

# SmileSplat: Generalizable Gaussian Splats for Unconstrained Sparse Images

Yanyan Li<sup>1</sup>Yixin Fang<sup>2</sup>Federico Tombari<sup>3,4</sup>Gim Hee Lee<sup>1</sup><sup>1</sup>National University of Singapore<sup>2</sup>Zhejiang university<sup>3</sup>Technical University of Munich<sup>4</sup>Google

## Abstract

*Sparse Multi-view Images can be Learned to predict explicit radiance fields via Generalizable Gaussian Splatting approaches, which can achieve wider application prospects in real-life when ground-truth camera parameters are not required as inputs. In this paper, a novel generalizable Gaussian Splatting method, **SmileSplat**, is proposed to reconstruct pixel-aligned Gaussian surfels for diverse scenarios only requiring unconstrained sparse multi-view images. First, Gaussian surfels are predicted based on the multi-head Gaussian regression decoder, which can be represented with less degree-of-freedom but have better multi-view consistency. Furthermore, the normal vectors of Gaussian surfel are enhanced based on high-quality of normal priors. Second, the Gaussians and camera parameters (both extrinsic and intrinsic) are optimized to obtain high-quality Gaussian radiance fields for novel view synthesis tasks based on the proposed Bundle-Adjusting Gaussian Splatting module. Extensive experiments on novel view rendering and depth map prediction tasks are conducted on public datasets, demonstrating that the proposed method achieves state-of-the-art performance in various 3D vision tasks. More information can be found on our project page (<https://yanyan-li.github.io/project/gs/smilesplat>).*

## 1. Introduction

Novel view rendering technology has become widely used in both industrial production and entertainment applications, enabling the generation of photo-realistic views of unknown environments based on estimated implicit or explicit radiance fields [20, 26]. To construct high-quality radiance fields, Structure-from-Motion (SfM) [30, 32] or Simultaneous Localization and Mapping (SLAM) [22, 48] systems are typically employed to generate 3D point clouds, along with camera intrinsics and extrinsics, from a set of

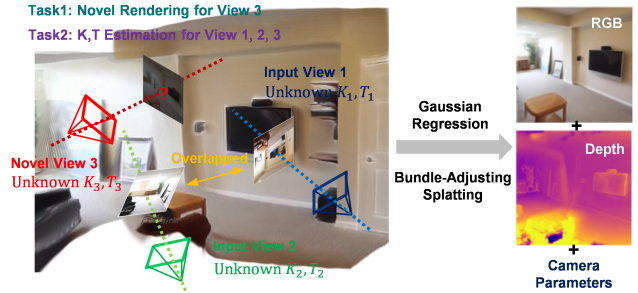


Figure 1. An example result of SmileSplat. It aims to render novel views and estimate camera parameters (intrinsic  $\mathbf{K}$  and extrinsic  $\mathbf{T}$ ) by using sparse views.

dense-view RGB images. However, when working with sparse-view images (as few as two), accurately estimating 3D point primitives and camera parameters using traditional systems becomes challenging. This issue is particularly pronounced in low-textured scenes or with fast camera motion, where the limited number of images fails to provide sufficient constraints for conventional optimization techniques.

Camera calibration [49], tracking [28], dense mapping [16, 19], and rendering [29] are classic problems in the fields of computer vision and computer graphics, typically addressed by traditional solutions based on multi-view geometry strategies [19, 28, 32]. Compared to traditional methods, Neural Radiance Fields (NeRFs) [26] and 3D Gaussian Splatting [20] have achieved photo-realistic novel view rendering performance. Given 2D posed images, NeRFs [26] address the problem by training neural networks to encode implicit 3D representations, enabling significant advances in photo-realistic novel view rendering applications. To improve efficiency in training and inference, 3D Gaussian Splatting methods use point-based representations instead of implicit neural parameterizations via differential optimization modules. Even though Gaussian Splatting SLAM approaches [25, 43, 51] have capabilities in tracking cameras in unknown scenes, they have limita-

tions in two main areas. On the one hand, a large number of views are required for training. On the other hand, camera intrinsic parameters and depth information are necessary to achieve robust rendering performance when dealing with monocular RGB inputs. When the input images collected in challenging scenes are sparse, it is extremely difficult to calibrate camera parameters and obtain the 3D structure of the unknown scenes. Therefore, the traditional Gaussian Splatting [20, 23] and neural radiance fields [2, 26] are doubly challenged by low-quality inputs and fewer constraints.

To improve the performance of novel view rendering in sparse-shot tasks, generalizable Gaussian splatting approaches [6, 45, 47] are explored in this area. These approaches can be classified into two types: **CamPara-Required** and **CamPara-Free** methods, depending on whether camera parameters are required as part of the inputs. For CamPara-Required methods, given ground truth camera poses and intrinsics, 3D Gaussians predicted by networks [4, 6] are rendered to novel viewpoints. To generate accurate Gaussian primitives, neural multi-head decoders first predict depth (point clouds), covariance, and opacity values. A forward-stream rendering module is then utilized to optimize the initial Gaussian parameters based on these predictions. To further eliminate the reliance on camera parameters, point clouds, rather than depth maps, are predicted by networks such as Dust3R [38] and Mast3R [21] in the canonical system. These point clouds are used to initialize 3D Gaussians [11], which are subsequently used to obtain camera poses. Then, Gaussian parameters will be optimized through the Gaussian Splatting module [20]. Since these two modules are separate, additional iterations are required to train the 3D Gaussians for such scenarios. Compared to CamPara-Required methods, CamPara-Free ones are more convenient for applications since they do not require an initialization step. However, these pioneering CamPara-Free approaches tend to directly integrate camera calibration based on predicted point clouds and the Gaussian Splatting module together, which limits the further high-quality achievement of generalizable Gaussian radiance fields.

In this paper, we propose a new generalizable Gaussian Splatting architecture, as shown in Figure 2, aimed at achieving high-fidelity novel view rendering performance for unconstrained sparse-view images. First, Gaussian surfels are predicted by a forward-stream neural network that utilizes a standard transformer encoder (Siamese ViT encoders [8] and cross-attention embedding blocks [39]) to detect geometric priors from images. These deep priors, along with the images, are then fed into the proposed Multi-head Gaussian Regression Decoder to predict pixel-aligned, generalizable 3D Gaussian surfel parameters in the canonical coordinate frame. Furthermore, we estimate an intrinsic matrix by considering geometric and photometric constraints based on the initial 3D Gaussians. To the best of our

knowledge, our method is the first to render images without requiring predefined intrinsic parameters. Next, the relative extrinsic matrix between images is predicted based on the estimated Gaussian surfels. To improve the consistency of the predicted Gaussian surfels, we propose a Gaussian Splatting Bundle Adjustment method to further refine the Gaussian parameters, intrinsics, and extrinsics. This refinement is based on photometric and geometric constraints, allowing us to establish scaled Gaussian Radiance Fields for unconstrained sparse images. The contributions of this paper are summarized as:

1. We present a versatile generalizable Gaussian Splatting architecture for un-calibrated and un-posed sparse-view images.
2. A camera parameter optimization module based on Gaussian Splatting is analyzed to achieve accurate motion estimation of sparse images.
3. A bundle-adjusting Gaussian Splatting algorithm is implemented to produce high-quality and scaled Gaussian radiance fields.
4. Extensive evaluations on variety of scenes demonstrate competitive novel view rendering performance. The code and generated dataset will be released to the community.

## 2. Related Work

**Radiance Fields for a Single Scene.** Radiance Fields have garnered significant attention in the field of 3D vision due to their ability to generate novel views of objects or scenes from arbitrary viewpoints. Neural Radiance Fields (NeRF) [26] is among the pioneering and most prominent methods for efficiently rendering high-quality novel views by representing 3D scenes implicitly via Multi-layer Perceptrons (MLPs). However, NeRF is hindered by slow training and inference speeds. Subsequent research [2, 3] has primarily focused on either enhancing rendering quality or improving computational efficiency. Recent advances have introduced explicit volume structures such as multi-resolution voxel grids [12, 37] or hash functions [27] to improve performance. Despite these improvements, per-pixel ray marching remains a bottleneck for rendering speed, while which is a critical issue in SLAM applications requiring real-time map interactions. In contrast, 3D Gaussian Splatting (3DGS) [20] employs anisotropic 3D Gaussians to represent radiance fields, combined with differentiable splatting for rendering. This method excels at quickly reconstructing complex real-world scenes with high detail, particularly capturing high-frequency details effectively. By iterating over rasterized primitives instead of marching along rays, 3DGS utilizes the natural sparsity of



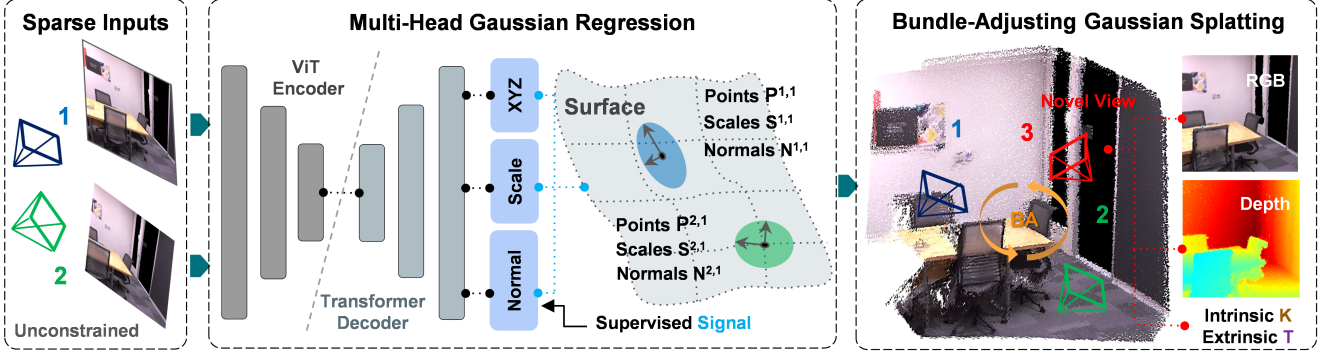


Figure 2. Architecture of SmileSplat. With sparse but overlapping views as input, the system consists of two main modules, Multi-Head Gaussian Regression and Bundle-Adjusting Gaussian Splatting, for achieving scaled Gaussian radiance fields.

3D scenes, offering a balance of high-fidelity representation and efficient rendering. Various studies have applied 3D Gaussians and differentiable rendering to static scene capture [7, 43], with recent works demonstrating superior results in dynamic scene capture [41, 44].

**Generalizable Gaussian Splatting.** Similar to multi-view stereo tasks, when two images with limited overlaps are fed into 3D Gaussian Splatting systems, several new challenges arise since traditional Gaussian Splatting systems [20] require a large number of images as inputs. SparseGS [42] addresses these challenges by using diffusion networks to remove outliers and different types of rendered depth maps to detect floaters. Instead of relying on initial point clouds generated from COLMAP [13] or SLAM systems [22], MVSplat [6] and COLMAP-free GS [13] estimate depth maps for both target and source RGB images, where these depth values are then transferred to point clouds, and the mean vectors of Gaussian ellipsoids are initialized based on these point clouds. Benefiting from the point clouds of Dust3R [38], the Gaussian Splatting process in InstantSplat [11] becomes more efficient compared to randomly generated point clouds used in MonoGS [25]. Based on the camera parameters estimated in Dust3r, InstantSplat optimizes 3D Gaussians and camera poses in the rendering process. However, the separation of modules means that there is a missed opportunity to further accelerate the convergence of Gaussian Splatting. Compared to InstantSplat which estimates camera parameters based on predicted point clouds from Dust3R, pixelSplat [4] predicts 3D Gaussians directly, which are defined in the Gaussian Splatting rasterization using ground truth camera parameters. Recently, pose-free novel view rendering methods [45] have estimated and optimized camera poses using the PnP algorithm [15] and Gaussian Splatting SLAM method [25], while incorporating intrinsic parameters into a deep feature token to predict scenes at a reasonable scale. Unlike these

approaches, our method predicts 3D Gaussians for uncalibrated and unposed images, estimating the unknown camera parameters based on the predicted Gaussians. All these estimated parameters and representations are then optimized via the bundle-adjusting Gaussian Splatting module to establish high-fidelity radiance fields.

### 3. Methodology

#### 3.1. Gaussian Surfel Prediction

Instead of using 3D Gaussians to represent 3D scenes, as used in PixelSplat [4] and MVSplat [6], the network in this paper predicts Gaussian surfels [43] for two main reasons. First, 3D Gaussians are not ideal for sparse-view training tasks due to limited constraints on Gaussian consistency [17]. Second, Gaussian surfels involve fewer parameters compared to the rotation matrix, and the more important thing is that they are easier for training compared to the covariance matrix, as several robust large language models [1] can be leveraged to train this head.

**Gaussian Surfel Representation.** Each 3D Gaussian surfel  $\mathcal{G}_i$  is parameterized by 12 values, consisting of three for its RGB color  $\mathbf{c}_i$ , three for its center position  $\boldsymbol{\mu}_i = [x_i \ y_i \ z_i]^T$ , two for its scale vector  $\mathbf{s}_i = [s_i^x \ s_i^y]^T$ , three for its normal vector  $\mathbf{n}_i = [n_i^x \ n_i^y \ n_i^z]^T$  representing its direction, and one for its opacity  $o_i \in [0, 1]$ . These parameters are expressed as:

$$\mathcal{G}_i = [\mathbf{c}_i \ \boldsymbol{\mu}_i \ \mathbf{s}_i \ \mathbf{n}_i \ o_i]. \quad (1)$$

The covariance matrix  $\boldsymbol{\Sigma}_i$  for the Gaussian surfel can then be defined as:

$$\boldsymbol{\Sigma}_i = \mathbf{R}_i \text{Diag}[(s_i^x)^2 \ (s_i^y)^2 \ 0] (\mathbf{R}_i)^T \quad (2)$$

where  $\mathbf{R}_i = [\mathbf{n}_1 \ \mathbf{n}_2 \ \mathbf{n}_i] \in SO(3)$  is an orthonormal matrix. Here,  $\mathbf{n}_1$  lies on the great circle of  $\mathbf{n}_i$ , and  $\mathbf{n}_2$  is obtained via the cross product  $\mathbf{n}_1 \times \mathbf{n}_i$ .  $Diag[\cdot]$  is a diagonal matrix.

**Standard Transformer Backbone.** Inspired by the architecture of CroCo [39] and Dust3D [38], we use the model that is pre-trained based on million-level images, where two ViT [8] encoders sharing weights are fed by multi-view stereo images  $I_1$  and  $I_2$ . Continually, a generic transformer network [40] equipped with self-attention and cross-attention strategies is used to deal with tokens from the Siamese ViT backbone.

**Multi-head Gaussian Regression Decoder.** In the Gaussian regression head, we predict 3D Gaussians in the coordinate system of image  $I_1$ , using four separate blocks to estimate the position  $\mathbf{P}$ , surface normal  $\mathbf{n}$ , opacity  $o$ , and scaling  $\mathbf{s}$  vectors. To reduce the parameter scale, we directly use the color of each pixel in the Gaussian surfel representation.

For the position prediction, we use an offset prediction block, leveraging a pretrained pointmap head from Dust3R [38]. The position is expressed as  $\mathbf{P} = \hat{\mathbf{P}} + \delta\mathbf{P}$ , where  $\hat{\mathbf{P}}$  is the predicted position and  $\delta\mathbf{P}$  is the offset, predicted by a 3-layer MLP network fed with the concatenation of pointmap features and the encoder embeddings.

For the surface normal prediction, we employ a U-Net architecture that predicts surface normals for each pixel based on the input image and deep embeddings. To achieve faster convergence of this block, we use the predicted surface normals  $\hat{\mathbf{n}}$  from a pre-trained surface normal network [1] to supervise the normal prediction head. The supervision is applied using the following loss function:

$$\mathcal{L}_n = |\mathbf{n} - \hat{\mathbf{n}}|_1 + |1 - \mathbf{n} \cdot \hat{\mathbf{n}}|_1. \quad (3)$$

To predict opacity, we note that points with low offsets are more likely to be on the accurate surface and should thus have higher opacity. Therefore, we feed the deep features from the backbone and features from the position offset block to predict the opacity of each 3D Gaussian.

For the shape prediction of each Gaussian, we already have its surface normal, then we estimate the scaling component of the Gaussian. Benefitting of the using representation, the scale vector is defined as  $\mathbf{s} = [s_1 \ s_2 \ 0]$ , which is predicted by a 2-layer MLP architecture with the inputs of pointmaps and deep embeddings. Therefore, the covariance matrix of the Gaussian surfel can be represented based on the Equation 2.

### 3.2. Camera Parameter Optimization Based on a Single View

**Intrinsic Parameter Estimation.** In methods like COLMAP [31, 33] and Dust3R [38], intrinsic parameters

are typically estimated by solving optimization problems that relate pixels to point clouds in camera coordinates. In contrast, we propose a more robust approach based on 3D Gaussian Splatting. Since the predicted Gaussians are pixel-aligned with the input image and located in the coordinate frame of the first image, they can be treated as being in the camera coordinate system. Thus, the unknown parameters are the intrinsic matrix, which can be optimized via the formulation:

$$\mathbf{K}^* = \min_{\mathbf{K}} \sum_{u=0}^W \sum_{v=0}^H \mathcal{L}_{c_1} \quad (4)$$

where  $\mathcal{L}_{c_1} = |I_1^{u,v}(\mathbf{K}, \mathcal{G}) - \bar{I}_1^{u,v}|_1$ ,  $\bar{I}_1^{u,v}$  represents the observed pixel value at position  $(u, v)$  in the input image, and  $I(\mathbf{K}, \mathcal{G})$  renders the image based on the intrinsic parameters  $\mathbf{K}$  and the initial Gaussian surfels  $\mathcal{G}$ .  $W$  and  $H$  represent the image width and height. Since the Gaussian primitives are defined in canonical coordinates, the intrinsic parameters can be optimized by minimizing Equation 4 while fixing the Gaussians.

The differential process for optimizing the intrinsic parameters is implemented within the Gaussian Splatting module, leveraging the CUDA library for efficient computation. The terms in Equation 4 are differentiable with respect to the camera intrinsic matrix  $\mathbf{K}$  via the chain rule:

$$\begin{aligned} \frac{\partial \mu_I}{\partial \mathbf{K}} &= \frac{\partial \mu_I}{\partial \mu_C} \frac{\partial \mu_C}{\partial \mathcal{P}} \frac{\partial \mathcal{P}}{\partial \mathbf{K}} \\ \frac{\partial \Sigma_I}{\partial \mathbf{K}} &= \frac{\partial \Sigma_I}{\partial \mathbf{J}} \frac{\partial \mathbf{J}}{\partial \mathbf{K}} \end{aligned} \quad (5)$$

where  $\mathbf{J}$  is the Jacobian of the affine approximation of the projective transformation  $\mathcal{P}$  [20],  $\mu_I$  and  $\Sigma_I$  are the 2D Gaussian on the image plane, while  $\mu_C$  and  $\Sigma_C$  are 3D Gaussian in the canonical coordinates. The detailed derivation of the gradients with respect to the camera intrinsic matrix  $\mathbf{K}$  is provided in the supplementary materials (see Section A.2 and A.3).

To assist with the convergence of the optimization process, we initialize the principal point at the center of the image plane, with  $c_x = 0.5W$  and  $c_y = 0.5H$ , respectively. The initial focal length is set to 1.2 times the image dimensions (width and height) to provide a reasonable starting estimate for the camera’s optical properties. It is important to note that the offset of the principal point is relatively small compared to the amount of adjustment needed for the camera’s focal length. Therefore, we use separate learning rates,  $l_c$  for the principal point parameters and  $l_f$  for the focal length, to ensure efficient and stable optimization.

**Extrinsic Parameter Estimation.** The camera extrinsic matrix  $\mathbf{T} = \begin{bmatrix} \mathbf{W} & \mathbf{t} \\ 0 & 1 \end{bmatrix} \in SE(3)$  consists of the rotation

matrix  $\mathbf{W} \in SO(3)$  and the translation vector  $\mathbf{t} \in \mathbb{R}^3$ , which describes the rigid transformation between the two camera views,  $I_1$  and  $I_2$ . In traditional camera pose estimation methods [28, 33], the common approach is to re-project 3D point clouds from world coordinates back into the image plane for camera tracking. However, in this section, we introduce the strategy for camera pose estimation based on the predicted Gaussian surfels.

To optimize the relative camera pose between the first and second views, we fix the intrinsic parameters and the Gaussians predicted by the first image, and iteratively refine the second view’s pose in canonical coordinates. The optimization is performed by minimizing the loss function defined as:

$$\mathbf{T}^* = \min_{\mathbf{T}} \sum_{u=0}^W \sum_{v=0}^H \mathcal{L}_{c_2}$$

where  $\mathcal{L}_{c_2} = |I_2^{u,v}(\mathbf{K}, \mathcal{G}, \mathbf{T}) - \bar{I}_2^{u,v}|_1$ , and  $\bar{I}_2$  is the observed image. The optimization of the camera extrinsics proceeds by calculating the gradients of the error term with respect to the camera pose parameters. The derivatives of the 2D Gaussian position  $\mu_I$  and covariance matrix  $\Sigma_I$  with respect to the camera extrinsics  $\mathbf{T}$  are computed as follows:

$$\begin{aligned} \frac{\partial \mu_I}{\partial \mathbf{T}} &= \frac{\partial \mu_I}{\partial \mu_C} \frac{\partial \mu_C}{\partial \mathbf{T}} \\ \frac{\partial \Sigma_I}{\partial \mathbf{T}} &= \frac{\partial \Sigma_I}{\partial \mathbf{J}} \frac{\partial \mathbf{J}}{\partial \mu_C} \frac{\partial \mu_C}{\partial \mathbf{T}} + \frac{\partial \Sigma_I}{\partial \mathbf{W}} \frac{\partial \mathbf{W}}{\partial \mathbf{T}} \end{aligned} \quad (6)$$

where  $\mathbf{T}$  is the camera extrinsic matrix, and the gradient update for the parameters  $\mathbf{T}$  occurs on the manifold  $\mathfrak{se}(3)$  via *Lie* algebra representations [25].

### 3.3. Bundle-Adjusting Gaussian Splatting

After obtaining the initial camera intrinsic parameters and relative camera poses (see Section 3.2), the Gaussian surfels are rasterized into corresponding depth maps  $D_1$  and  $D_2$  using the alpha-blending algorithm. The depth map  $D_2$  from the second viewpoint is transformed to the first viewpoint using the following warping operation:

$$D_{1, \text{warp}} = \Pi(D_2, \mathbf{K}, \mathbf{T}) \quad (7)$$

where  $\Pi(\cdot)$  represents the warp operation that projects the depth map  $D_2$  from the second camera view to the first camera view using the estimated intrinsic matrix  $\mathbf{K}$  and the relative camera pose  $\mathbf{T}$ .

Then, a geometric constraint is established to enforce better multi-view consistency of the Gaussian radiance fields in terms of geometry. This constraint is formulated as follows:

$$\mathcal{L}_{geo} = |D_{1, \text{warp}} - D_1|. \quad (8)$$

And the photometric error between rendered and observed images is defined as

$$\mathcal{L}_{pho} = (1.0 - \lambda_{sm})\mathcal{L}_{c_i} + \lambda_{sm}|1.0 - \mathcal{L}_{ssim_i}| \quad (9)$$

where  $i \in [1, 2]$ , and  $\mathcal{L}_{ssim_i} = \text{ssim}(I_i^{u,v} - \bar{I}_i^{u,v})$  is to compute structural similarity index measure distance between two images.

To jointly optimize the scaled radiance fields, including Gaussian surfels, as well as the camera intrinsic and extrinsic parameters, for the input sparse views, we define a comprehensive loss function to supervise the refinement process:

$$\mathcal{G}^*, \mathbf{K}^*, \mathbf{T}^* = \min_{\mathcal{G}, \mathbf{K}, \mathbf{T}} \sum_{u=0}^W \sum_{v=0}^H (\lambda_1 \mathcal{L}_{pho_1} + \lambda_2 \mathcal{L}_{pho_2} + \lambda_3 \mathcal{L}_{geo}) \quad (10)$$

where  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  are weighting parameters that balance the contributions of the respective loss terms. This joint optimization allows for the simultaneous refinement of the Gaussian surfels, camera intrinsics, and extrinsics, improving the overall performance of the sparse-view radiance fields reconstruction.

To clarify the refinement stage, we illustrate the strategy of this stage used in the experiments, as shown in Algorithm 1. The following procedure outlines the iterative process for intrinsics, Gaussian primitives, extrinsics, and joint optimization:

---

**Algorithm 1** Strategy for achieving scaled Gaussian radiance fields

---

**Require:**  $\mathcal{G}$ ,  $I^{c_1}$ , and  $I^{c_2}$

**Ensure:** Iteration  $n = 100$

```

1: for  $j = 0$  to 10 do
2:    $\mathbf{K}$  optimization
3: end for
4: for  $j = 10$  to 20 do
5:    $\mathcal{G}$  optimization
6: end for
7: for  $j = 20$  to 40 do
8:    $\mathbf{T}$  optimization
9:   if  $\mathbf{T} \rightarrow \mathbf{T}^*$  then
10:    Finish
11:   end if
12: end for
13: for  $j = 40$  to 100 do
14:    $\mathcal{G}, \mathbf{K}, \mathbf{T}$  optimization
15: end for
```

---

## 4. Experiments

In this section, we present both quantitative and qualitative evaluations comparing the proposed method with state-of-

the-art algorithms, focusing on performance in the novel view rendering task.

#### 4.1. Implementation Details and Baselines

SmileSplat consists of two main modules: Gaussian prediction and radiance field bundle adjustment. For the first module, which handles Gaussian surfel prediction, we utilize a neural network built on the PyTorch library. The encoder is based on a Vision Transformer (ViT) model with a patch size of 16, while the decoder is based on a ViT-base architecture. This module includes four separate heads for Gaussian parameter registration, designed to generate Gaussians for an image resolution of  $224 \times 224$ . The specific size and scale of each head are provided in the supplementary material. The second module, radiance field bundle adjustment, incorporates both intrinsic and extrinsic parameters. This module is implemented using the CUDA library to leverage GPU acceleration for efficient computations. For training, we use two GeForce RTX 3090 GPUs. During the inference stage, only a single GPU is used in our experiments. Throughout all experiments, we maintain a consistent learning rate of 0.0002 for Gaussian optimization. The weighting  $\lambda_{ssim}$  for the photometric loss is 0.6. Additionally, the weights  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  are set to 0.05, 0.05 and 0.01, respectively.

The state-of-the-art approaches are compared against with the proposed method in the novel view synthesis task, which are categorized as two types, where the first set of methods, including pixelNeRF [47], AttnRend [9], pixelSplat [5], and MVSplat [6], needs camera parameters in their training and testing process, while the second set of approaches, DUS3R [38], MAST3R [21], Splatt3R [35], and NoPoSplat [45], have capabilities to estimate intrinsic and extrinsic.

#### 4.2. Datasets and Metrics

Following the experimental setup of pioneering approaches [4, 6] in this domain, we evaluate the proposed method on two large-scale datasets. The first dataset, RealEstate10K (Re10K) [50], is collected from real-estate sequences on YouTube. Based on the training/testing split, the dataset contains 29,144 scenes for training and 7286 scenes for testing, respectively. The second dataset, ACID[24], focuses on nature scenes collected from the viewpoints of aerial drones. To evaluate the reconstruction capabilities of the proposed method, we further assess the trained model across additional datasets, including Replica [36] and ICL-NUIM [14].

For quantitative evaluation, we report the rendering performance using standard image quality metrics, including Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM), and Learned Perceptual Image Patch Similarity (LPIPS). Specifically, the first two metrics

(PSNR and SSIM) evaluate the color-wise similarity and structural similarity between the rendered and observed images. The third metric, LPIPS, compares the feature-level similarity between two images, using features extracted by a pre-trained neural network (e.g., VGG-Net [34]) rather than directly comparing the pixel values of the images. To analyze the relationship between rendering performance and the overlap of input images, the visual overlap  $\gamma$  between two input images is computed based on a dense feature matching method [10]. The overlap is then classified into three levels: Small ( $\gamma \leq 0.3\%$ ), Medium ( $0.3\% \leq \gamma \leq 0.55\%$ ), and Large ( $\gamma \geq 0.55\%$ ), following the method of [45].

#### 4.3. Novel View Rendering

As shown in Table 1, state-of-the-art methods, both CamPara-Free and CamPara-Required, are compared with the proposed approach, SmileSplat, in the novel view rendering task. For CamPara-required methods, Gaussian Splatting approaches, such as pixelSplat and MVSplat, demonstrate superior rendering quality compared to implicit representation methods like pixelNeRF. For CamPara-free methods, DUS3R [38] and MAST3R [21] predict pixel-aligned point clouds based on the input images, and other intrinsic and extrinsic parameters can be estimated and optimized using conventional multi-view geometry algorithms. Building upon the architectures of these methods, the approaches Splatt3R [35] and NoPoSplat [45] introduce additional heads to estimate the parameters of 3D Gaussian ellipsoids, which significantly improve rendering performance from 14.49 to 23.08 for inputs with small visual overlaps. Unlike the Re10K dataset, which consists of real-estate sequences, the ACID dataset focuses on natural scenarios. However, the trends observed in Table 1 and Figure 3 can also be seen in the supplementary.

#### 4.4. Cross-Dataset Generalization.

As shown in Table 3, both NoPoSplat and the proposed SmileSplat demonstrate superior robustness and accuracy compared to other CamPara-Free and CamPara-Required approaches. In this section, we continue to evaluate the zero-shot performance of NoPoSplat and SmileSplat by directly applying their models to the Replica sequences [36]. It is important to note that both models were trained on the Re10K dataset without any further fine-tuning on the Replica dataset.

Table 3 lists four sequences, including *Office 0*, *Office 1*, *Room 0*, and *Room 1*, which are used to compare the novel view rendering performance. Similar to the settings used in the Re10K testing, two images are fed into the systems for Gaussian radiance field reconstruction, and a novel view is used for testing. Based on the number of input images, the sequences are categorized into three classes: Small (7 im-



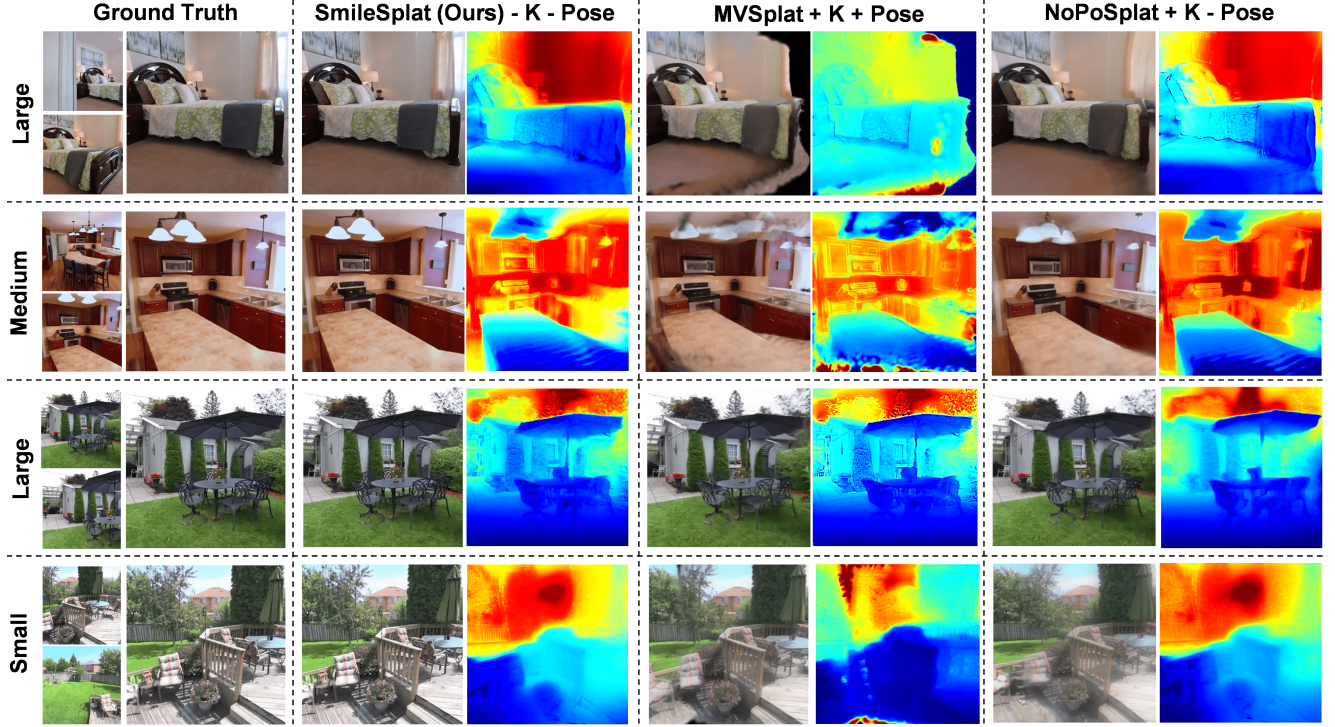


Figure 3. Comparison of novel view rendering and depth prediction results on Re10k [50] and ACID [24] datasets with different settings. The rendering results include the RGB in the left column and the depth in the right column.  $\pm K$  and  $\pm Pose$  denote whether intrinsics and extrinsics are free or not.

Method		pixelNeRF [47]	AttnRend [9]	pixelSplat [5]	MVSplat [6]	DUST3R [38]	MASt3R [21]	Splatt3R [35]	NoPoSplat [45]	SmileSplat
Setting	K	✓	✓	✓	✓	×	×	×	✓	×
	Pose	✓	✓	✓	✓	×	×	×	×	×
Small	PSNR ↑	18.41	19.15	20.26	20.35	14.10	13.53	14.35	22.51	<b>26.81</b>
	SSIM ↑	0.601	0.663	0.717	0.724	0.432	0.407	0.475	0.585	<b>0.856</b>
	LPIPS ↓	0.526	0.368	0.266	0.250	0.468	0.494	0.472	0.462	<b>0.071</b>
Medium	PSNR ↑	19.93	22.53	23.71	23.77	15.41	14.94	15.52	24.89	<b>26.60</b>
	SSIM ↑	0.632	0.763	0.809	0.812	0.451	0.436	0.502	0.616	<b>0.843</b>
	LPIPS ↓	0.480	0.269	0.181	0.173	0.432	0.451	0.425	0.160	<b>0.071</b>
Large	PSNR ↑	20.86	25.89	27.15	27.40	16.42	16.02	15.81	<b>27.41</b>	27.26
	SSIM ↑	0.639	0.845	0.879	0.884	0.453	0.444	0.483	<b>0.883</b>	0.850
	LPIPS ↓	0.458	0.186	0.122	0.116	0.402	0.418	0.421	0.119	<b>0.063</b>

Table 1. Quantitative comparisons of novel view rendering on the Re10K [50]. ✓ and × in the setting rows show whether corresponding camera parameters are required or not during the inference.

ages), Medium (12 images), and Large (20 images). Within each class, the proposed SmileSplat consistently robustness performance. Specifically, in the *Room 1* sequence, the PSNR result for NoPoSplat is 27.64, which is improved by 25% to 34.60 with our method. However, when directly comparing the rendering quality of NoPoSplat, as shown in Figure 4, the rendered images still exhibit high fidelity. To better understand the gap between the quantitative and qualitative results, we compute the photometric distance between the rendered images and the corresponding ground truth images, as shown in Figure 4. The results reveal that the alignment between the rendered images and the ground

truth images has notable issues, especially in the *Office* sequence. In textured regions, the photometric error between the rendered and ground truth images is significantly higher for NoPoSplat. Since the camera pose of the reference image is estimated by these methods themselves, the observed phenomenon suggests that the proposed SmileSplat method achieves better alignment and overall performance than NoPoSplat. More detailed camera pose experiments are provided in the supplementary materials.

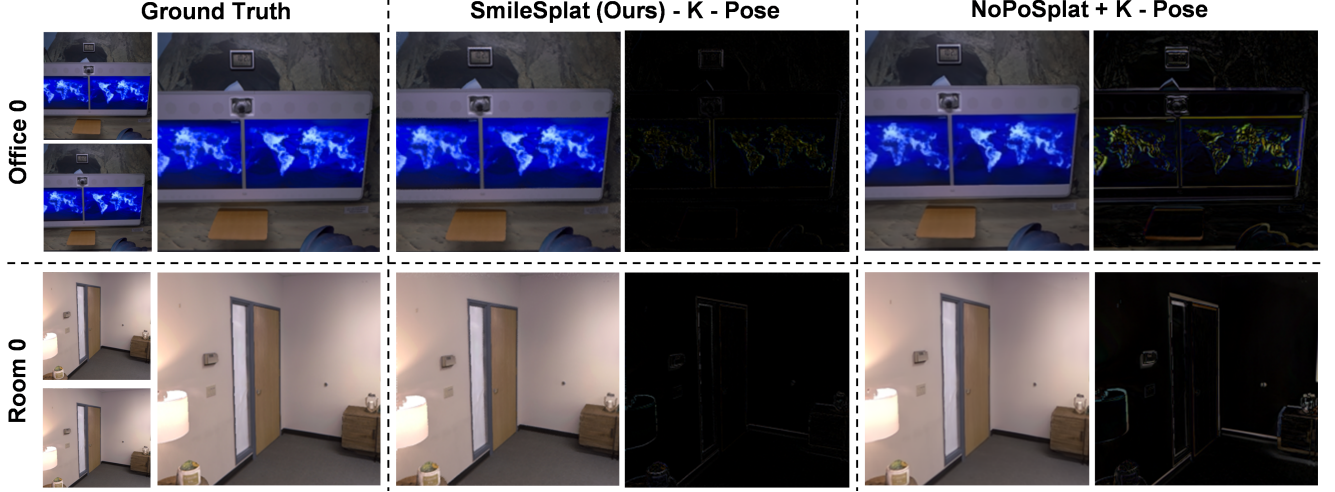


Figure 4. Comparisons of novel view rendering on Replica [36] dataset. The rendering results include the RGB in the left column and the difference with the ground truth in the right column.  $\pm K$  and  $\pm$  Pose denote whether intrinsics and extrinsics are free or not.

Method		View 1			View 2			View 3		
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
$\mathcal{P}_h + \text{Cam}_a$		23.98	0.769	0.227	15.80	0.486	0.355	18.63	0.582	0.296
$\mathcal{P}_h + \text{Cam}_b$		23.67	0.763	0.183	16.72	0.538	0.283	17.85	0.559	0.266
$\mathcal{G}_h + \text{Cam}_a + \mathcal{L}_{pho}: \mathcal{O}_G$		28.43	0.911	0.083	18.61	0.743	0.234	18.63	0.744	0.234
$\mathcal{G}_h + \text{Cam}_b + \mathcal{L}_{pho}: \mathcal{O}_G$		29.06	0.927	0.070	20.64	0.807	0.188	19.27	0.608	0.240

Initial	$\mathcal{P}_h$	Pointmap	point cloud prediction of Dust3R
	$\mathcal{G}_h$	$[\mu \text{ n o s}]$	Gaussian representation prediction of our method
	$\text{Cam}_a$	K, R, t	Intrinsics and extrinsics of Dust3R
	$\text{Cam}_b$	K, R, t	Intrinsics and extrinsics of ICP

Optimization	$\mathcal{L}_{pho}$	$\mathcal{L}_{\{pho\}}$	photometric constraint between RGB images
	$\mathcal{O}_G$	$[\mu \text{ n s o}]$	Gaussian surfel optimization
	$\mathcal{O}_K$	$f_x, f_y, c_z, c_y$	intrinsic optimization
	$\mathcal{O}_{pose}$	Lie algebra	camera pose optimization

Table 2. Ablation studies for testing different modules in the novel view rendering task based on the Re10K [50] dataset. For Gaussian prediction testing, all three views are novel views, while View 3 is the novel view when the first two views are used for reference Gaussian surfels. More settings are provided in Supplementary.

Method		NoPoSplat [45]				SmileSplat (Ours)			
		Small	Medium	Large	Ave.	Small	Medium	Large	Ave.
Office 0	PSNR	32.64	28.61	30.99	30.44	<b>35.62</b>	<b>33.82</b>	<b>31.65</b>	<b>33.71</b>
	SSIM	0.878	0.845	0.869	0.861	<b>0.953</b>	<b>0.942</b>	<b>0.923</b>	<b>0.940</b>
	LPIPS	0.166	0.164	0.172	0.167	<b>0.036</b>	<b>0.049</b>	<b>0.071</b>	<b>0.052</b>
Office 1	PSNR	32.65	28.59	30.95	30.42	<b>34.22</b>	<b>33.27</b>	<b>31.79</b>	<b>33.39</b>
	SSIM	0.877	0.844	0.869	0.861	<b>0.934</b>	<b>0.926</b>	<b>0.918</b>	<b>0.928</b>
	LPIPS	0.166	0.164	0.172	0.167	<b>0.056</b>	<b>0.062</b>	<b>0.077</b>	<b>0.063</b>
Room 0	PSNR	28.43	30.37	32.60	30.47	<b>31.11</b>	<b>31.76</b>	<b>30.89</b>	<b>31.33</b>
	SSIM	0.684	0.723	0.723	0.702	<b>0.894</b>	<b>0.913</b>	<b>0.895</b>	<b>0.902</b>
	LPIPS	0.157	0.150	0.150	0.155	<b>0.025</b>	<b>0.016</b>	<b>0.024</b>	<b>0.021</b>
Room 1	PSNR	27.64	27.73	29.50	28.29	<b>34.60</b>	<b>33.64</b>	<b>32.39</b>	<b>33.83</b>
	SSIM	0.836	0.134	0.729	0.780	<b>0.961</b>	<b>0.956</b>	<b>0.949</b>	<b>0.957</b>
	LPIPS	0.133	0.777	0.140	0.135	<b>0.015</b>	<b>0.019</b>	<b>0.025</b>	<b>0.018</b>

Table 3. Cross-Dataset generalization tests on the Replica [36] dataset.

#### 4.5. Ablation Study

In this section, we analyze the performance of different modules proposed in our method. As shown in Table 2,

various settings of modules are integrated and tested on the Re10K dataset benchmark. First, we render the predicted Gaussians at three different viewpoints using two different initial camera parameter estimation methods. Since the Gaussian parameters are predicted at the coordinates of View 1, camera pose errors result in degraded performance when rendering from the other two viewpoints. When we optimize the Gaussian parameters based on feedback from the photometric loss of View 1 and View 2, the rendering quality for these two views improves significantly. However, this optimization does not have a strong positive impact on the novel view rendering (View 3), as the camera pose of View 3 cannot be accurately estimated. This suggests that while feedback from a few views helps improve performance for those specific viewpoints, the lack of accurate pose estimation for novel views limits the overall benefit.

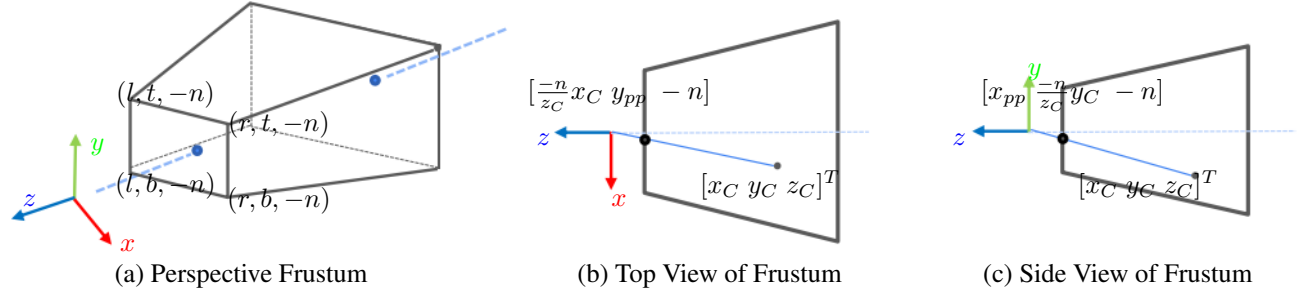


Figure 5. 3D Perspective frustum and its different views in 2D space.

## 5. Discussion and Conclusion

In this paper, we introduce a robust and generalizable Gaussian Splatting method for unconstrained sparse images. To improve multi-view consistency in 3D Gaussian models, we estimate the Gaussian surfel parameters using our multi-head Gaussian registration framework. To avoid reliance on ground truth camera intrinsics and extrinsics, these parameters are optimized directly in the canonical coordinate frame based on the bundle-adjusting Gaussian Splatting module, to achieve accurate scaled Gaussian radiance fields.

### A. Fundamental Theory in RGB and Depth Rendering

#### A.1. Novel View Rendering

During the Gaussian Splatting optimization process, the rendered RGB and depth maps used in Section 3.2, are computed based on the color  $\mathbf{c}_i$  (or depth) and opacity  $\alpha_i$  of the involved Gaussians  $\mathcal{G}_i$ :

$$\begin{cases} \mathcal{C}_p = \sum_{i \in N} \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \\ \mathcal{D}_p = \sum_{i \in N} z_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \end{cases} \quad (11)$$

where  $z_i$  is the distance to the position of the Gaussian surfel along the camera ray, and the  $\alpha_i$  can be computed based on the 2D covariance  $\Sigma_I \in \mathbb{R}^{2 \times 2}$  and the opacity value  $o_i$  following the 3DGS algorithm [20].

#### A.2. Projection Process in Gaussian Splatting

In perspective projection, the 3D position  $\mu_C$  (the mean vector of the Gaussian in camera coordinates) within a truncated pyramid frustum, as illustrated in Figure 5(a), can be mapped to normalized device coordinates (NDC) through the following two steps: (1) projecting  $\mu_C$  onto the projection plane, and (2) mapping the projection onto the normalized device coordinates.

Based on the top view, as shown in Figure 5(b), of the truncated pyramid frustum, the following relationship can be derived:

$$\begin{cases} x_{pp} = -\frac{n}{z_C} x_C \\ y_{pp} = -\frac{n}{z_C} y_C \end{cases} \quad (12)$$

where  $x_{pp}$  and  $y_{pp}$  represent the coordinates on the projection plane, and  $\mu_C = [x_C \ y_C \ z_C]^T$  denotes the 3D position in camera coordinates. In the truncated pyramid frustum, the range of the  $x$ -coordinate spans from  $l$  to  $r$ , while the ranges for the  $y$ - and  $z$ -coordinates are  $[b, t]$  and  $[-n, -f]$ , respectively.

When mapping the range  $[l \rightarrow r]$  of the projection plane in the  $x$ -direction to the normalized device coordinates (NDC) range  $[-1 \rightarrow 1]$ , the mapping function can be defined as a linear transformation:

$$x_{ndc} = ax_{pp} + b \quad (13)$$

where  $a = \frac{1-(-1)}{r-l}$  and  $b$  can be computed as  $b = 1 - \frac{2r}{r-l}$  when  $x_{pp}$  is set to  $r$ . Therefore, the mapping relationship between  $x_{pp}$  and  $x_{ndc}$  is given by:

$$x_{ndc} = \frac{2}{r-l} x_{pp} - \frac{r+l}{r-l}. \quad (14)$$

Similarly, the mapping relationship in the  $y$ -direction can be expressed as:

$$y_{ndc} = \frac{2}{t-b} y_{pp} - \frac{t+b}{t-b} \quad (15)$$

where  $t$  and  $b$  represent the upper and lower bounds of the projection plane in the  $y$ -direction, respectively, as shown in Figure 5.

By substituting  $[x_C \ y_C]$  from Equation 12 in place of  $[x_{pp} \ y_{pp}]$  in Equations 14 and 15, the mapping function from camera coordinates to NDC can be expressed as:

$$\begin{bmatrix} x_{ndc} \\ y_{ndc} \end{bmatrix} = \begin{bmatrix} -\frac{2}{r-l} \frac{n}{z_C} x_C - \frac{r+l}{r-l} \\ -\frac{2}{t-b} \frac{n}{z_C} y_C - \frac{t+b}{t-b} \end{bmatrix} \quad (16)$$

which can then be simplified as functions of

$$\begin{bmatrix} x_{ndc} \\ y_{ndc} \end{bmatrix} = \begin{bmatrix} -(\frac{2n}{r-l}x_C + \frac{r+l}{r-l}z_C)\frac{1}{z_C} \\ -(\frac{2n}{t-b}y_C + \frac{t+b}{t-b}z_C)\frac{1}{z_C} \end{bmatrix}. \quad (17)$$

By setting  $w_{clip} = -z_C$ , the Equation 17 can be rewritten as:

$$\begin{bmatrix} x_{ndc} \\ y_{ndc} \\ z_{ndc} \end{bmatrix} = \begin{bmatrix} x_{clip}/w_{clip} \\ y_{clip}/w_{clip} \\ z_{clip}/w_{clip} \end{bmatrix} \quad (18)$$

$$\begin{aligned} \begin{bmatrix} x_{clip} \\ y_{clip} \\ z_{clip} \\ w_{clip} \end{bmatrix} &= \begin{bmatrix} \frac{2n}{r-l}x_C + \frac{r+l}{r-l}z_C \\ \frac{2n}{t-b}y_C + \frac{t+b}{t-b}z_C \\ -\frac{(f+n)}{f-n}z_C + \frac{-2fn}{f-n} \\ -z_C \end{bmatrix} \\ &= \begin{bmatrix} \frac{2n}{r-l} & 0 & \frac{r+l}{r-l} & 0 \\ 0 & \frac{2n}{t-b} & \frac{t+b}{t-b} & 0 \\ 0 & 0 & -\frac{(f+n)}{f-n} & \frac{-2fn}{f-n} \\ 0 & 0 & -1 & 0 \end{bmatrix} \begin{bmatrix} x_C \\ y_C \\ z_C \\ 1 \end{bmatrix} \end{aligned} \quad (19)$$

which can be denoted in a more familiar form as:

$$\mu_{clip} = \mathcal{P}[\mu_C^T 1]^T \quad (20)$$

where  $\mu_{clip} = [x_{clip} \ y_{clip} \ z_{clip} \ w_{clip}]^T$  and  $\mathcal{P}$  is represented by

$$\begin{bmatrix} \frac{2n}{r-l} & 0 & \frac{r+l}{r-l} & 0 \\ 0 & \frac{2n}{t-b} & \frac{t+b}{t-b} & 0 \\ 0 & 0 & -\frac{(f+n)}{f-n} & \frac{-2fn}{f-n} \\ 0 & 0 & -1 & 0 \end{bmatrix}.$$

Furthermore, we explore the relationship between the intrinsic matrix and the projection matrix based on the camera's field of view (FoV). First, we assume that the viewing volume is symmetric, which allows us to derive the following:

$$\begin{cases} r+l=0 \\ t+b=0 \end{cases} \begin{cases} r = \tan(FoV_x/2) \cdot n \\ t = \tan(FoV_y/2) \cdot n \end{cases} \quad (21)$$

here  $\tan(FoV_x/2)$  and  $\tan(FoV_y/2)$  are  $\frac{W}{2f_x}$  and  $\frac{H}{2f_y}$ , respectively. Additionally,  $f_x$  and  $f_y$  are the focal length parameters of the intrinsic matrix  $\mathbf{K}$ , given by:  $\mathbf{K} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$ .

Therefore the Equation 19 can be represented. We express the transformation from camera coordinates to normalized clip space  $\mathcal{P}$  [46, 52],

$$\mathcal{P} = \begin{bmatrix} \frac{2f_x}{W} & 0 & 0 & 0 \\ 0 & \frac{2f_y}{H} & 0 & 0 \\ 0 & 0 & -\frac{(f+n)}{f-n} & \frac{-2fn}{f-n} \\ 0 & 0 & -1 & 0 \end{bmatrix} \quad (22)$$

here,  $W$  and  $H$  represent the width and height of the image.

After obtaining the positions in the NDC cube, the screen coordinates are computed by applying a viewport transformation from the normalized device coordinates, which maps them into rendering pixel coordinates:

$$\mu_I = \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} \frac{W}{2}(x_{ndc} + 1) + c_x \\ \frac{H}{2}(y_{ndc} + 1) + c_y \end{bmatrix}. \quad (23)$$

Therefore, the Equation 23 can be substituted as the mapping function  $\begin{bmatrix} u \\ v \end{bmatrix} = \phi(\mu_C)$ , and the function  $\phi(\mathbf{x})$  for a position  $\mathbf{x} = [x \ y \ z]^T$  can be denoted as:

$$\begin{aligned} \phi(\mathbf{x}) &= \begin{bmatrix} \frac{W}{2} & 0 & \frac{W}{2} + c_x \\ 0 & \frac{H}{2} & \frac{H}{2} + c_y \end{bmatrix} \begin{bmatrix} x_{ndc} \\ y_{ndc} \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} \frac{W}{2} & 0 & \frac{W}{2} + c_x \\ 0 & \frac{H}{2} & \frac{H}{2} + c_y \end{bmatrix} \begin{bmatrix} x_{clip}/w_{clip} \\ y_{ndc}/w_{clip} \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} \frac{W}{2} & 0 & \frac{W}{2} + c_x \\ 0 & \frac{H}{2} & \frac{H}{2} + c_y \end{bmatrix} \begin{bmatrix} \frac{2f_x}{W}x/z \\ \frac{2f_y}{H}y/z \\ 1 \end{bmatrix}. \end{aligned} \quad (24)$$

### A.3. Affine Approximation of Projective Matrix

As is well-known, Gaussians are closed [52] under affine mappings. However, the projection transformations discussed are not affine, and the perspective projection of a 3D Gaussian does not directly yield a 2D Gaussian [18, 23].

To address this, we establish a local affine approximation  $\phi_k(\mathbf{x})$  of the projective transformation, defined by the first two terms of the Taylor expansion of  $\phi$  at the point  $\mu$ :

$$\phi(\mathbf{x}) \approx \phi(\mu) + \mathbf{J}(\mathbf{x} - \mu) \quad (25)$$

where  $\mu$  represents the center of the Gaussian in camera coordinates. The Jacobian  $\mathbf{J}$  is given by the partial derivatives of  $\phi$  evaluated at  $\mu$ ,  $\mathbf{J} = \frac{\partial \phi}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mu}$ .

Therefore, to approximate the projection of the covariance matrix  $\Sigma$  into pixel space, we use a first-order Taylor expansion around the point  $\mu$  in the camera frame. Specifically, we compute the affine transformation  $\mathbf{J} \in \mathbb{R}^{2 \times 3}$  as follows:

$$\mathbf{J} = \begin{bmatrix} f_x/z_C & 0 & -f_x x_C/z_C^2 \\ 0 & f_y/z_C & -f_y y_C/z_C^2 \end{bmatrix}. \quad (26)$$

In the world space, the 3D Gaussian function  $\mathcal{G}$  is defined as:

$$\mathcal{G}_{\Sigma}(\mathbf{x} - \mu) = \exp\left(\frac{-1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right) \quad (27)$$

where  $\Sigma$  and  $\mu$  are parameters for the Gaussian distribution.



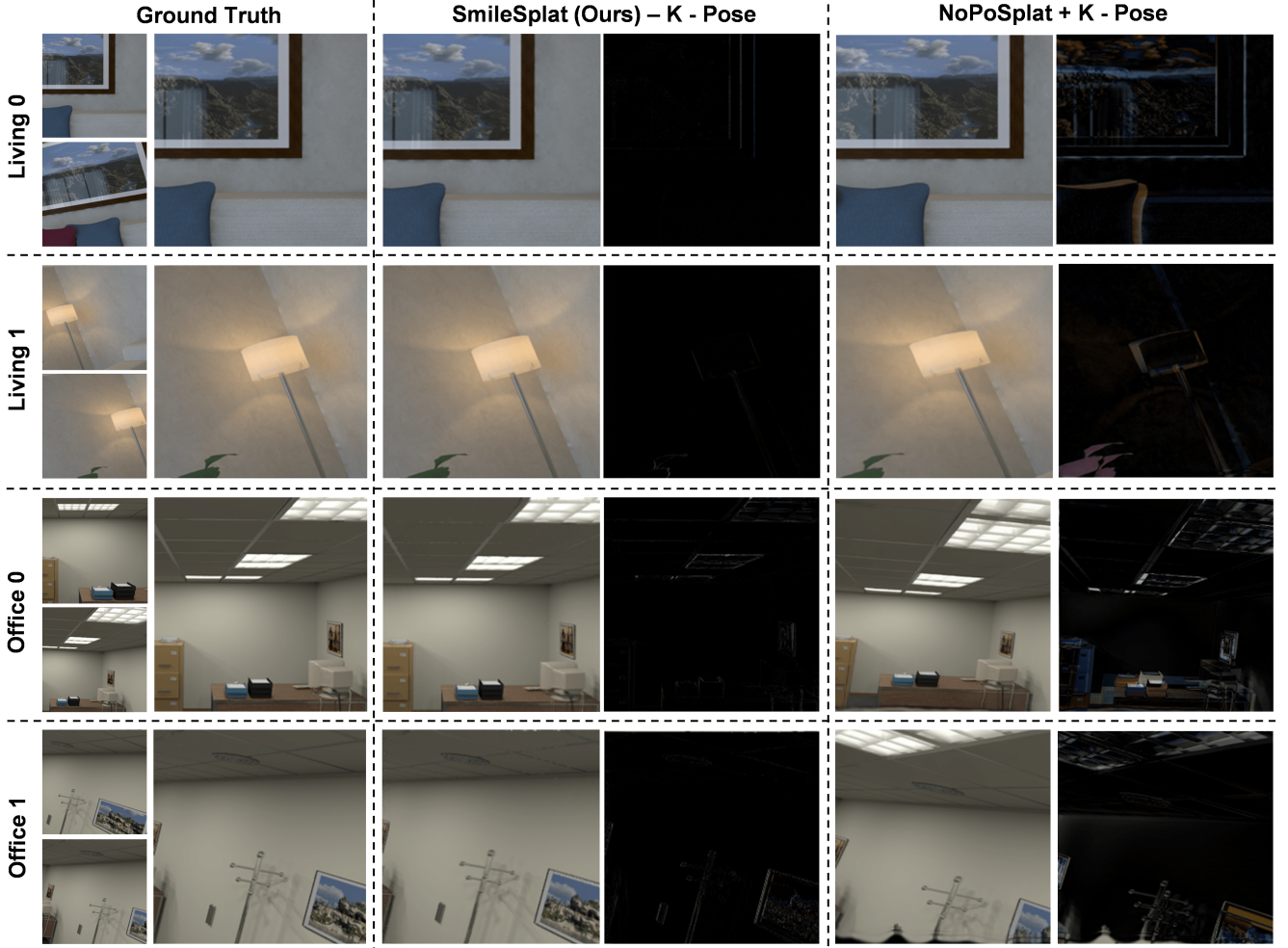


Figure 6. Comparisons of novel view rendering on the ICL-NUIM [14] dataset are presented. The rendering results include the RGB images in the left column and the differences from the ground truth in the right column.  $\pm K$  and  $\pm$  Pose indicate whether the intrinsics and extrinsics are required or not, respectively.

When we transform the 3D Gaussian from the world coordinates to the camera coordinates via the camera pose  $\mathbf{T} = \begin{bmatrix} \mathbf{W} & \mathbf{t} \\ 0 & 1 \end{bmatrix} \in SE(3)$ , the transformed Gaussian can be represented as  $\mathcal{G}_{\mathbf{W}\Sigma\mathbf{W}^T}(\mathbf{W}\mathbf{x} - \mathbf{W}\boldsymbol{\mu})$ .

When we project the Gaussian function into screen space using the affine approximation defined in Equation 25, the covariance matrix of the 2D Gaussian can be expressed as  $\mathbf{JW}\Sigma\mathbf{W}^T\mathbf{J}^T$ . The optimization for the intrinsic matrix can then be derived based on the relationships represented in Equations 23 and 26.

Method		NoPoSplat [45]				SmileSplat (Ours)			
		Small	Medium	Large	Ave.	Small	Medium	Large	Ave.
Living 0	PSNR	26.80	25.01	25.71	25.84	<b>36.59</b>	<b>35.11</b>	<b>33.21</b>	<b>34.97</b>
	SSIM	0.747	0.746	0.766	0.753	<b>0.945</b>	<b>0.971</b>	<b>0.913</b>	<b>0.943</b>
	LPIPS	0.298	0.276	0.256	0.277	<b>0.018</b>	<b>0.021</b>	<b>0.029</b>	<b>0.022</b>
Living 1	PSNR	27.07	25.66	25.40	26.04	<b>35.44</b>	<b>34.89</b>	<b>34.38</b>	<b>34.90</b>
	SSIM	0.701	0.748	0.761	0.736	<b>0.932</b>	<b>0.947</b>	<b>0.943</b>	<b>0.941</b>
	LPIPS	0.254	0.261	0.266	0.260	<b>0.020</b>	<b>0.025</b>	<b>0.025</b>	<b>0.023</b>
Office 0	PSNR	26.17	27.07	25.64	26.29	<b>36.21</b>	<b>34.10</b>	<b>30.06</b>	<b>33.46</b>
	SSIM	0.778	0.775	0.821	0.791	<b>0.973</b>	<b>0.985</b>	<b>0.930</b>	<b>0.963</b>
	LPIPS	0.234	0.240	0.285	0.253	<b>0.023</b>	<b>0.012</b>	<b>0.053</b>	<b>0.029</b>
Office 1	PSNR	26.65	25.21	21.30	24.39	<b>35.66</b>	<b>34.27</b>	<b>32.89</b>	<b>34.27</b>
	SSIM	0.788	0.771	0.727	0.762	<b>0.921</b>	<b>0.977</b>	<b>0.931</b>	<b>0.943</b>
	LPIPS	0.259	0.249	0.291	0.266	<b>0.020</b>	<b>0.019</b>	<b>0.037</b>	<b>0.025</b>

Table 4. Cross-Dataset generalization tests on the ICL-NUIM [14] dataset.

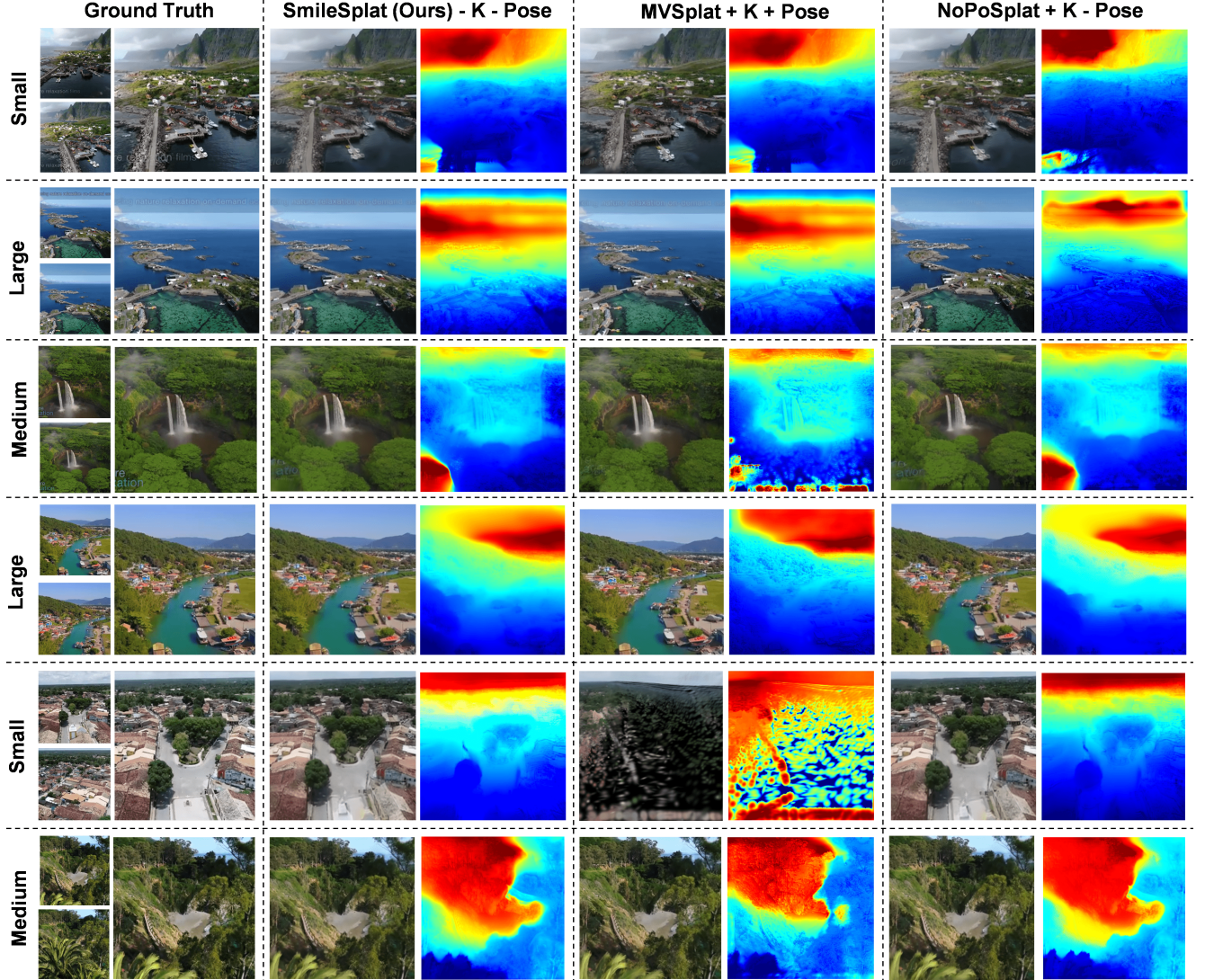


Figure 7. Comparisons of novel view rendering on the ACID [24] dataset are presented. The rendering results include the RGB images in the left column and the differences from the ground truth in the right column.  $\pm K$  and  $\pm$  Pose indicate whether the intrinsics and extrinsics are required or not, respectively.

## B. Comparison of Novel View Rendering

### B.1. Qualitative and quantitative Results on the ICL-NUIM and ACID dataset

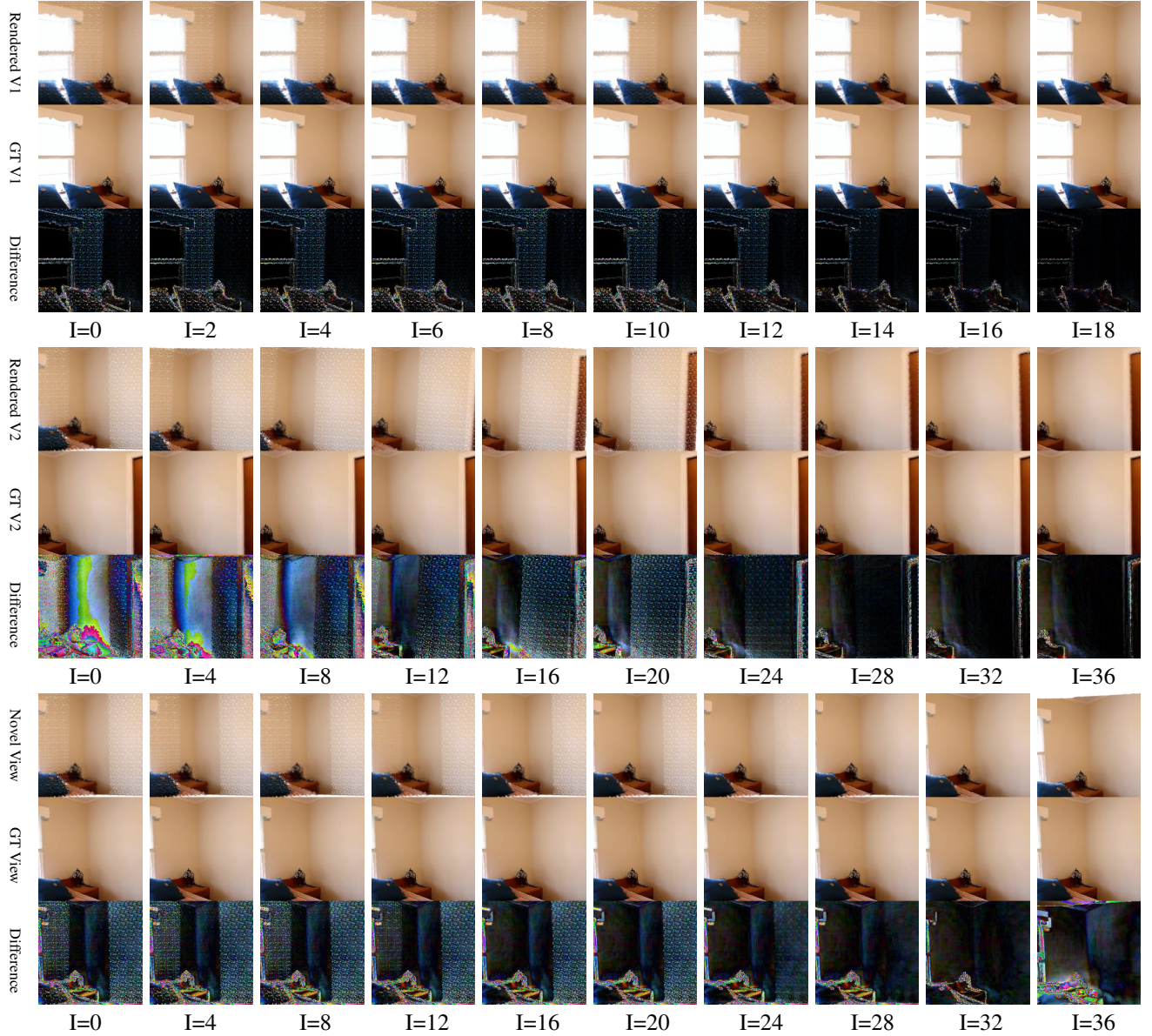
As mentioned in Section 4.2, both qualitative and quantitative results on the ICL-NUIM [14] dataset are evaluated for NoPoSplat [45] and our method.

As shown in Table 4, four sequences from living room and office room scenes are selected to evaluate these two approaches. Similar to the cross-dataset evaluation in Table 3, the inputs are grouped into three categories: Small overlap, Medium overlap, and Large overlap. For the average values, our method achieves a PSNR of over 33.4, while

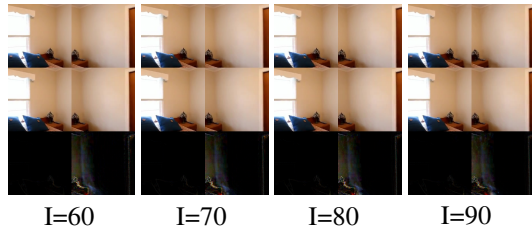
the best PSNR performance of NoPoSplat is below 26.3, demonstrating an improvement of up to 26%. Furthermore, this improvement is also reflected in two additional metrics.

In Figure 6, the reference images (two views) are shown in the left column. The third row (*Office 0*) illustrates the example of small-overlap pairs, while the second row (*Living 1*) shows the example of large-overlap inputs. By comparing the second and third rows, we can observe that our method demonstrates robust performance for both small and large overlap scenes. However, the rendering result of NoPoSplat in the third row contains more outliers, as shown in the difference map. For the Medium-overlap case (first row), alignment issues are noticeable in both the rendered

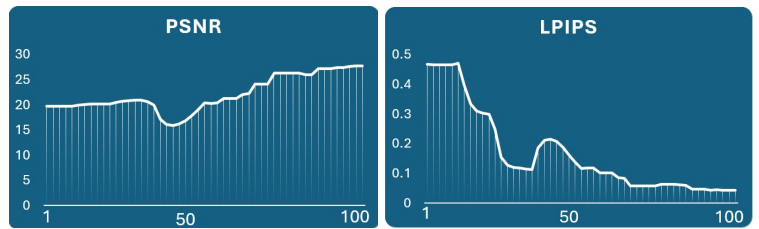




(a) Optimization (the first 36 iterations) for intrinsic and extrinsic parameters.



(b) Bundle adjusting Gaussian Splatting for two views.



(c) The curves of PSNR and LPIPS of novel view rendering during the refinement process

Figure 8. The entire optimization process. To clearly visualize the difference maps, we have magnified the error values by a **factor of 5**.

image and the difference map. In contrast, our method, SmileSplat, produces more accurate rendering results, ben-

efiting from the proposed bundle-adjusting Gaussian Splatting module, which refines both the estimated intrinsics and

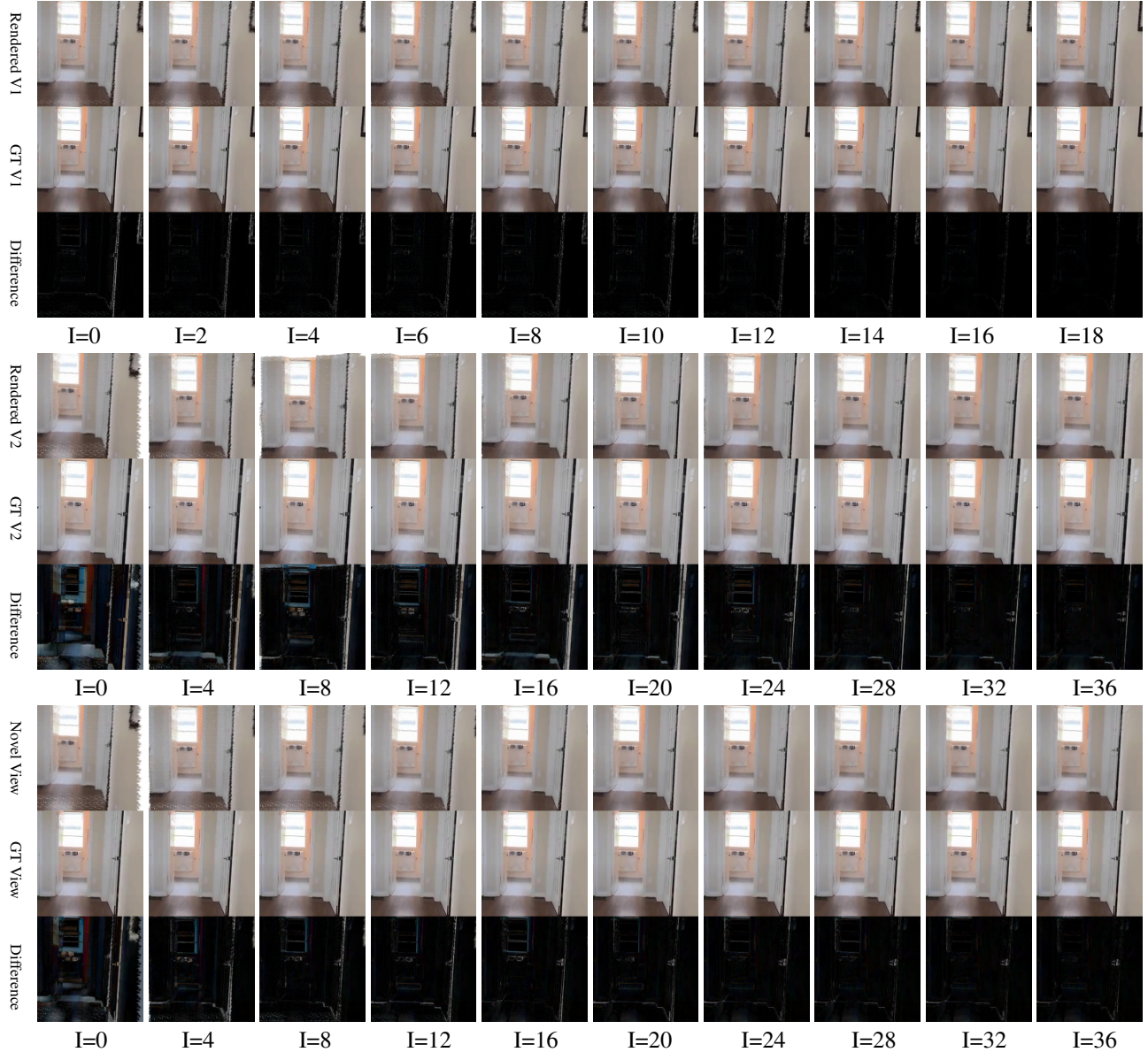


Figure 9. Optimization process (the first 36 iterations) for intrinsic and extrinsic parameters.

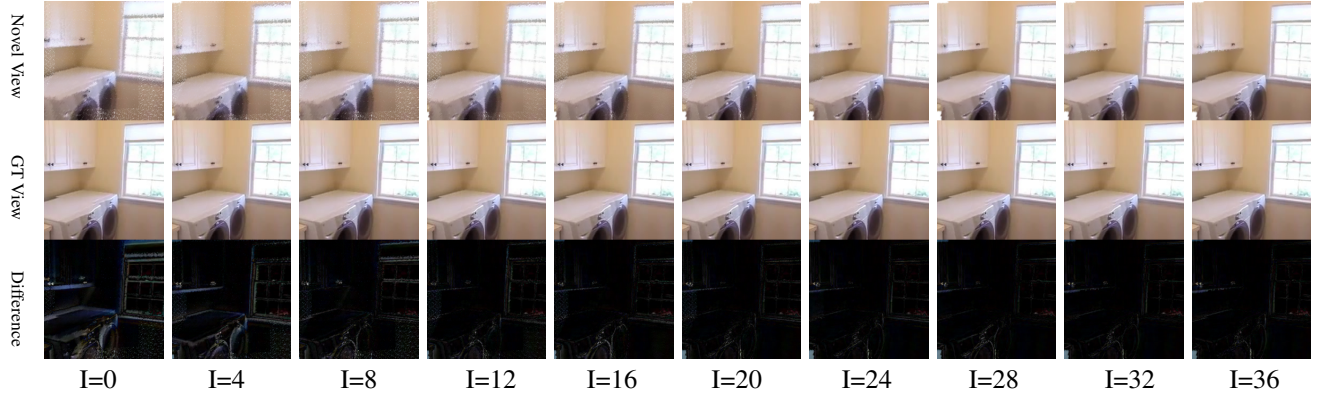
extrinsics.

As shown in Figure 7, many natural scenes from the ACID [24] dataset are used for testing the performance of the proposed method and state-of-the-art approaches, MVSplat (CamPara-Required) and NoPoSplat.

As shown in Figure 7, a variety of natural scenes from the ACID dataset are used to evaluate the performance of the proposed method, along with two state-of-the-art approaches: MVSplat (CamPara-Required) and NoPoSplat. These scenes, which cover diverse real-world environments, serve as a comprehensive benchmark for comparing the accuracy and robustness of each method under different con-

ditions, such as varying scene complexity, lighting, and object types. Compared to man-made or indoor scenarios, the scenes in the wild present entirely different challenges. However, the performance of the proposed method in both RGB and depth rendering tasks remains robust and accurate. For example, in the third row, two views are captured in a waterfall environment. The rendered RGB image produced by NoPoSplat shows misalignment issues, whereas our method maintains accurate rendering with better alignment and detail preservation.

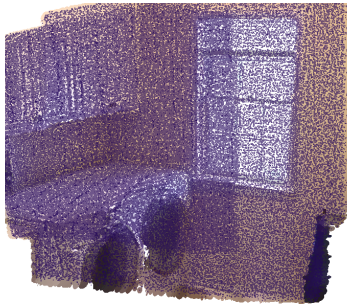




(a) First 36 iterations in refinement for View 1, View 2, and the novel viewpoint.



(b) 3D Gaussian map



(c) Visual overlaps



(d) Multi-view visualization.

Figure 10. Visualization for the optimization process and Gaussian Splatting fields.

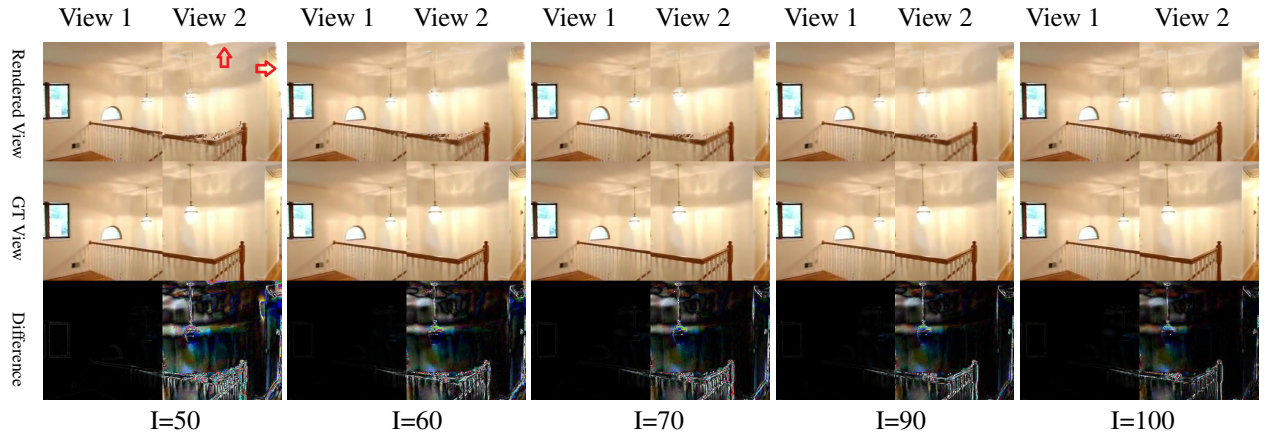


Figure 11. Bundle adjusting Gaussian Splatting refinement for View 1 and View 2. For clearly visualizing difference maps, we have magnified the error values by a **factor of 5**.

## C. Ablation studies

### C.1. Iterative Optimization for Intrinsic and Relative Pose Estimation

In the refinement stage, we first optimize the intrinsic matrix by **rendering View 1 for the first 10 iterations**, as shown in Figures 8 and 9. During this stage, only the intrinsic ma-

trix is refined, and the Gaussians are not optimized. As a result, we are still able to detect the photometric distance between the rendered and observed images. The quality of the rendered image improves further over the next 10 iterations of rendering View 1.

Once the optimization of the intrinsic matrix is completed, it is fixed during the second view rendering process. In the first 20 iterations of rendering View 2, neither the in-

intrinsic parameters nor the Gaussians are updated based on the photometric error generated from View 2. As shown in the process for handling View 2, the values in the difference map are large at the beginning due to the noisy initial camera pose. However, the photometric distance progressively decreases as the relative pose estimation converges. When comparing the rendered results at iterations  $I = 8$  and  $I = 16$ , we observe that the large error regions shift, indicating that it is difficult to eliminate all errors by solely adjusting the camera pose. Therefore, we continue to optimize the Gaussians in subsequent iterations, taking the camera poses into account.

In particular, View 2 detects some newly revealed regions, as shown in Figure 8 (rows of View 2). These regions can be further improved by optimizing the Gaussians over a few more iterations. In these novel view rendering steps, following the approach outlined by methods such as [45], we optimize only the camera poses, without incorporating the observed image in the Gaussian optimization process.

### C.2. Bundle Adjusting Gaussian Splatting based on View 1 and View 2

As shown in Figures 8(b), 10, and 11, the results from the last 50 iterations of bundle adjustment are presented. In particular, Figure 11 (Red Arrow) shows that at the beginning of the optimization stage, some regions are still misaligned due to the initial inaccuracy of the camera pose. As the optimization progresses, the camera pose is refined incrementally, and we can observe that these regions continue to improve, leading to better rendering results by the 100<sup>th</sup> iteration.

### C.3. Limitation

As shown in Figure 11, it is important to note a limitation of this module, as seen in this figure. For the stair guardrail region, while the rendering quality improves, we cannot achieve perfect results due to the limitations of the initial geometry, which is not accurate for the stair guardrail. Additionally, given that the views from View 1 and View 2 are limited, we are unable to significantly improve the geometry in this region, which is also a future direction for further improvement by leveraging additional multi-view geometric constraints into the system.

## References

- [1] Gwangbin Bae and Andrew J. Davison. Rethinking inductive biases for surface normal estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3, 4
- [2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. 2
- [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19697–19705, 2023. 2
- [4] David Charatan, Sizhe Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. *arXiv preprint arXiv:2312.12337*, 2023. 2, 3, 6
- [5] David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19457–19467, 2024. 6, 7
- [6] 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. *arXiv preprint arXiv:2403.14627*, 2024. 2, 3, 6, 7
- [7] Kai Cheng, Xiaoxiao Long, Kaizhi Yang, Yao Yao, Wei Yin, Yuexin Ma, Wenping Wang, and Xuejin Chen. Gaussianpro: 3d gaussian splatting with progressive propagation. *arXiv preprint arXiv:2402.14650*, 2024. 3
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2, 4
- [9] Yilun Du, Cameron Smith, Ayush Tewari, and Vincent Sitzmann. Learning to render novel views from wide-baseline stereo pairs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4970–4980, 2023. 6, 7
- [10] Johan Edstedt, Qiyu Sun, Georg Bökman, Mårten Wadenbäck, and Michael Felsberg. Roma: Robust dense feature matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19790–19800, 2024. 6
- [11] Zhiwen Fan, Wenyan Cong, Kairun Wen, Kevin Wang, Jian Zhang, Xinghao Ding, Danfei Xu, Boris Ivanovic, Marco Pavone, Georgios Pavlakos, et al. Instantplat: Unbounded sparse-view pose-free gaussian splatting in 40 seconds. *arXiv preprint arXiv:2403.20309*, 2024. 2, 3
- [12] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. 2
- [13] Yang Fu, Sifei Liu, Amey Kulkarni, Jan Kautz, Alexei A Efros, and Xiaolong Wang. Colmap-free 3d gaussian splatting. *arXiv preprint arXiv:2312.07504*, 2023. 3
- [14] Ankur Handa, Thomas Whelan, John McDonald, and Andrew J Davison. A benchmark for rgb-d visual odometry, 3d reconstruction and slam. In *2014 IEEE international confer-*

- ence on Robotics and automation (ICRA), pages 1524–1531. IEEE, 2014. 6, 11, 12
- [15] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. 3
- [16] Heiko Hirschmüller. Accurate and efficient stereo processing by semi-global matching and mutual information. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, pages 807–814. IEEE, 2005. 1
- [17] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for geometrically accurate radiance fields. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 3
- [18] Letian Huang, Jiayang Bai, Jie Guo, Yuanqi Li, and Yanwen Guo. On the error analysis of 3d gaussian splatting and an optimal projection strategy. *CoRR*, 2024. 10
- [19] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, 2006. 1
- [20] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 1, 2, 3, 4, 9
- [21] Vincent Leroy, Johann Cabon, and Jerome Revaud. Grounding image matching in 3d with mast3r, 2024. 2, 6, 7
- [22] Yanyan Li, Raza Yunus, Nikolas Brasch, Nassir Navab, and Federico Tombari. Rgb-d slam with structural regularities. In *2021 IEEE international conference on Robotics and automation (ICRA)*, pages 11581–11587. IEEE, 2021. 1, 3
- [23] Yanyan Li, Chenyu Lyu, Yan Di, Guangyao Zhai, Gim Hee Lee, and Federico Tombari. Geogaussian: Geometry-aware gaussian splatting for scene rendering. In *European Conference on Computer Vision*, pages 441–457. Springer, 2025. 2, 10
- [24] Andrew Liu, Richard Tucker, Varun Jampani, Ameesh Makadia, Noah Snavely, and Angjoo Kanazawa. Infinite nature: Perpetual view generation of natural scenes from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14458–14467, 2021. 6, 7, 12, 14
- [25] Hidenobu Matsuki, Riku Murai, Paul HJ Kelly, and Andrew J Davison. Gaussian splatting slam. *arXiv preprint arXiv:2312.06741*, 2023. 1, 3, 5
- [26] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2
- [27] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–15, 2022. 2
- [28] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015. 1, 5
- [29] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4531–4540, 2019. 1
- [30] Linfei Pan, Dániel Baráth, Marc Pollefeys, and Johannes L Schönberger. Global structure-from-motion revisited. 1
- [31] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 4
- [32] Johannes L Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 1
- [33] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 4, 5
- [34] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 6
- [35] Brandon Smart, Chuanxia Zheng, Iro Laina, and Victor Adrian Prisacariu. Splatt3r: Zero-shot gaussian splatting from uncalibrated image pairs. *arXiv preprint arXiv:2408.13912*, 2024. 6, 7
- [36] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 6, 8
- [37] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022. 2
- [38] Shuzhe Wang, Vincent Leroy, Johann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. *arXiv preprint arXiv:2312.14132*, 2023. 2, 3, 4, 6, 7
- [39] Philippe Weinzaepfel, Vincent Leroy, Thomas Lucas, Romain Brégier, Johann Cabon, Vaibhav Arora, Leonid Antsfeld, Boris Chidlovskii, Gabriela Csurka, and Jérôme Revaud. Croco: Self-supervised pre-training for 3d vision tasks by cross-view completion. *Advances in Neural Information Processing Systems*, 35:3502–3516, 2022. 2, 4
- [40] Philippe Weinzaepfel, Thomas Lucas, Vincent Leroy, Johann 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17969–17980, 2023. 4
- [41] 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. *arXiv preprint arXiv:2310.08528*, 2023. 3

- [42] Haolin Xiong, Sairisheek Muttukuru, Rishi Upadhyay, Pradyumna Chari, and Achuta Kadambi. Sparsegs: Real-time 360  $\{\deg\}$  sparse view synthesis using gaussian splatting. *arXiv preprint arXiv:2312.00206*, 2023. 3
- [43] Chi Yan, Delin Qu, Dong Wang, Dan Xu, Zhigang Wang, Bin Zhao, and Xuelong Li. Gs-slam: Dense visual slam with 3d gaussian splatting. *arXiv preprint arXiv:2311.11700*, 2023. 1, 3
- [44] Zeyu Yang, Hongye Yang, Zijie Pan, Xiatian Zhu, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. *arXiv preprint arXiv:2310.10642*, 2023. 3
- [45] Botao Ye, Sifei Liu, Haofei Xu, Xueting Li, Marc Pollefeys, Ming-Hsuan Yang, and Songyou Peng. No pose, no problem: Surprisingly simple 3d gaussian splats from sparse unposed images. *arXiv preprint arXiv:2410.24207*, 2024. 2, 3, 6, 7, 8, 11, 12, 16
- [46] Vickie Ye, Ruilong Li, Justin Kerr, Matias Turkulainen, Brent Yi, Zhuoyang Pan, Otto Seiskari, Jianbo Ye, Jeffrey Hu, Matthew Tancik, et al. gsplat: An open-source library for gaussian splatting. *arXiv preprint arXiv:2409.06765*, 2024. 10
- [47] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 2, 6, 7
- [48] Jing Yuan, Shuhao Zhu, Kaitao Tang, and Qinxuan Sun. Orbedm: An rgb-d slam approach fusing orb triangulation estimates and depth measurements. *IEEE Transactions on Instrumentation and Measurement*, 71:1–15, 2022. 1
- [49] Zhengyou Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22(11):1330–1334, 2000. 1
- [50] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. 6, 7, 8
- [51] Zunjie Zhu, Youxu Fang, Xin Li, Chengang Yan, Feng Xu, Chau Yuen, and Yanyan Li. Robust gaussian splatting slam by leveraging loop closure. *arXiv preprint arXiv:2409.20111*, 2024. 1
- [52] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa splatting. *IEEE Transactions on Visualization and Computer Graphics*, 8(3):223–238, 2002. 10