0.106 | 24 |



Fig. 4: **Qualitative comparison on human-only data.** Our method produces more detailed human appearances and can recover more reasonable geometry.

sparse views, such representation can hardly handle thin structure object, *e.g.* knife in Fig. 3. Although 4D-GS accelerates the optimization process from the original 3D-GS by decomposing spatial-temporal deformation into multi-resolution planes, such decomposition produces blurry re-sults under fast movements, see Fig. 1 and Fig. 3. Thanks to determinant stereo-matching and geometry regularization, our generated geometry is consistent from the two source views so that our rendering results are more decent with fewer floating artifacts.

TABLE 3: **Quantitative ablation study on synthetic human-only data.** We report PSNR, SSIM and LPIPS metrics for evaluating the rendering quality, while the end-point-error (EPE) and the ratio of pixel error in 1 pix level for measuring depth accuracy. ✓ denotes training with depth supervision.

| Model | Dep. Sup. | Rendering | | | Depth | |
|---|---|---|---|---|---|---|
| | | PSNR↑ | SSIM↑ | LPIPS↓ | EPE↓ | 1 pix↑ |
| GPS-Gaussian | ✓ | **25.05** | **0.886** | 0.121 | **1.494** | **65.94** |
| *w/o* Joint Train. | ✓ | 23.97 | 0.862 | **0.115** | 1.587 | 63.71 |
| *w/o* Depth Enc. | ✓ | 23.84 | 0.858 | 0.204 | 1.496 | 65.87 |
| GPS-Gaussian | | 24.22 | 0.874 | 0.145 | 4.066 | 33.38 |
| GPS-Gaussian+ | | 24.41 | 0.878 | 0.137 | 3.133 | 36.21 |

TABLE 4: **Quantitative ablation study on our captured human-scene data.** We report PSNR, SSIM and LPIPS metrics for evaluating the rendering quality.

| Model | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| GPS-Gaussian | 31.84 | 0.959 | 0.063 |
| GPS-Gaussian+ | **33.74** | **0.971** | **0.041** |
| *w/o* Geometry Reg. | 32.04 | 0.961 | 0.061 |
| *w/o* Epipolar Att. | 32.00 | 0.960 | 0.059 |
| *w/o* Depth Res. | 33.31 | 0.969 | 0.043 |

### 5.3.2 Results on Human-Only Data

We illustrate comparisons on two synthetic datasets and our collected real-world data in Table 2. Our method outperforms all baselines on all metrics and achieves a much faster rendering speed. Once occlusion occurs, some target regions under the novel view are invisible in one or both of the source views. ENeRF and IBRNet are not able to render reasonable results due to depth ambiguity. The unreliable geometric proxy in these cases also makes FloRen produce blurry outputs even if it employs the depth and flow refining networks. In our method, the efficient stereo-matching strategy and the geometry regularization help to alleviate the adverse effects caused by occlusion. In addition, it takes several minutes for 3D-GS parameter optimization for a single frame and produces noisy renderings of novel views, see Fig. 4, from such sparse views. Also, we demonstrate the effectiveness of our method on thin structures, such as the hockey stick and robe in Fig. 4.

## 5.4 Ablation Studies

In this part, we evaluate the effectiveness of our proposed components in GPS-Gaussian and GPS-Gaussian+ through ablation studies. The efficacy of joint training and depth encoder proposed in GPS-Gaussian are validated on aforementioned synthetic human-only data. As depth information is accessible in synthetic data, we further evaluate depth (identical to disparity) estimation, other than rendering metrics, with the end-point-error (EPE) and the ratio of pixel error in 1 pix level, following [29]. Furthermore, the additional components in GPS-Gaussian+ are evaluated on our captured human-scene data because our data includes single and multiple characters in different scenes. We also conduct a comparison between GPS-Gaussian and GPS-Gaussian+ on both datasets in order to illustrate the effi-



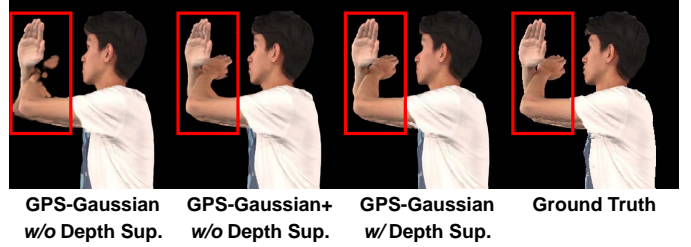| GPS-Gaussian | GPS-Gaussian+ | GPS-Gaussian | Ground Truth |
|---|---|---|---|
| *w/o* Depth Sup. | *w/o* Depth Sup. | *w/* Depth Sup. | |

Fig. 5: **Qualitative ablation study on GPS-Gaussian/GPS-Gaussian+ with different supervision settings.** We show the effectiveness of the integration in GPS-Gaussian+ when neglecting depth supervision.



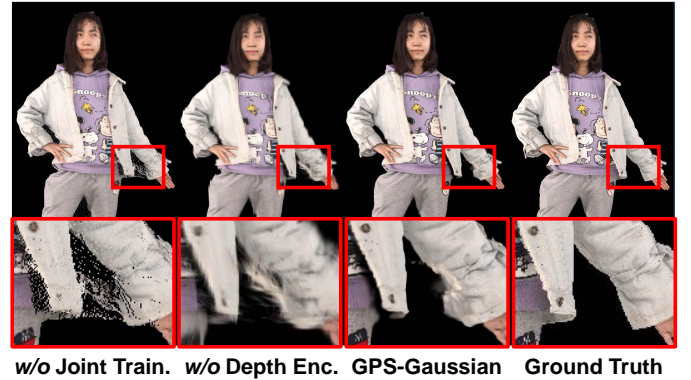| *w/o* Joint Train. | *w/o* Depth Enc. | GPS-Gaussian | Ground Truth |
|---|---|---|---|

Fig. 6: **Qualitative ablation study on designed components of GPS-Gaussian.** We show the effectiveness of the joint training and the depth encoder in the full pipeline. The proposed designs make the rendering results more visually appealing with fewer artifacts and less blurry.

ciency of adaptive integration under the setting of training without depth supervision.

### 5.4.1 Effects of Depth Supervision

A consistent geometry prior is essential to rendering reasonable images in novel views, especially for explicit Gaussian Splatting. Although the rendering metrics in Table 3 are not dramatically degraded for GPS-Gaussian when neglecting depth supervision, the geometry metrics are not competitive with respect to the models using depth supervision. In addition, a wrong geometry prior could produce unreasonable rendering results in occluded regions, see Fig. 5. Therefore, how to exploit accurate geometry proxy in the case of the lack of depth supervision is the key to our extensions.

### 5.4.2 Effects of Joint Training Mechanism

For GPS-Gaussian trained on synthetic human-only data, we compare jointly training both depth estimation and rendering modules with separately training them. We design a model substituting the differentiable Gaussian rendering with point cloud rendering at a fixed radius. Since the point cloud renderer is no longer differentiable, the rendering quality is merely based on the accuracy of depth estimation while the rendering loss could not conversely promote the depth estimator. The rendering results in Fig. 6 witness floating artifacts due to the depth ambiguity in the margin area of the source views where the depth value changes drastically. In Table 3, joint training makes a more robust depth estimator with a 5% improvement in EPE.
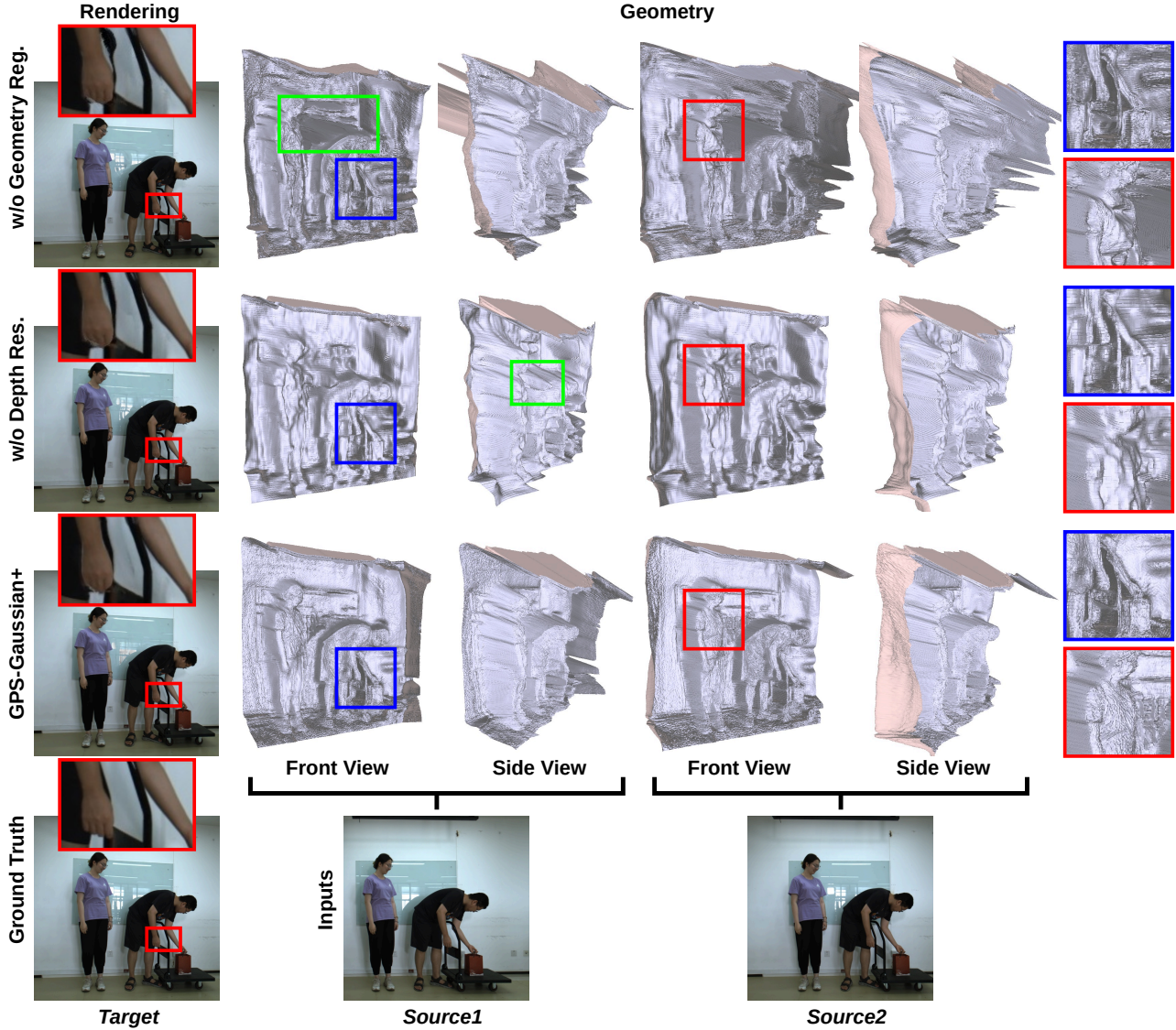
Fig. 7: **Qualitative ablation study on designed components in GPS-Gaussian+ for geometry.** We show the effectiveness of the geometry regularization and the depth residual in the full pipeline for geometry reconstruction.

### 5.4.3 Effects of Depth Encoder

We claim that merely using image features is insufficient for predicting Gaussian parameters. Herein, we ablate the depth encoder from our model, thus the Gaussian parameter decoder only takes as input the image features to predict $\mathcal{M}_r, \mathcal{M}_s, \mathcal{M}_\alpha$ simultaneously. As shown in Fig. 6, the ablated model fails to recover the details of human appearance, leading to blurry rendering results. The scale of Gaussian points is impacted by comprehensive factors including depth, texture and surface roughness, see Sec. 5.6. The absence of spatial awareness degrades the regression of scaling map $\mathcal{M}_s$, which deteriorates the visual perception reflected on LPIPS, even with a comparable depth estimation accuracy, as shown in Table 3.

### 5.4.4 Effects of Geometry Regularization

In GPS-Gaussian+ trained on real captured data without depth supervision, geometry regularization is designed to preserve geometry consistency between Gaussian points of the two source views. Due to the lack of depth supervision, marginal regions in Fig. 8 and regions with view-dependent reflection in Fig. 7 are hardly reconstructed when missing our proposed geometry regularization. Such geometric constraints can boost the unsupervised depth learning of the two source views to reach a mutual optimum. A good geometry prior also improves the rendering results, as reported in Table 4.

### 5.4.5 Effects of Depth Residual

Without depth supervision, GPS-Gaussian generates scarcely reasonable geometry, as shown in the fourth row of Table 3. Although GPS-Gaussian+ integrates the aforementioned adaptions, the geometry of the second row in Fig. 7 is still not acceptable. This problem is caused by two reasons. First, a downsampling operator is used in the stereo-matching module to minimize time cost. Second, the differentiability of point position in Gaussian Splatting is not satisfactory enough. By using our proposed depth residual map, our model recovers more details and corrects geometry artifacts, see the third row of Fig. 7.

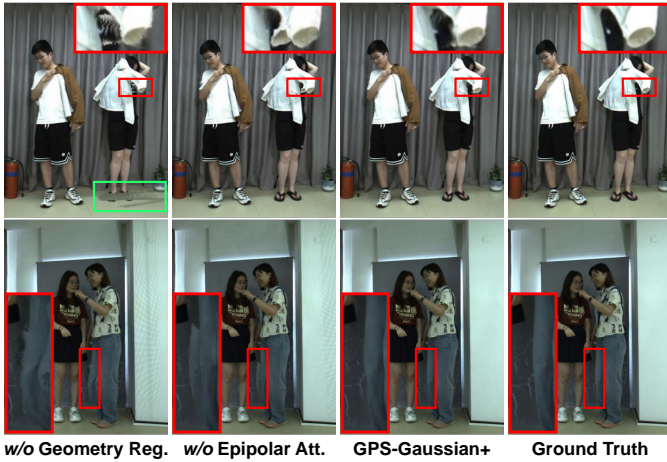**w/o Geometry Reg.**  **w/o Epipolar Att.**  **GPS-Gaussian+**  **Ground Truth**

Fig. 8: **Qualitative ablation study on designed components in GPS-Gaussian+ for rendering.** We show the effectiveness of the geometry regularization and the epipolar attention in the full pipeline.

### 5.4.6  Effects of Epipolar Attention

Compared with GPS-Gaussian, GPS-Gaussian+ incorporates epipolar attention into the feature extraction module to achieve a solid stereo-matching result with only rendering loss. Epipolar attention allows the encoder to exchange useful information between source views so that we can build a compact cost volume for stereo-matching. Even if the disparity (identical to depth) is not accessible during training, our proposed model corrects floating artifacts caused by wrong matching in Fig. 8. Since we apply such attention mechanism only along epipolar line, the time cost is on par with GPS-Gaussian, see Table 5.

In general, our adaptive integration is designed to compensate for the absence of depth supervision. In Table 3, GPS-Gaussian+ without depth supervision achieves competitive rendering results against GPS-Gaussian with depth supervision and produces reasonable geometry results. Moreover, the adaptive integration works better on full-scene images than masked human-only images. When the background is concerned, GPS-Gaussian+ can largely improve all rendering metrics with respect to GPS-Gaussian, as shown in Table 4.

### 5.5  Visualization of Opacity Maps

We discover that the joint regression with Gaussian parameters eliminates the outliers by predicting an extremely low opacity for the Gaussian points centered at these positions. The visualization of opacity maps is shown in Fig. 9. Since the depth prediction works on low resolution and upsampled to full image resolution, the drastically changed depth in the margin areas causes ambiguous predictions (*e.g.* the front and rear placed legs and the crossed arms in Fig. 9). These ambiguities lead to rendering noise on novel views when using a point cloud rendering technique. Thanks to the learned opacity map, the low opacity values make the outliers invisible in novel view rendering results, as shown in Fig. 9 (e).
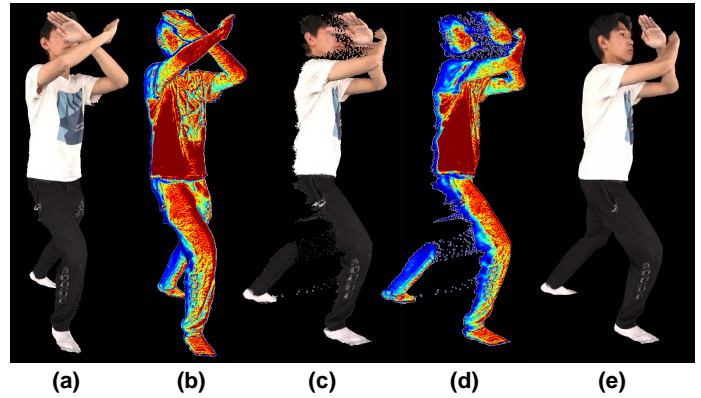


Fig. 9: **Visualization of opacity maps.** (a) One of the source view images. (b) The predicted opacity map related to (a). (c)/(d) The directly projected color/opacity map at novel viewpoint. (e) Novel view rendering results. A cold color in (b) and (d) represents an opacity value near 0, while a hot color near 1. The low opacity values predicted for the outliers make them invisible.



**(a) Input Image**  **(b) Depth**  **(c) Scaling Map**  **(d) Gaussian Points**
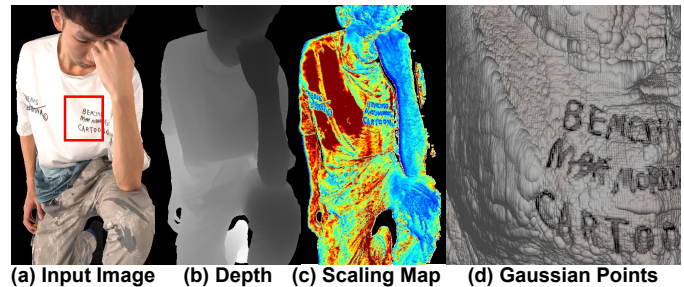
Fig. 10: **Visualization of scaling map and the shape of Gaussian points.** (a) One of the source view images. (b) The depth of (a). (c) The scaling map shown in heat map, where a hotter color represents a larger value. (d) The zoom-in Gaussian points of the boxed area in (a). The depth and scaling map are normalized.

### 5.6  Visualization of Scaling Maps

The visualization of the scaling map (mean of three axes) in Fig. 10 (c) indicates that the Gaussian points with lower depth roughly have smaller scales than the distant ones. However, the scaling property is also impacted by comprehensive factors. For example, as shown in Fig. 10 (c) and (d), fine-grained textures or high-frequency geometries lead to small-scaled Gaussians.

### 5.7  Robustness to Random Camera Views

We evaluate the robustness of our method to the randomly placed source-view cameras in the first row of Fig. 11. The model trained under a uniformly placed 8-camera setup in Sec. 5 shows a strong generalization capability to random camera setup with a pitch in range of $[-20°, +20°]$ and yaw in range of $[-25°, +25°]$ for human-only data. In Fig. 11 (f) and (h), our method achieves reasonable renderings of novel views with a pitch angle of about $\pm10°$ for human-scene data, even without any supervision of views with pitch angles during training.

**(a) Source View 1**  **(b) Source View 2**  **(c) Output**  **(d) Ground Truth**

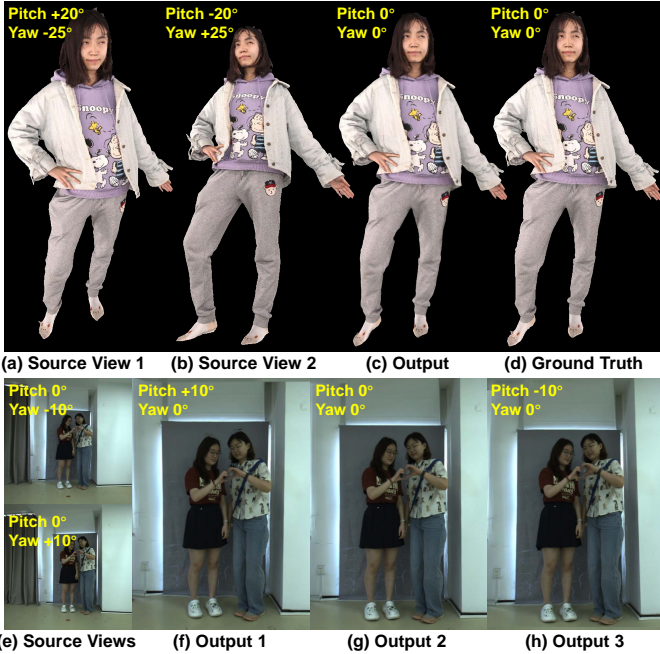**(e) Source Views**  **(f) Output 1**  **(g) Output 2**  **(h) Output 3**

Fig. 11: **Results on random camera views.** In the first row, (a) and (b) are the source view images with an extreme pitch and yaw, (c) is the novel view rendering result and (d) is the novel view ground truth. In the second row, we show (e) source views and target-view renderings with pitch angles of (f) $+10°$, (g) $0°$ and (h) $-10°$.

## 5.8 Robustness to Unseen Scenes

We further evaluate the generalization ability of our approach on unseen scene data in Fig. 12. We use the aforementioned model (Sec. 5.1) trained on our captured data under three backgrounds without any fine-tuning. Even if the background in Fig. 12 is totally unseen during training, our method is able to generate reasonable renderings.

## 5.9 Comparison on Run-Time

We conduct all experiments of our method and other baseline methods on the same machine with an RTX 3090 GPU of 24GB memory, except memory-consuming MVSplat [13]. Even if we prepare another machine with a V100 GPU of 32GB memory, MVSplat can only be fed with input images of $512 \times 512$ resolution during training. In Table 5, the overall run-time can be generally divided into two parts: one correlating to the source views and the other concerning the desired novel view. The source view correlated computation in FloRen [62] refers to coarse geometry initialization while the key components, the depth and flow refinement networks, operate on novel viewpoints. IBRNet [50] uses transformers to aggregate multi-view cues at each sampling point aggregated to the novel view image plane, which is time-consuming. ENeRF [12] constructs two cascade cost volumes on the target viewpoint, then predicts the target view depth followed by a depth-guided sampling for volume rendering. Once the target viewpoint changes, these methods need to recompute the novel view correlated modules. However, the computation on source views dominates the run-time of GPS-Gaussian+, which includes binocular depth estimation and Gaussian parameter map regression.

TABLE 5: **Run-time comparison.** We report the run-time correlated to the source views and each novel view on an RTX 3090 GPU. Input resolution is $512 \times 512$ for MVSplat, while all other methods take two $1024 \times 1024$ source images as input. Our methods can render multiple novel views concurrently in real-time.

| Methods | Source View Processing (ms) | Novel View Rendering (ms/view) |
|---|---|---|
| FloRen [62] | 14 | 11 |
| IBRNet [50] | 5 | 4000 |
| ENeRF [12] | 11 | 125 |
| MVSplat [13] | 120 | 1.5 |
| GPS-Gaussian | 27 | 1.9 |
| GPS-Gaussian+ | 30 | 1.9 |



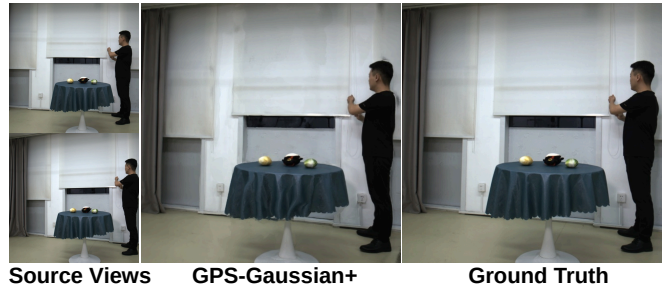**Source Views**  **GPS-Gaussian+**  **Ground Truth**

Fig. 12: **Result on unseen scenes.** In the case of unseen background during training, our method achieves reasonable rendering without any finetuning.

Similar to our method, MVSplat [13] spends the majority of run-time on the multi-view stereo process of source views, which includes multiple CNN and transformer operations. As reported in Table 5, it takes only 1.9 ms to render the 3D Gaussians to the desired novel view of human-scene data for GPS-Gaussian+, while this can be reduced to 0.8 ms when rendering human-only data with fewer Gaussian points. This allows us to render multiple novel views simultaneously, which caters to a wider range of applications such as holographic displays. Suppose that $n = 10$ novel views are required concurrently, it takes our method $T = T_{src} + n \times T_{novel} = 49ms$ to synthesize, while $135ms$ for MVSplat and $1261ms$ for ENeRF. In a real-world capture system, we should also consider I/O process and human matting, thus the frame rate is slightly degraded in Table 1 and Table 2.

## 6 LIVE-DEMO SYSTEMS

We prepare a machine equipped with an RTX 3090 GPU to run our algorithm and build two capture systems to shoot live demo. For human-scene data, our capture system consists of ten cameras positioned on a 1.6-meter beam, a piece of illumination equipment and two synchronizers, as shown in Fig. 13(a). In Fig. 13(b), we position all cameras in a circle of a 2-meter radius to capture human-only data. For live demos, capture and rendering processes are run on the same machine. We only use the cameras with red circles, 4 cameras in Fig. 13(a) and 6 cameras in Fig. 13(b), as input source views. Our method enables real-time high-quality rendering, even for challenging human-scene, human-object and multi-human interactions.
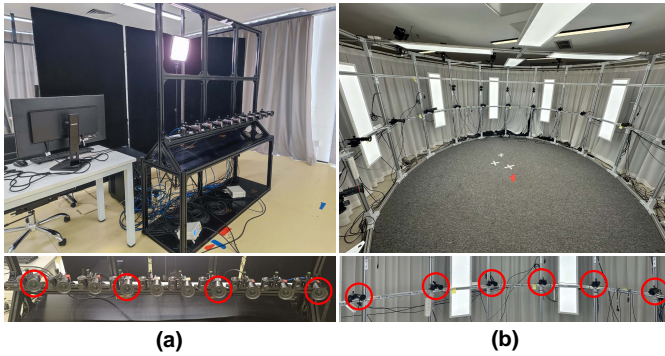
Fig. 13: **Capture systems.** We show (a) forward-facing camera rig capturing human-scene data and (b) human-centered camera stage for human-only data.

## 7 DISCUSSION

**Conclusion.** In this paper, we present GPS-Gaussian+, a feed-forward rendering method for both human-only data and human-scene data. By leveraging stereo-matching and pixel-wise Gaussian parameter map regression, our method takes a significant step towards a real-time photo-realistic human-centered free-viewpoint video system from sparse views. When lacking depth supervision during training, a regularization term and depth residual module are designed to ensure geometry consistency and high-frequency details. An adaptive integration, epipolar attention, is proposed in GPS-Gaussian+ to improve stereo-matching accuracy with only rendering supervision. We demonstrate that our GPS-Gaussian+ notably improves both quantitative and qualitative results compared with baseline methods and extends original GPS-Gaussian from human-only synthesis to more scalable and general scenarios of human-centered scenes.

**Limitations.** We notice some ghost artifacts on the white wall or the light yellow ground in the supplementary video. This is mainly because less textured regions could increase the difficulty of stereo-matching. Capturing more data covering more complex backgrounds to expand the diversity of the training scenes is a general solution to this issue. To achieve this, a portable system composed of mobile phones (*e.g.* Mobile-Stage dataset [69]) could break through the limitations of the fixed in-door capture system. Another feasible solution is using the massive monocular videos of static scenes captured by moving cameras (*e.g.* RealEstate10k [98]) to pre-train the network. However, since our method requires accurate camera calibration and strict synchronization, additional effort is required when making it practical to leverage the aforementioned data.

## REFERENCES

[1] Jason Lawrence, Danb Goldman, Supreeth Achar, Gregory Major Blascovich, Joseph G Desloge, Tommy Fortes, Eric M Gomez, Sascha Häberling, Hugues Hoppe, Andy Huibers, et al. Project starline: A high-fidelity telepresence system. *ACM TOG*, 40(6):1–16, 2021. 1

[2] 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 telepresence system with low-budget and high-authenticity using sparse rgb cameras. In *SIGGRAPH*, pages 1–12, 2024. 1

[3] Shenchang Eric Chen and Lance Williams. View interpolation for image synthesis. In *SIGGRAPH*, pages 279–288, 1993. 1

[4] Byong Mok Oh, Max Chen, Julie Dorsey, and Frédo Durand. Image-based modeling and photo editing. In *SIGGRAPH*, pages 433–442, 2001. 1

[5] Bennett Wilburn, Neel Joshi, Vaibhav Vaish, Eino-Ville Talvala, Emilio Antunez, Adam Barth, Andrew Adams, Mark Horowitz, and Marc Levoy. High performance imaging using large camera arrays. *ACM TOG*, 24(3):765–776, 2005. 1

[6] 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*, pages 405–421, 2020. 1, 2, 3

[7] 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*, pages 9054–9063, 2021. 1, 2

[8] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*, 41(4):1–15, 2022. 1, 3

[9] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, pages 5501–5510, 2022. 1, 3

[10] Ruizhi Shao, Zerong Zheng, Hanzhang Tu, Boning Liu, Hongwen Zhang, and Yebin Liu. Tensor4d: Efficient neural 4d decomposition for high-fidelity dynamic reconstruction and rendering. In *CVPR*, pages 16632–16642, 2023. 1, 2, 3

[11] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *CVPR*, pages 4578–4587, 2021. 1, 2

[12] Haotong Lin, Sida Peng, Zhen Xu, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Efficient neural radiance fields for interactive free-viewpoint video. In *SIGGRAPH Asia*, pages 1–9, 2022. 1, 3, 6, 7, 8, 12

[13] Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. In *ECCV*, 2024. 1, 2, 3, 7, 12

[14] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM TOG*, 42(4):1–14, 2023. 1, 2, 3, 5, 7, 8

[15] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. In *CVPR*, pages 20310–20320, 2024. 1, 3, 7

[16] Marc Levoy and Turner Whitted. The use of points as a display primitive. 1985. 1

[17] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Surface splatting. In *SIGGRAPH*, pages 371–378, 2001. 1

[18] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *CVPR*, pages 7467–7477, 2020. 1, 3

[19] Christoph Lassner and Michael Zollhofer. Pulsar: Efficient sphere-based neural rendering. In *CVPR*, pages 1440–1449, 2021. 1, 3

[20] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *ECCV*, pages 696–712, 2020. 1, 3

[21] Ruslan Rakhimov, Andrei-Timotei Ardelean, Victor Lempitsky, and Evgeny Burnaev. Npbg++: Accelerating neural point-based graphics. In *CVPR*, pages 15969–15979, 2022. 1, 3

[22] Georgios Kopanas, Thomas Leimkühler, Gilles Rainer, Clément Jambon, and George Drettakis. Neural point catacaustics for novel-view synthesis of reflections. *ACM TOG*, 41(6):1–15, 2022. 1

[23] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In *3DV*, 2024. 2, 3

[24] Stanislaw Szymanowicz, Chrisitian Rupprecht, and Andrea Vedaldi. Splatter image: Ultra-fast single-view 3d reconstruction. In *CVPR*, pages 10208–10217, 2024. 2

[25] David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In *CVPR*, pages 19457–19467, 2024. 2, 3

[26] Tianqi Liu, Guangcong Wang, Shoukang Hu, Liao Shen, Xinyi Ye, Yuhang Zang, Zhiguo Cao, Wei Li, and Ziwei Liu. Fast

generalizable gaussian splatting reconstruction from multi-view stereo. In *ECCV*, 2024. 2, 3

[27] Jiankun Li, Peisen Wang, Pengfei Xiong, Tao Cai, Ziwei Yan, Lei Yang, Jiangyu Liu, Haoqiang Fan, and Shuaicheng Liu. Practical stereo matching via cascaded recurrent network with adaptive correlation. In *CVPR*, pages 16263–16272, 2022. 2

[28] Philippe Weinzaepfel, Thomas Lucas, Vincent Leroy, Yohann Cabon, Vaibhav Arora, Romain Brégier, Gabriela Csurka, Leonid Antsfeld, Boris Chidlovskii, and Jérôme Revaud. Croco v2: Improved cross-view completion pre-training for stereo matching and optical flow. In *ICCV*, pages 17969–17980, 2023. 2

[29] Lahav Lipson, Zachary Teed, and Jia Deng. Raft-stereo: Multilevel recurrent field transforms for stereo matching. In *3DV*, pages 218–227, 2021. 2, 5, 6, 9

[30] 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 *CVPR*, pages 19680–19690, 2024. 2, 3

[31] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*, pages 4460–4470, 2019. 2

[32] 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 *ICCV*, pages 2304–2314, 2019. 2, 3

[33] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *CVPR*, pages 84–93, 2020. 2

[34] Yang Hong, Juyong Zhang, Boyi Jiang, Yudong Guo, Ligang Liu, and Hujun Bao. Stereopifu: Depth aware clothed human digitization via stereo vision. In *CVPR*, pages 535–545, 2021. 2

[35] 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*, pages 1–9, 2022. 2

[36] Fuqiang Zhao, Wei Yang, Jiakai Zhang, Pei Lin, Yingliang Zhang, Jingyi Yu, and Lan Xu. Humannerf: Efficiently generated human radiance field from sparse inputs. In *CVPR*, pages 7743–7753, 2022. 2

[37] 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 *CVPR*, pages 16210–16220, 2022. 2

[38] Chen Guo, Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Vid2avatar: 3d avatar reconstruction from videos in the wild via self-supervised scene decomposition. In *CVPR*, pages 12858–12868, 2023. 2

[39] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, pages 165–174, 2019. 2

[40] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 34:27171–27183, 2021. 2

[41] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. In *ICCV*, pages 3295–3306, 2023. 2

[42] Boyao Zhou, Di Meng, Jean-Sébastien Franco, and Edmond Boyer. Human body shape completion with implicit shape and flow learning. In *CVPR*, pages 12901–12911, 2023. 2

[43] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *IEEE TPAMI*, 44(6):3170–3184, 2021. 2

[44] Yuliang Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J Black. Icon: Implicit clothed humans obtained from normals. In *CVPR*, pages 13286–13296, 2022. 2

[45] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *ECCV*, pages 523–540. Springer, 2020. 2

[46] Julian Chibane, Gerard Pons-Moll, et al. Neural unsigned distance fields for implicit function learning. *NeurIPS*, 33:21638–21652, 2020. 2

[47] Ruizhi Shao, Hongwen Zhang, He Zhang, Mingjia Chen, Yan-Pei Cao, Tao Yu, and Yebin Liu. Doublefield: Bridging the neural surface and radiance fields for high-fidelity human reconstruction and rendering. In *CVPR*, pages 15872–15882, 2022. 2

[48] Jianchuan Chen, Wentao Yi, Liqian Ma, Xu Jia, and Huchuan Lu. Gm-nerf: Learning generalizable model-based neural radiance fields from multi-view images. In *CVPR*, pages 20648–20658, 2023. 2, 3

[49] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. Neuman: Neural human radiance field from a single video. In *ECCV*, pages 402–418, 2022. 2

[50] 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 *CVPR*, pages 4690–4699, 2021. 3, 7, 8, 12

[51] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *ICCV*, pages 14124–14133, 2021. 3

[52] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *ICCV*, pages 5752–5761, 2021. 3

[53] Ruilong Li, Hang Gao, Matthew Tancik, and Angjoo Kanazawa. Nerfacc: Efficient sampling accelerates nerfs. In *ICCV*, pages 18537–18546, 2023. 3

[54] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *ECCV*, pages 623–640, 2020. 3

[55] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *CVPR*, pages 12216–12225, 2021. 3

[56] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM TOG*, 38(4):1–12, 2019. 3

[57] Moustafa Meshry, Dan B Goldman, Sameh Khamis, Hugues Hoppe, Rohit Pandey, Noah Snavely, and Ricardo Martin-Brualla. Neural rerendering in the wild. In *CVPR*, pages 6878–6887, 2019. 3

[58] Francesco Pittaluga, Sanjeev J Koppal, Sing Bing Kang, and Sudipta N Sinha. Revealing scenes by inverting structure from motion reconstructions. In *CVPR*, pages 145–154, 2019. 3

[59] Ricardo Martin-Brualla, Rohit Pandey, Shuoran Yang, Pavel Pidlypenskyi, Jonathan Taylor, Julien Valentin, Sameh Khamis, Philip Davidson, Anastasia Tkach, Peter Lincoln, et al. Lookingood: enhancing performance capture with real-time neural re-rendering. *ACM TOG*, 37(6):1–14, 2018. 3

[60] Phong Nguyen-Ha, Nikolaos Sarafianos, Christoph Lassner, Janne Heikkilä, and Tony Tung. Free-viewpoint rgb-d human performance capture and rendering. In *ECCV*, pages 473–491, 2022. 3

[61] Ang Cao, Chris Rockwell, and Justin Johnson. Fwd: Real-time novel view synthesis with forward warping and depth. In *CVPR*, pages 15713–15724, 2022. 3

[62] Ruizhi Shao, Liliang Chen, Zerong Zheng, Hongwen Zhang, Yuxiang Zhang, Han Huang, Yandong Guo, and Yebin Liu. Floren: Real-time high-quality human performance rendering via appearance flow using sparse rgb cameras. In *SIGGRAPH Asia*, pages 1–10, 2022. 3, 7, 8, 12

[63] Yebin Liu, Qionghai Dai, and Wenli Xu. A point-cloud-based multiview stereo algorithm for free-viewpoint video. *IEEE TVCG*, 16(3):407–418, 2009. 3

[64] Boyao Zhou, Jean-Sébastien Franco, Federica Bogo, Bugra Tekin, and Edmond Boyer. Reconstructing human body mesh from point clouds by adversarial gp network. In *ACCV*, 2020. 3

[65] Qianli Ma, Jinlong Yang, Siyu Tang, and Michael J Black. The power of points for modeling humans in clothing. In *ICCV*, pages 10974–10984, 2021. 3

[66] Siyou Lin, Hongwen Zhang, Zerong Zheng, Ruizhi Shao, and Yebin Liu. Learning implicit templates for point-based clothed human modeling. In *ECCV*, pages 210–228, 2022. 3

[67] Yufeng Zheng, Wang Yifan, Gordon Wetzstein, Michael J Black, and Otmar Hilliges. Pointavatar: Deformable point-based head avatars from videos. In *CVPR*, pages 21057–21067, 2023. 3

[68] Hongwen Zhang, Siyou Lin, Ruizhi Shao, Yuxiang Zhang, Zerong Zheng, Han Huang, Yandong Guo, and Yebin Liu. Closet: Modeling clothed humans on continuous surface with explicit template decomposition. In *CVPR*, pages 501–511, 2023. 3

[69] Zhen Xu, Sida Peng, Haotong Lin, Guangzhao He, Jiaming Sun, Yujun Shen, Hujun Bao, and Xiaowei Zhou. 4k4d: Real-time 4d

view synthesis at 4k resolution. In *CVPR*, pages 20029–20040, 2024. 3, 13

[70] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. In *CVPR*, pages 5438–5448, 2022. 3

[71] Shih-Yang Su, Timur Bagautdinov, and Helge Rhodin. Npc: Neural point characters from video. In *ICCV*, pages 14795–14805, 2023. 3

[72] Liangxiao Hu, Hongwen Zhang, Yuxiang Zhang, Boyao Zhou, Boning Liu, Shengping Zhang, and Liqiang Nie. Gaussianavatar: Towards realistic human avatar modeling from a single video via animatable 3d gaussians. In *CVPR*, pages 634–644, 2024. 3

[73] Ruizhi Shao, Jingxiang Sun, Cheng Peng, Zerong Zheng, Boyao Zhou, Hongwen Zhang, and Yebin Liu. Control4d: Efficient 4d portrait editing with text. In *CVPR*, pages 4556–4567, 2024. 3

[74] Yuelang Xu, Benwang Chen, Zhe Li, Hongwen Zhang, Lizhen Wang, Zerong Zheng, and Yebin Liu. Gaussian head avatar: Ultra high-fidelity head avatar via dynamic gaussians. In *CVPR*, pages 1931–1941, 2024. 3

[75] Zhe Li, Zerong Zheng, Lizhen Wang, and Yebin Liu. Animatable gaussians: Learning pose-dependent gaussian maps for high-fidelity human avatar modeling. In *CVPR*, pages 19711–19722, 2024. 3

[76] Xian Liu, Xiaohang Zhan, Jiaxiang Tang, Ying Shan, Gang Zeng, Dahua Lin, Xihui Liu, and Ziwei Liu. Humangaussian: Text-driven 3d human generation with gaussian splatting. In *CVPR*, pages 6646–6657, 2024. 3

[77] Jiakai Sun, Han Jiao, Guangyuan Li, Zhanjie Zhang, Lei Zhao, and Wei Xing. 3dgstream: On-the-fly training of 3d gaussians for efficient streaming of photo-realistic free-viewpoint videos. In *CVPR*, pages 20675–20685, 2024. 3

[78] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, and Zhaoyang Lv. Neural 3d video synthesis from multi-view video. In *CVPR*, pages 5521–5531, 2022. 3, 6, 7

[79] Benjamin Attal, Jia-Bin Huang, Christian Richardt, Michael Zollhoefer, Johannes Kopf, Matthew O'Toole, and Changil Kim. Hyperreel: High-fidelity 6-dof video with ray-conditioned sampling. In *CVPR*, pages 16610–16620, 2023. 3

[80] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *CVPR*, pages 130–141, 2023. 3

[81] Zhan Li, Zhang Chen, Zhong Li, and Yi Xu. Spacetime gaussian feature splatting for real-time dynamic view synthesis. In *CVPR*, pages 8508–8520, 2024. 3

[82] Yuheng Jiang, Zhehao Shen, Penghao Wang, Zhuo Su, Yu Hong, Yingliang Zhang, Jingyi Yu, and Lan Xu. Hifi4g: High-fidelity human performance rendering via compact gaussian splatting. In *CVPR*, pages 19734–19745, 2024. 3

[83] Xin Suo, Yuheng Jiang, Pei Lin, Yingliang Zhang, Minye Wu, Kaiwen Guo, and Lan Xu. Neuralhumanfvv: Real-time neural volumetric human performance rendering using rgb cameras. In *CVPR*, pages 6226–6237, 2021. 3

[84] Youngjoong Kwon, Baole Fang, Yixing Lu, Haoye Dong, Cheng Zhang, Francisco Vicente Carrasco, Albert Mosella-Montoro, Jianjin Xu, Shingo Takagi, Daeil Kim, et al. Generalizable human gaussians for sparse view synthesis. In *ECCV*, 2024. 3

[85] Kaiwen Guo, Feng Xu, Tao Yu, Xiaoyang Liu, Qionghai Dai, and Yebin Liu. Real-time geometry, albedo, and motion reconstruction using a single rgb-d camera. *ACM TOG*, 36(3):1–13, 2017. 3

[86] 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 *CVPR*, pages 5746–5756, 2021. 3, 6, 8

[87] Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural human performer: Learning generalizable radiance fields for human performance rendering. *NeurIPS*, 34:24741–24752, 2021. 3

[88] Xiao Pan, Zongxin Yang, Jianxin Ma, Chang Zhou, and Yi Yang. Transhuman: A transformer-based human representation for generalizable neural human rendering. In *ICCV*, pages 3544–3555, 2023. 3

[89] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM TOG*, 34(6):1–16, 2015. 3

[90] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa splatting. *IEEE TVCG*, 8(3):223–238, 2002. 3

[91] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM TOG*, 38(6):1–14, 2019. 3

[92] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 4

[93] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, pages 402–419, 2020. 5

[94] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 13(4):600–612, 2004. 6

[95] Twindom, 2020. https://web.twindom.com. 6, 8

[96] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pages 586–595, 2018. 6

[97] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2018. 6

[98] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: learning view synthesis using multiplane images. *ACM TOG*, 37(4):1–12, 2018. 13