

Self-Ensembling Gaussian Splatting for Few-Shot Novel View Synthesis

Chen Zhao¹ Xuan Wang² Tong Zhang¹ Saqib Javed¹ Mathieu Salzmann^{1,3}
¹EPFL ²Ant Group ³Swiss Data Science Center
chen.zhao@epfl.ch

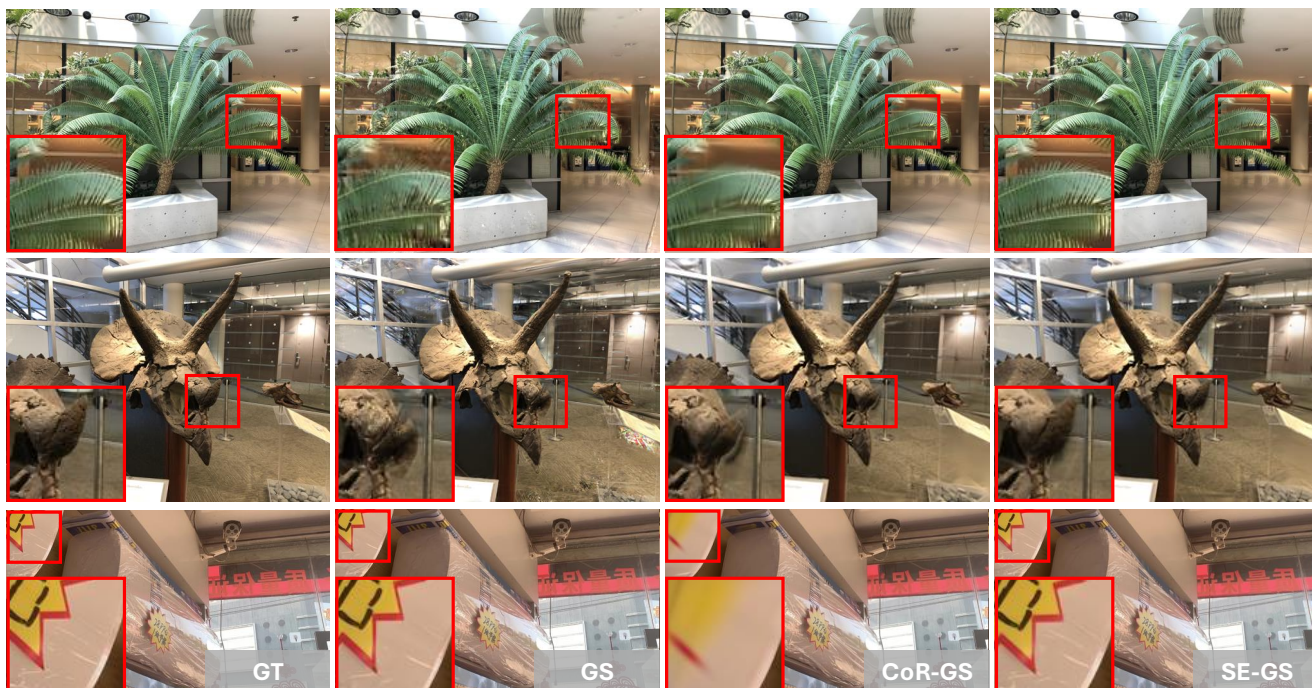


Figure 1. **Qualitative results of our SE-GS and state-of-the-art approaches.** The models are trained on sparse views and the images rendered from novel views are shown. As highlighted in the zoomed-in patches, our SE-GS captures finer details and produces fewer artifacts for novel views when trained on few-shot images.

Abstract

3D Gaussian Splatting (3DGS) has demonstrated remarkable effectiveness for novel view synthesis (NVS). However, the 3DGS model tends to overfit when trained with sparse posed views, limiting its generalization ability to novel views. In this paper, we alleviate the overfitting problem, presenting a Self-Ensembling Gaussian Splatting (SE-GS) approach. Our method encompasses a Σ -model and a Δ -model. The Σ -model serves as an ensemble of 3DGS models that generates novel-view images during inference. We achieve the self-ensembling by introducing an uncertainty-aware perturbation strategy at the training state. We complement the Σ -model with the Δ -model, which is dynamically perturbed based on the uncertainties of novel-view renderings across different training steps. The perturbation yields diverse temporal samples in the Gaus-

sian parameter space without additional training costs. The geometry of the Σ -model is regularized by penalizing discrepancies between the Σ -model and these temporal samples. Therefore, our SE-GS conducts an effective and efficient regularization across a large number of 3DGS models, resulting in a robust ensemble, the Σ -model. Our experimental results on the LLFF, Mip-NeRF360, DTU, and MVImgNet datasets show that our approach improves NVS quality with few-shot training views, outperforming existing state-of-the-art methods. The code is released at: [project page](#).

1. Introduction

Novel view synthesis (NVS) is a critical task [49] in computer vision and graphics, playing a pivotal role in applications such as virtual reality [8], augmented reality [50],

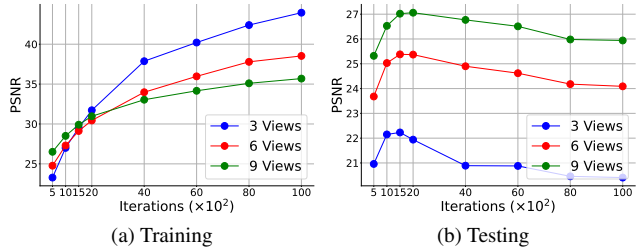


Figure 2. **Overfitting in 3D Gaussian Splatting with sparse training views.** 2a and 2b illustrate the performance of 3DGS on training and testing views, respectively. Each curve represents the PSNR values across different training iterations.

and 3D content generation [16, 42]. The objective of NVS is to generate photo-realistic images from previously unseen viewpoints. Typically, NVS starts by constructing a 3D representation [28, 39] from a set of existing 2D observations. In recent years, 3D Gaussian Splatting (3DGS) [6, 21, 46] has emerged as a powerful representation, integrating the advantages of both explicit [35] and implicit [28] representations. This approach enables efficient novel view generation and yields promising synthesized results with densely sampled observations that cover a wide range of viewpoints. However, 3DGS tends to overfit the available views when only a limited number of images are provided. To illustrate this issue, in Fig. 2, we evaluate the 3DGS model trained on sparse images with different numbers of iterations. The performance on the training data consistently improves as the number of iterations increases, while the testing results deteriorate after 2000 iterations. Moreover, the overfitting problem becomes more noticeable with fewer training views, such as when using only 3 views.

To mitigate overfitting, self-ensembling [11, 22, 29] has been highlighted as an effective strategy in tasks such as detection [48] and segmentation [25]. Nevertheless, its use for NVS remains unstudied, and how to exploit it in a 3DGS formalism is an open question. Therefore, in this paper, we bridge this gap, introducing a new 3DGS method that enhances the quality of novel view synthesis with sparse training views via *self-ensembling*. Our method consists of two models, a Σ -model and a Δ -model. The Σ -model constitutes an ensemble of 3DGS models, which is used to generate novel views during inference. A straightforward way to achieve ensembling would be to utilize multiple 3DGS models that correspond to different samples in the Gaussian parameter space, thereby increasing the likelihood of reaching a more robust solution. Nevertheless, as we will demonstrate in Sec. 4, the computational cost of training multiple 3DGS models limits the ability to produce sufficiently diverse samples. By contrast, we complement the Σ -model with a single Δ -model and obtain diverse samples by perturbing the Δ -model. Specifically, instead of applying ran-

dom perturbation, which often results in invalid samples, we perturb the Δ -model based on uncertainties derived from the training data. We store images generated by rendering the Δ -model from a set of pseudo views at different training iterations in buffers and calculate pixel-level uncertainties across images within each buffer. We then add random noise to the Gaussian parameters of the Δ -model associated with pixels that have high uncertainty scores. Consequently, our perturbation strategy generates uncertainty-aware samples without introducing significant training overhead.

Building upon the uncertainty-aware samples derived from the Δ -model, we construct an ensemble by leveraging a self-ensembling mechanism. We train the Σ -model over the training views without the aforementioned perturbation and regularize the geometry of the Σ -model by minimizing its discrepancies with the perturbed Δ -model. These discrepancies are measured via a photometric loss between images synthesized from the pseudo views. Since the regularization is performed based on diverse samples in the Gaussian parameter space, the self-ensembling process enhances the robustness of the Σ -model, thereby improving its generalization ability to novel views. Moreover, our self-ensembling technique is independent of any additional ground-truth signals as the regularization is carried out in a self-supervised manner.

We conduct experiments on multiple datasets, including LLFF [27], DTU [20], Mip-NeRF360 [3], and MVImgNet [45], with limited numbers of training views. Our SE-GS achieves the best performance across all datasets in the sparse-view setting, surpassing the state-of-the-art approaches. Furthermore, we perform a comprehensive analysis of our method to demonstrate its effectiveness and efficiency. To the best of our knowledge, we are the first to explore the potential of the self-ensembling mechanism in 3DGS for few-shot novel view synthesis. The code will be publicly available upon acceptance.

2. Related Work

Radiance field modeling. Recently, Neural Radiance Field (NeRF) [3, 28, 32, 41] has brought groundbreaking advancements to novel view synthesis. NeRF learns the density and color information of a scene or object from multi-view posed images, enabling the synthesis of realistic images from novel views via volume rendering [9]. However, most NeRF-based methods are inefficient [18] during training and inference because they require numerous MLP queries in the rendering process. This limitation makes it challenging to support tasks with real-time requirements. In this context, 3D Gaussian Splatting [21] is developed as an efficient alternative for radiance field modeling. 3DGS relies on explicit representations [35] that allow for faster rendering and training [10]. While initially introduced as a scene-specific model, some 3DGS variants [5, 7, 26, 40]

explore the generalization ability of 3DGS towards novel scenes, optimizing the Gaussian parameters based on learnable cost volumes [15]. In this paper, we stick to the vanilla 3DGS rendering pipeline, focusing on the scene-level radiance field modeling.

3DGS with few-shot images. The insufficiency of training views is a primary factor leading to artifacts in novel view synthesis. Even when training views are relatively dense, some regions may still be observed from few-shot images, resulting in quality degradation in the synthesized views. To handle such an under-constrained problem, some approaches incorporate additional ground-truth signals into the 3DGS pipeline. For instance, DNGaussian [24], CoherentGS [31], and FSGS [53] utilize a monocular depth estimator [34] to predict depth maps, enabling the 3DGS model to be trained with both a photometric and a depth loss. Other methods, such as [43] and [17], leverage multi-view stereo techniques [37] to generate novel-view images as ground truth. However, ground-truth data obtained from off-the-shelf methods is inevitably noisy, which potentially impacts the training of 3DGS. As noted in [47] and observed in our experiments, when the number of training views increases, the performance of these methods is sometimes worse than the vanilla 3DGS. To overcome this problem, Zhang *et al.* [47] propose training multiple 3DGS models with a cross-model regularization term. However, training additional 3DGS models incurs a significant computational cost, making it impractical to scale up to a large number of 3DGS models for stronger regularization.

Ensembling. Ensembling has been evidenced as a powerful technique in machine learning [1, 14, 52] to improve model robustness and generalization by aggregating predictions from multiple models. Traditionally, ensembles [13, 23, 33] are created by training independent models and averaging their outputs or applying a voting mechanism. To improve the efficiency, some self-ensembling methods, such as temporal ensembling [22] and consistency regularization [11], leverage variations of a single model across training iterations to build an ensemble. Motivated by the success of self-ensembling, we present the first Self-Ensembling Gaussian Splatting approach, in which we generate diverse samples in the Gaussian parameter space through an uncertainty-aware perturbation mechanism.

3. Method

3.1. Preliminaries and Motivations

3D Gaussian Splatting [21] represents a scene as a collection of Gaussians. These Gaussians are splatted onto the 2D image plane during rendering. Formally, we denote the set of N Gaussians as $\{G_i, i = 1, 2, \dots, N\}$. Each Gaussian G_i is defined as $G_i = (\mu_i, \Sigma_i, \mathbf{h}_i, o_i)$, where μ_i is the 3D position, Σ_i is the covariance matrix, \mathbf{h}_i represents spherical

harmonics (SH) coefficients associated with the Gaussian, and o_i indicates opacity. Each Gaussian contributes to a 3D point \mathbf{x} according to the 3D Gaussian distribution

$$\mathcal{N}_{3d}^i(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\mu_i)^T \Sigma_i^{-1}(\mathbf{x}-\mu_i)}. \quad (1)$$

The parameters of the Gaussians are then learned from the input views. To ensure that Σ_i remains positive semi-definite throughout optimization, it is decomposed into two learnable components as $\Sigma_i = \mathbf{R}\mathbf{S}\mathbf{S}^T\mathbf{R}^T$, where \mathbf{R} is a rotation matrix and \mathbf{S} stands for a scaling matrix. For RGB rendering, the projection from 3D space to a 2D image plane relies on the projection matrix \mathbf{W} to compute the projected 2D covariance matrix

$$\Sigma_i' = \mathbf{J}\mathbf{W}\Sigma_i\mathbf{W}^T\mathbf{J}^T, \quad (2)$$

where \mathbf{J} denotes the Jacobian of the affine approximation of the projection. The color of a pixel is obtained by performing alpha-blending as

$$\mathbf{c} = \sum_{i=1}^M \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (3)$$

where M is the number of Gaussians covering the pixel, \mathbf{c}_i denotes the color of the Gaussian derived from the SH coefficients, and α_i is computed from the 2D covariance matrices and opacity scores. For better convergence, the Gaussian parameters are initialized based on 3D points obtained using structure-from-motion (SFM) techniques [35, 36]. These parameters are optimized via gradient descent with a photometric loss [21] where the posed training images serve as ground truth.

In this paper, we focus on sparse-view scenarios, where only few-shot training images are provided. We denote the 3DGS model trained on these images as \mathcal{G} , which represents a sample in the Gaussian parameter space. In this setting, \mathcal{G} is prone to overfitting the training data and thereby getting trapped in a suboptimal solution, which ultimately degrades the quality of the synthesized novel views. To handle the overfitting problem, a promising solution is to train a set of 3DGS models $\{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_k\}$ that correspond to different samples. By applying regularization across this set, each model is guided away from specific local optima, facilitating a more thorough exploration of the parameter space. As we will report in Sec. 4, the performance on testing data improves as k grows. However, training multiple 3DGS models incurs a significant computational cost, which becomes unaffordable with a large k . This limitation prevents the method from scaling up for better performance. Consequently, we propose a novel training strategy based on self-ensembling [11, 22], which enhances 3DGS for few-shot NVS with negligible additional training overhead. The pipeline is illustrated in Fig. 3 and the details will be elaborated in the remainder of this section.

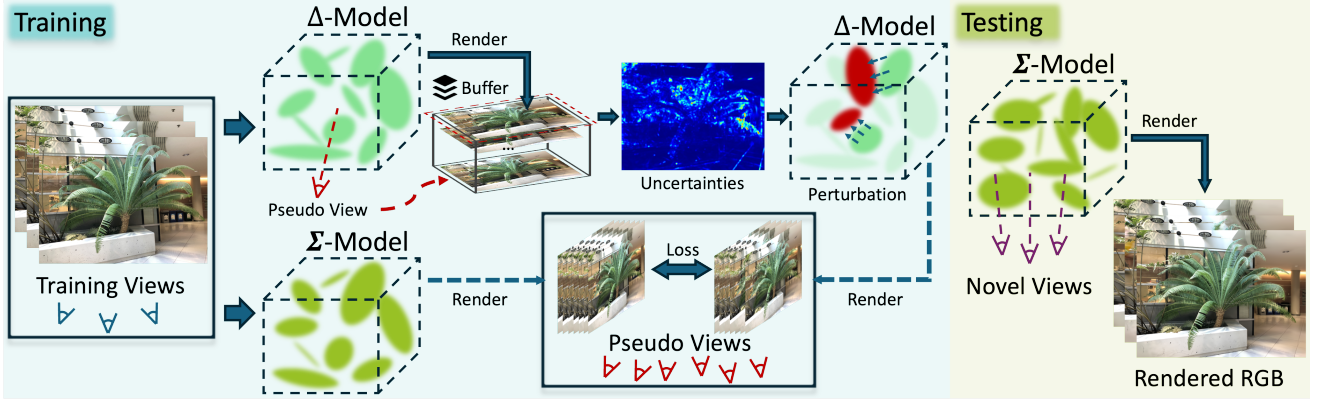


Figure 3. **Pipeline of the presented SE-GS.** We tackle the overfitting problem in sparse-view scenarios by incorporating a self-ensembling mechanism into the 3DGS framework. Our self-ensembling builds upon diverse samples in the Gaussian parameter space. These samples are derived from the Δ -model via an uncertainty-aware perturbation strategy. Specifically, we store images rendered from pseudo views at different training steps in buffers, from which we compute pixel-level uncertainties. We then perturb the Gaussians overlapping the pixels with high uncertainties, as highlighted as red ellipses. Self-ensembling over the resulting samples is achieved by training a Σ -model with a regularization that penalizes the discrepancies of the Σ -model and the perturbed Δ -model. During inference, novel view synthesis is performed using the Σ -model.

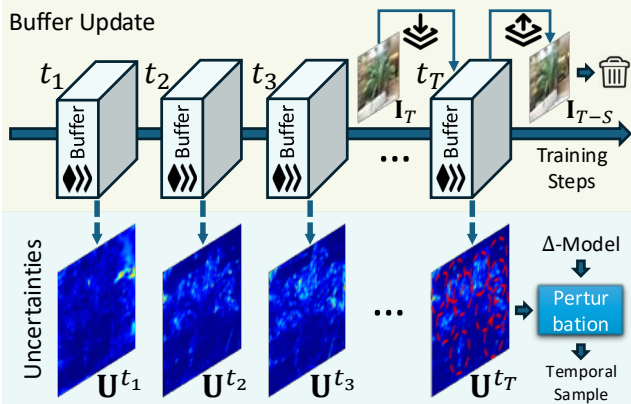


Figure 4. **Buffer update during training.** For each sampled pseudo view, we dynamically update the buffer that stores the images rendered at different training steps. For instance, at training step t_T , the oldest image I_{T-S} in the buffer is popped, and the new image I_T is pushed into the buffer. An uncertainty map U^{tT} is computed based on the current buffer, which is then employed to compute a perturbation that results in a new 3DGS model.

3.2. Uncertainty-Aware Perturbation

In contrast to generating multiple samples by training separate 3DGS models, we dynamically create temporal samples from a single Δ -model during training. Note that, in this paper, *temporal* refers to the dynamic nature of samples that vary across training iterations. Specifically, we train the Δ -model on the available posed images, following the same optimization and density control strategies introduced in 3DGS [21]. At each training iteration, the current Δ -

model represents a specific sample in the Gaussian parameter space based on the information contained in the training data. A temporal 3DGS model is then sampled from the neighboring region of the current Δ -model in the parameter space by perturbing the Δ -model as

$$\hat{G}_\Delta^t = G_\Delta^t + \delta_t, \quad (4)$$

where G_Δ^t stands for a Gaussian in the Δ -model at training step t , δ_t indicates the perturbation, and \hat{G}_Δ^t denotes a perturbed Gaussian in the temporal model. To achieve the perturbation, a naive approach is to sample random noise from a normal distribution, i.e., $\delta_t \in \mathcal{N}(\mu_t, \sigma_t^2)$. However, such perturbation is independent of the training data, and this method thus tends to yield invalid samples.

Therefore, we present an uncertainty-aware perturbation strategy, leveraging the statistics of renderings. Our approach starts by creating M pseudo views through interpolation between the training views. Specifically, given two cameras sampled from the training views, the camera extrinsics of a pseudo view are computed as

$$\hat{\mathbf{R}} = \text{SLERP}(\mathbf{R}_1, \mathbf{R}_2, \beta), \quad (5)$$

$$\hat{\mathbf{c}} = \beta \mathbf{c}_1 + (1 - \beta) \mathbf{c}_2, \quad (6)$$

$$\hat{\mathbf{T}} = -\hat{\mathbf{R}}\hat{\mathbf{c}}, \quad (7)$$

where $(\mathbf{R}_1, \mathbf{R}_2)$ represent the rotations of the sampled cameras, $(\mathbf{c}_1, \mathbf{c}_2)$ indicate the sampled camera centers, SLERP denotes the spherical linear interpolation [4], β is a scalar that controls the interpolation, and $(\hat{\mathbf{R}}, \hat{\mathbf{T}})$ stand for the camera parameters of the pseudo view. As illustrated in

Fig. 3, for each pseudo view, we render an RGB image using the current Δ -model and store the renderings at different training steps in a buffer. We then compute a pixel-wise uncertainty map

$$\mathbf{U} = \sqrt{\frac{1}{S} \sum_{i=1}^S (\mathbf{I}_i - \bar{\mathbf{I}})^2}, \quad (8)$$

of the same resolution as the rendered image, i.e., $\mathbf{U} \in \mathbb{R}^{H \times W}$, where \mathbf{I}_i represents an image in the buffer, $\bar{\mathbf{I}}$ indicates the mean of these images, and S is the buffer size. This uncertainty estimation is carried out over all M sampled pseudo views in parallel. In practice, to enhance robustness, we perform local smoothing over each uncertainty map, which yields

$$\hat{\mathbf{U}}(i, j) = \frac{1}{k^2} \sum_{m, n} \mathbf{U}(i + m, j + n), \quad (9)$$

where $k = 5$ is the kernel size. Subsequently, we apply an uncertainty-aware perturbation to each of the Gaussians in the Δ -model as

$$\hat{G}_{\Delta}^t = G_{\Delta}^t + \delta_t h(G_{\Delta}^t, \hat{\mathcal{U}}^t), \quad (10)$$

where $\hat{\mathcal{U}}^t$ is the set of the uncertainty maps computed from the buffers at the current training step, and $h(\cdot, \cdot)$ is an indicator function defined as

$$h(G_{\Delta}^t, \hat{\mathcal{U}}^t) = \begin{cases} 1 & \text{if } \max_{(i, j) \in \mathcal{P}(G_{\Delta}^t)} u_{ij} \geq \tau \\ 0 & \text{if } \max_{(i, j) \in \mathcal{P}(G_{\Delta}^t)} u_{ij} < \tau, \end{cases} \quad (11)$$

with τ a predefined threshold, and u_{ij} the uncertainty score of a pixel sampled from the set of pixels \mathcal{P} overlapping with the 2D splats of G_{Δ}^t . Please refer to the supplementary material for details on defining τ and identifying \mathcal{P} . In short, the reliability of each Gaussian is connected with the uncertainties of the renderings, and only the unreliable Gaussians, characterized by high uncertainty scores, are perturbed. In practice, we add the perturbation to the 3D positions, 3D rotations, scales, and opacities of the Δ -model. In particular, since 3D rotation is not continuous when expressed as a 3D matrix, we perturb the 6D continuous representation [51] of the rotation, which yields

$$\hat{\mathbf{R}}_{\Delta}^t = f^{-1}(f(\mathbf{R}_{\Delta}^t) + \delta_t^R h(G_{\Delta}^t, \hat{\mathcal{U}}^t)), \quad \delta_t^R \in \mathbb{R}^6, \quad (12)$$

where $f(\cdot)$ denotes the mapping from rotation matrix to 6D representation. Fig. 4 illustrates the dynamic updates of the buffer throughout the training process. As the buffer is updated, the uncertainty map varies accordingly, resulting in diverse temporal samples from the Δ -model.

It is worth noting that the number of Gaussians in the Δ -model varies during training due to density control, making

it challenging to compute uncertainties directly in the Gaussian parameter space across different training steps. Our approach naturally handles this challenge as the statistics on 2D renderings provide a consistent ground to measure uncertainties, regardless of the varying number of Gaussians.

3.3. Self-Ensembling in 3DGS

Given the temporal samples around the Δ -model, the next step in our pipeline is to construct an ensemble from these samples. Concretely, we train a separate 3DGS model, named the Σ -model, on the posed images without perturbations. To prevent the Σ -model from being trapped at a suboptimal solution, we guide the training of the Σ -model with a regularization [47] formulated as

$$\mathcal{L}_r = (1 - \lambda) \|\mathbf{I}_{\Sigma}^t - \mathbf{I}_{\Delta}^t\|_1 + \lambda \mathcal{L}_{\text{D-SSIM}}(\mathbf{I}_{\Sigma}^t, \mathbf{I}_{\Delta}^t), \quad (13)$$

where $\lambda = 0.2$ is a predefined weight, $\mathcal{L}_{\text{D-SSIM}}$ denotes a D-SSIM term [21], and $(\mathbf{I}_{\Sigma}^t, \mathbf{I}_{\Delta}^t)$ represent the images rendered from a pseudo view using the current Σ -model and Δ -model, respectively. We also utilize a co-pruning strategy [47] to enhance the regularization. Therefore, the Σ -model aggregates information from different samples in the parameter space, serving as an ensemble of 3DGS models. The final loss function during training is defined as

$$\mathcal{L} = \mathcal{L}_{\text{RGB}} + \gamma \mathcal{L}_r, \quad (14)$$

where $\gamma = 1$ by default and \mathcal{L}_{RGB} represents a photometric loss over the training views. During testing, we keep the Σ -model for novel view synthesis. Notably, compared with the previous approach [47], our SE-GS inherently encodes a large number of 3DGS models without significant computational overhead, thereby enabling effective regularization efficiently. Furthermore, the regularization is carried out on pseudo views in a self-supervised manner, making it independent of additional information such as depth [24, 53].

4. Experiments

4.1. Setup

Datasets. We conduct experiments on four datasets, i.e., LLFF [27], DTU [20], Mip-NeRF360 [3], and MVImgNet [45]. On LLFF, DTU, and Mip-NeRF360, we follow the experimental setup introduced in [47], using the same training/testing splits and downsampling rates. As suggested in [47], we mask the background when assessing the quality of novel view synthesis on the DTU dataset. Since LLFF, DTU, and Mip-NeRF360 offer a limited number of scenarios, we extend our experiments to a large-scale dataset, i.e., MVImgNet. We randomly sample 50 scenes from this dataset and resize the longest side of each image to 512. The Gaussian parameters for all evaluated methods are initialized based on the point clouds obtained

Method	PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow			AVGE \downarrow		
	3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view
Mip-NeRF [2]	16.11	22.91	24.88	0.401	0.756	0.826	0.460	0.213	0.160	0.215	0.090	0.066
DietNeRF [19]	14.94	21.75	24.28	0.370	0.717	0.801	0.496	0.248	0.183	0.240	0.105	0.073
RegNeRF [30]	19.08	23.10	24.86	0.587	0.760	0.820	0.336	0.206	0.161	0.149	0.086	0.067
FreeNeRF [44]	19.63	23.73	25.13	0.612	0.779	0.827	0.308	0.195	0.160	0.134	0.075	0.064
SparseNeRF [38]	19.86	-	-	0.624	-	-	0.328	-	-	0.127	-	-
3DGS [21]	19.22	23.80	25.44	0.649	0.814	0.860	0.229	0.125	0.096	0.120	0.066	0.051
DNGaussian [24]	19.12	-	-	0.591	-	-	0.294	-	-	0.132	-	-
FSGS [53]	20.43	24.09	25.31	0.682	0.823	0.860	0.248	0.145	0.122	0.104	0.066	0.054
CoR-GS [47]	20.45	24.49	26.06	0.712	0.837	0.874	0.196	0.115	0.089	0.101	0.060	0.046
SE-GS	20.79	24.78	26.36	0.724	0.839	0.878	0.183	0.110	0.084	0.100	0.059	0.045

Table 1. **Results on LLFF with 3, 6, and 9 training views.** We highlight the best, second-best, and third-best results in red, orange, and yellow, respectively.

Method	PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow			AVGE \downarrow		
	3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view
3DGS [21]	17.67	23.69	26.80	0.804	0.894	0.941	0.158	0.086	0.050	0.112	0.058	0.031
MVSplat [7]	17.33	16.34	16.22	0.598	0.596	0.587	0.279	0.296	0.304	-	-	-
DNGaussian [24]	18.57	-	-	0.776	-	-	0.178	-	-	0.110	-	-
FSGS [53]	17.84	23.68	26.17	0.822	0.905	0.941	0.161	0.096	0.064	0.110	0.056	0.036
CoR-GS [47]	18.65	24.39	27.38	0.835	0.910	0.950	0.140	0.074	0.045	0.099	0.049	0.028
SE-GS	19.24	25.28	28.08	0.857	0.924	0.958	0.132	0.073	0.043	0.089	0.045	0.026

Table 2. **Results on DTU with 3, 6, and 9 training views.** We use red, orange, and yellow to indicate the best, second-best, and third-best results, respectively. Object masks are used for all evaluated methods to remove background when conducting the evaluation.

from COLMAP [36] to ensure a fair comparison. Note that COLMAP fails on certain DTU scenes in the sparse-view setting. In these cases, we instead randomly initialize the point cloud. Please refer to the supplementary material for more details on the setup.

Implementation details. We train our SE-GS for 10,000 iterations on the LLFF, DTU, and MVImgNet datasets, and for 30,000 iterations on Mip-NeRF360. In our uncertainty-aware perturbation mechanism, we employ $M = 24$ image buffers with a buffer size of $S = 3$ and perturb the Δ -model every 500 iterations. The noise δ_i in Eq. 10 is sampled from a normal distribution with a mean of 0, and the standard deviation is adaptively adjusted based on the magnitude of the Gaussian parameters. More details are provided in the supplementary material.

4.2. Quantitative Results

We provide the quantitative results of the evaluated approaches on the LLFF, DTU, Mip-NeRF360, and MVImgNet datasets in Table 1, Table 2, Table 3, and Table 4, respectively. The best, second-best, and third-best results in each column are highlighted in red, orange, and yellow, respectively. For MVSplat [7], we employ the pretrained model released by the authors and evaluate it on DTU. Notably, the per-scene optimized 3DGS methods significantly outperform MVSplat, highlighting the advan-

tage of scene-level radiance field modeling for high-quality novel view synthesis. For these scene-level approaches, the performance of both NeRF and 3DGS methods consistently declines as the number of training views decreases, underscoring the challenge of few-shot novel view synthesis. In this context, our method surpasses all competitors across the four datasets. Note that in the sparse-view scenario, the available information in the training data is limited, making it more challenging to improve performance compared to the dense-view setting.

Specifically, on the LLFF and DTU datasets, we train the methods using 3, 6, and 9 views. Our method achieves an improvement of 0.34 in terms of PSNR with 3 views on LLFF. The improvement is larger on DTU, reaching a gain of 0.59. Since the images in the Mip-NeRF360 dataset are captured from diverse viewpoints, we sample more training views, i.e., 12 and 24, following the setting used in [47]. On this dataset, our SE-GS achieves the highest PSNR and SSIM, as well as the lowest LPIPS and AVGE, demonstrating superior NVS quality under the challenges of sparse views with large-scale viewpoint variations. Additionally, to further assess the robustness across diverse scenarios, we conduct an experiment on 50 scenes sampled from the MVImgNet dataset. As reported in Table 4, our method yields consistently better results than previous state-of-the-art approaches when trained with 5, 7, and 10 views. Our

Method	12-view				24-view			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	AVGE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	AVGE \downarrow
3DGS [21]	18.52	0.523	0.415	0.167	22.80	0.708	0.276	0.096
FSGS [53]	18.80	0.531	0.418	0.163	23.28	0.715	0.274	0.093
CoR-GS [47]	19.52	0.558	0.418	0.151	23.39	0.727	0.271	0.091
SE-GS	19.91	0.596	0.400	0.143	23.74	0.745	0.265	0.086

Table 3. **Results on Mip-NeRF360 with 12 and 24 training views.** The first, second, and third-best results are marked in red, orange, and yellow, respectively.

Method	PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow		
	5-view	7-view	10-view	5-view	7-view	10-view	5-view	7-view	10-view
3DGS [21]	26.48	29.06	31.82	0.819	0.878	0.924	0.301	0.228	0.184
FSGS [53]	26.88	29.12	31.68	0.848	0.886	0.922	0.358	0.292	0.250
CoR-GS [47]	27.62	29.77	32.09	0.842	0.889	0.928	0.302	0.235	0.190
SE-GS	27.88	30.19	32.68	0.857	0.902	0.937	0.277	0.214	0.177

Table 4. **Results on MVImgNet with 5, 7, and 10 training views.** Results colored in red, orange, and yellow denote the best, second-best, and third-best performances, respectively.

SE-GS also demonstrates better compatibility with relatively dense views. For instance, the PSNR improvement increases from 0.26 to 0.59 as the number of training views varies from 5 to 10.

The methods that leverage auxiliary data terms, such as DNGaussian [24] and FSGS [53], sometimes outperform 3DGS. However, in some cases, such as with 9 views on LLFF, the PSNR of FSGS becomes worse than 3DGS. Since the data acquired through off-the-shelf approaches might be unreliable, the quality of the NVS results is consequently affected. Compared with FSGS and DNGaussian, CoR-GS [47] exhibits better stability, outperforming 3DGS in more scenarios. This finding highlights the promising potential of regularization in the sparse-view setting. However, the improvement is still not consistent, as CoR-GS shows worse LPIPS than 3DGS when trained with 10 views on the MVImgNet dataset. By contrast, our SE-GS beats 3DGS in all cases, demonstrating superior stability. Our method is capable of avoiding specific local optima as the Σ -model aggregates information from diverse temporal samples in the presented self-ensembling mechanism.

4.3. Qualitative Results

Fig. 5 presents qualitative comparisons of 3DGS, CoR-GS, and our SE-GS, showing renderings from novel views. The models are trained on the DTU and MVImgNet datasets with 3 and 5 posed images, respectively. Our SE-GS produces fewer visual artifacts than the other methods, achieving better robustness against the challenge of few-shot training views. Moreover, our method captures finer details, par-

ticularly in areas with complex and repeated textures, such as the pineapple and watch strap. This demonstrates the effectiveness of our self-ensembling strategy in improving both the visual quality and stability of novel view neural rendering under sparse-view conditions.

4.4. Analysis

To shed more light on the effectiveness and efficiency of our SE-GS, we perform a comprehensive analysis of the introduced self-ensembling mechanism. The analysis starts by comparing our approach with CoR-GS trained with varying numbers of Gaussian models. The detailed results on LLFF with 3 training views are shown in Fig. 6. As shown by Fig. 6a, the PSNR of CoR-GS increases as more Gaussian models are trained. Since adding more Gaussian models enhances the regularization effect in CoR-GS, this observation highlights the potential of improving 3DGS in the sparse-view setting by strengthening the regularization process. The PSNR reaches a peak with 6 Gaussian models, indicating an upper bound of CoR-GS. Our method perturbs the Δ -model based on the uncertainties computed from dynamically updated image buffers, which results in more diverse samples of Gaussian models compared to CoR-GS. This leads to more effective regularization, as reflected by the better PSNR score.

Moreover, Fig. 6b shows the training speed measured as the number of training iterations completed per second. The speed of CoR-GS significantly drops as the number of Gaussian models increases. This limitation poses a challenge in scaling CoR-GS up to a large number of Gaus-



Figure 5. **Qualitative results.** The methods are trained on sparse views and the renderings of novel views are illustrated. The images are selected from the DTU and MVImgNet datasets.

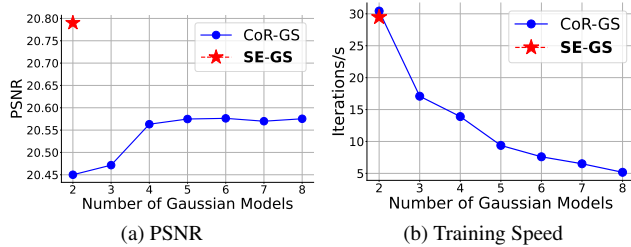


Figure 6. **Comparison with CoR-GS.** CoR-GS [47] is trained with different numbers of 3DGS models. The number varies from 2 to 8. The PSNR on LLFF with 3 training views is reported to evaluate the NVS quality. Iterations/s refers to the number of training iterations completed per second, indicating the training speed.

sian models. By contrast, in our pipeline, only two models are trained, i.e., the Δ -model and Σ -model. The temporal samples are generated through perturbations. We achieve a comparable training speed to CoR-GS with 2 models, showing that the perturbation process incurs a negligible additional cost. Note that only the Σ -model is retained for NVS during inference, so the testing speed of SE-GS is the same as that of the original 3DGS. Consequently, our method is

Method	Random	Gradient-Aware	Uncertainty-Aware
PSNR \uparrow	18.63	19.18	20.34

Table 5. **Effectiveness of the uncertainty-aware perturbations.** Random indicates that we add random perturbations to all Gaussians in the Δ -model. Gradient-Aware denotes an alternative in which we perturb the Gaussians based on gradients. All methods are trained on LLFF with 3 views for 2000 iterations and the PSNR is utilized as the metric.

capable of efficiently enhancing the performance of 3DGS with sparse training views.

Finally, we assess the effectiveness of our uncertainty-aware perturbation strategy, comparing it with two alternatives. Specifically, as introduced in Sec. 3, a straightforward method for perturbing the Δ -model is to add random noise to the parameters of all Gaussians in the model. We denote this approach as *Random*. Building upon this baseline, an alternative approach, referred to as *Gradient-Aware*, identifies unreliable Gaussians by analyzing the gradients. Specifically, we replace the indicator function in Eq. 11 with one based on gradient magnitudes. Gaussians with gradient magnitudes exceeding a threshold are perturbed. In practice, we adopt the same threshold as used in density control. We train these two alternatives and our approaches on LLFF with 3 views for 2000 iterations, which challenges the methods in terms of convergence speed. As listed in Table 5, *Random* performs the worst since it leads to many invalid Gaussian samples, negatively impacting convergence speed. Note that our uncertainty-aware strategy aggregates information from multiple training steps, while the gradient-aware method can only utilize the gradients at the current training step. The improvement in PSNR shown in Table 5 demonstrates the superior effectiveness of our approach.

5. Conclusion

In this paper, we have presented a new self-ensembling mechanism that enhances the novel view synthesis of 3DGS with sparse training images. We have tackled the challenge of overfitting by training an ensemble named Σ -model with the guidance of diverse temporal models sampled from the Gaussian parameter space. To obtain such samples, we have introduced an effective strategy that perturbs a Δ -model based on uncertainties computed from dynamically updated buffers of pseudo-view renderings. We have conducted experiments on LLFF, DTU, Mip-NeRF360, and MVImgNet, as well as a comprehensive analysis of our approach. The experimental results have demonstrated the effectiveness, stability, and efficiency of our SE-GS. In future work, we plan to incorporate an outlier identifier into our perturbation mechanism to reduce the impact of unreliable samples in regularization and

facilitate the self-ensembling process.

Acknowledgment. This work was funded in part by the Swiss National Science Foundation via the Sinergia grant CRSII5-180359 and the Swiss Innovation Agency (Innosuisse) via the BRIDGE Discovery grant 40B2-0_194729.

Self-Ensembling Gaussian Splatting for Few-Shot Novel View Synthesis

Supplementary Material

Details on Uncertainty-Aware Perturbation

Recall that in Eq. 11 of the main paper, we associate the reliability of Gaussians with the pixel-level uncertainty scores. In practice, for each uncertainty map $\hat{\mathbf{U}}$, we define τ as

$$\tau = \max(\hat{\mathbf{U}}_{\text{sorted}} \left[\left\lceil r * |\hat{\mathbf{U}}_{\text{sorted}}| \right\rceil \right], \theta), \quad (15)$$

where $\hat{\mathbf{U}}_{\text{sorted}}$ is obtained by sorting $\hat{\mathbf{U}}$ in descending order, $|\hat{\mathbf{U}}_{\text{sorted}}|$ denotes the number of uncertainty scores in the map, $\lceil \cdot \rceil$ indicates the ceiling function, $\hat{\mathbf{U}}_{\text{sorted}}[\cdot]$ accesses the value at the specified index. Here, $r = 0.05$ is a ratio scalar and $\theta = 0.01$ serves as a minimum tolerance. Subsequently, we identify the Gaussians that overlap with the pixels having uncertainty scores greater than τ . For each Gaussian in the Δ -model, we splat it to the 2D image planes corresponding to uncertainty maps. If any of these pixels are covered by the 2D splats of this Gaussian, the indicator function in Eq. 11 returns 1, marking the Gaussian as unreliable.

In Eq. 10, we sample the noise δ_t from a normal distribution. As an example, we elaborate on the perturbation process for the 3D position, which is similar for the other Gaussian parameters. In this context, we define δ_t as $\delta_t \in \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ where $\mathbf{I} \in \mathbb{R}^{3 \times 3}$ is an identity matrix. In practice, we adaptively adjust σ based on the magnitude of the parameters, as formulated by

$$\sigma = \omega \frac{1}{N} \sum_{i=1}^N \|\mu_i\|_1, \quad (16)$$

where ω is a scalar that controls the noise level and $\|\cdot\|_1$ denotes the L1 norm of the Gaussian parameter. The value of ω decays from 0.08 to 0.02, following a decay function introduced in [12].

Setup on DTU

As we mentioned in the main paper, COLMAP fails in some scenes on DTU when using sparse-view images. In our experiments, we randomly initialize the point cloud in these scenes. Specifically, we use random initialization for *scan8*, *scan40*, and *scan110* with 3 training views, and for *scan21* with 6 training views.

Ablation Studies

To shed more light on the impact of noise during the perturbation, we conduct comprehensive ablation studies on the perturbation interval and noise level.

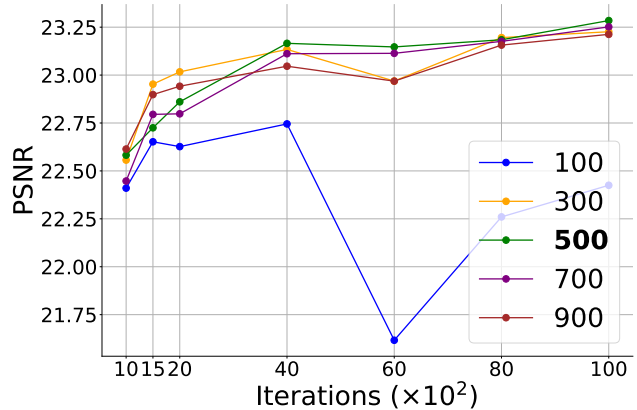


Figure 7. **Ablation study on the perturbation interval.** PSNR values on testing views throughout training are reported. The interval varies from 100 to 900.

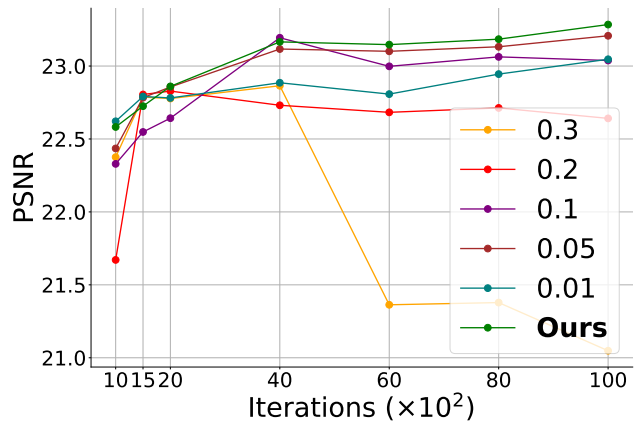


Figure 8. **Ablation study on the noise level.** We evaluate our method using different values of ω , each corresponding to a distinct noise level.

In our experiments, we perturb the Δ -model every 500 iterations by default. To evaluate the effect of different intervals, we test our method on the LLFF dataset with 3 training views, using intervals ranging from 100 to 900. As shown in Fig. 7, the curves represent the PSNR values obtained on testing views at different training iterations. Perturbing the Δ -model too frequently, e.g., with an interval of 100, significantly degrades performance, indicating a negative impact on our method. Conversely, large intervals, such as 900, reduce the number of temporal samples, thereby weakening the self-ensembling effect and leading to decreased performance compared to the default setting.

Moreover, we analyze the effect of the noise level by

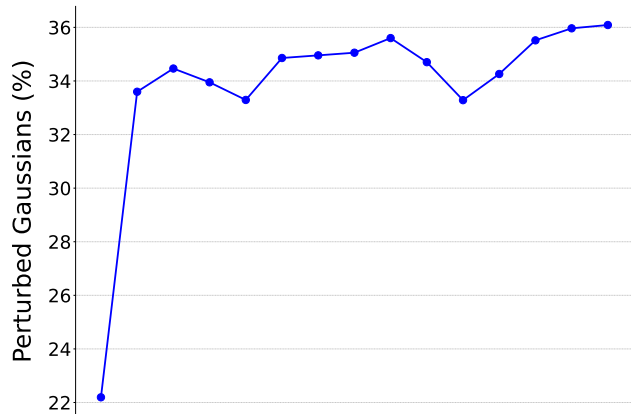


Figure 9. **Number of perturbed Gaussians throughout training.** The curve illustrates the percentage of perturbed Gaussians in the Δ -model at different training iterations.

conducting experiments with varying values of σ in Eq. 16. Specifically, we adjust the value of ω from 0.01 to 0.3, corresponding to increasing noise levels. We train our method with each specific noise level on the LLFF dataset using 3 training views. The results on testing views at different training iterations are shown in Fig. 8. Strong perturbations, such as those with ω values of 0.3 and 0.2, result in significantly worse performance compared to other settings. In these cases, the temporal models are sampled too far from the Δ -model in the Gaussian parameter space, reducing the consistency among temporal samples. The performance is also limited when a small ω , such as 0.01, is used. In this scenario, all temporal models closely resemble the Δ -model, diminishing the benefits of self-ensembling. In contrast to using a fixed ω , we adopt a decay function, where ω dynamically varies from 0.08 to 0.02 during training. This strategy achieves a good trade-off between the consistency and diversity of the temporal samples.

To better understand the presented perturbation mechanism, in Fig. 9, we illustrate the percentage of perturbed Gaussians in the Δ -model. Note that the threshold τ in Eq. 15 is defined adaptively. Therefore, the percentage dynamically varies throughout training. This dynamic behavior balances the diversity and reliability of the temporal models, thereby enhancing the self-ensembling process.

More Visualization Results

We provide more visualizations including pseudo-view renderings of the Σ -model and the Δ -model, uncertainty maps, and NVS results. For more details, please refer to the uploaded videos. Notably, the videos of pseudo-view renderings and uncertainty maps showcase the results across different training iterations.

References

- [1] Arsenii Ashukha, Alexander Lyzhov, Dmitry Molchanov, and Dmitry Vetrov. Pitfalls of in-domain uncertainty estimation and ensembling in deep learning. *arXiv preprint arXiv:2002.06470*, 2020. [3](#)
- [2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5855–5864, 2021. [6](#)
- [3] 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](#), [5](#)
- [4] Samuel R Buss and Jay P Fillmore. Spherical averages and applications to spherical splines and interpolation. *ACM Transactions on Graphics*, 20(2):95–126, 2001. [4](#)
- [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. [2](#)
- [6] Guikun Chen and Wenguan Wang. A survey on 3d gaussian splatting. *arXiv preprint arXiv:2401.03890*, 2024. [2](#)
- [7] 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](#), [6](#)
- [8] Hochul Cho, Jangyoon Kim, and Woontack Woo. Novel view synthesis with multiple 360 images for large-scale 6-dof virtual reality system. In *IEEE Conference on Virtual Reality and 3D User Interfaces*, pages 880–881. IEEE, 2019. [1](#)
- [9] Robert A Drebin, Loren Carpenter, and Pat Hanrahan. Volume rendering. *ACM Siggraph Computer Graphics*, 22(4):65–74, 1988. [2](#)
- [10] Zhiwen Fan, Wenyan Cong, Kairun Wen, Kevin Wang, Jian Zhang, Xinghao Ding, Danfei Xu, Boris Ivanovic, Marco Pavone, Georgios Pavlakos, et al. Instantsplat: Unbounded sparse-view pose-free gaussian splatting in 40 seconds. *arXiv preprint arXiv:2403.20309*, 2024. [2](#)
- [11] Geoffrey French, Michal Mackiewicz, and Mark Fisher. Self-ensembling for visual domain adaptation. *arXiv preprint arXiv:1706.05208*, 2017. [2](#), [3](#)
- [12] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qin-hong 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. [1](#)
- [13] Mudasir A Ganaie, Minghui Hu, Ashwani Kumar Malik, Muhammad Tanveer, and Ponnuthurai N Suganthan. Ensemble deep learning: A review. *Engineering Applications of Artificial Intelligence*, 115: 105151, 2022. [3](#)
- [14] Timur Garipov, Pavel Izmailov, Dmitrii Podoprikin, Dmitry P Vetrov, and Andrew G Wilson. Loss surfaces, mode connectivity, and fast ensembling of dnns. *Advances in Neural Information Processing Systems*, 31, 2018. [3](#)
- [15] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2495–2504, 2020. [3](#)
- [16] Antoine Guédon and Vincent Lepetit. Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5354–5363, 2024. [2](#)
- [17] Liang Han, Junsheng Zhou, Yu-Shen Liu, and Zhizhong Han. Binocular-guided 3d gaussian splatting with view consistency for sparse view synthesis. *arXiv preprint arXiv:2410.18822*, 2024. [3](#)
- [18] Tao Hu, Shu Liu, Yilun Chen, Tiancheng Shen, and Jiaya Jia. Efficientnerf efficient neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12902–12911, 2022. [2](#)
- [19] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, 2021. [6](#)
- [20] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanaes. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 406–413, 2014. [2](#), [5](#)
- [21] 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. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#)
- [22] Samuli Laine and Timo Aila. Temporal ensem-

- bling for semi-supervised learning. *arXiv preprint arXiv:1610.02242*, 2016. 2, 3
- [23] Stefan Lee, Senthil Purushwalkam, Michael Cogswell, David Crandall, and Dhruv Batra. Why m heads are better than one: Training a diverse ensemble of deep networks. *arXiv preprint arXiv:1511.06314*, 2015. 3
- [24] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. Dngaussian: Optimizing sparse-view 3d gaussian radiance fields with global-local depth normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20775–20785, 2024. 3, 5, 6, 7
- [25] Xiaomeng Li, Lequan Yu, Hao Chen, Chi-Wing Fu, Lei Xing, and Pheng-Ann Heng. Transformation-consistent self-ensembling model for semisupervised medical image segmentation. *IEEE Transactions on Neural Networks and Learning Systems*, 32(2):523–534, 2020. 2
- [26] Tianqi Liu, Guangcong Wang, Shoukang Hu, Liao Shen, Xinyi Ye, Yuhang Zang, Zhiguo Cao, Wei Li, and Ziwei Liu. Mvsgaussian: Fast generalizable gaussian splatting reconstruction from multi-view stereo. In *European Conference on Computer Vision*, pages 37–53. Springer, 2024. 2
- [27] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (ToG)*, 38(4):1–14, 2019. 2, 5
- [28] 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. 2
- [29] Duc Tam Nguyen, Chaithanya Kumar Mummadi, Thi Phuong Nhung Ngo, Thi Hoai Phuong Nguyen, Laura Beggel, and Thomas Brox. Self: Learning to filter noisy labels with self-ensembling. *arXiv preprint arXiv:1910.01842*, 2019. 2
- [30] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022. 6
- [31] Avinash Paliwal, Wei Ye, Jinhui Xiong, Dmytro Kotovenko, Rakesh Ranjan, Vikas Chandra, and Nima Khademi Kalantari. Coherentgs: Sparse novel view synthesis with coherent 3d gaussians. *arXiv preprint arXiv:2403.19495*, 2, 2024. 3
- [32] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. 2
- [33] Rahul Rahaman et al. Uncertainty quantification and deep ensembles. *Advances in Neural Information Processing Systems*, 34:20063–20075, 2021. 3
- [34] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12179–12188, 2021. 3
- [35] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 3
- [36] 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. 3, 6
- [37] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 519–528. IEEE, 2006. 3
- [38] Guangcong Wang, Zhaoxi Chen, Chen Change Loy, and Ziwei Liu. Sparsenerf: Distilling depth ranking for few-shot novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9065–9076, 2023. 6
- [39] 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. *arXiv preprint arXiv:2106.10689*, 2021. 2
- [40] Yunsong Wang, Tianxin Huang, Hanlin Chen, and Gim Hee Lee. Freesplat: Generalizable 3d gaussian splatting towards free-view synthesis of indoor scenes. *arXiv preprint arXiv:2405.17958*, 2024. 2
- [41] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf-: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021. 2
- [42] Yaniv Wolf, Amit Bracha, and Ron Kimmel. Surface reconstruction from gaussian splatting via novel stereo views. *arXiv preprint arXiv:2404.01810*, 2024. 2
- [43] Wangze Xu, Huachen Gao, Shihe Shen, Rui Peng, Jianbo Jiao, and Ronggang Wang. Mvpgs: Excavating multi-view priors for gaussian splatting from sparse

input views. *arXiv preprint arXiv:2409.14316*, 2024. 3

- [44] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering with free frequency regularization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8254–8263, 2023. 6
- [45] Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan, Chenming Zhu, Zhangyang Xiong, Tianyou Liang, et al. Mvimngnet: A large-scale dataset of multi-view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9150–9161, 2023. 2, 5
- [46] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19447–19456, 2024. 2
- [47] Jiawei Zhang, Jiahe Li, Xiaohan Yu, Lei Huang, Lin Gu, Jin Zheng, and Xiao Bai. Cor-gs: sparse-view 3d gaussian splatting via co-regularization. In *European Conference on Computer Vision*, pages 335–352. Springer, 2024. 3, 5, 6, 7, 8
- [48] Wu Zheng, Weiliang Tang, Li Jiang, and Chi-Wing Fu. Se-ssd: Self-ensembling single-stage object detector from point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14494–14503, 2021. 2
- [49] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *European Conference on Computer Vision*, pages 286–301. Springer, 2016. 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. 1
- [51] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5745–5753, 2019. 5
- [52] Zhi-Hua Zhou, Jianxin Wu, and Wei Tang. Ensembling neural networks: many could be better than all. *Artificial intelligence*, 137(1-2):239–263, 2002. 3
- [53] Zehao Zhu, Zhiwen Fan, Yifan Jiang, and Zhangyang Wang. Fsgs: Real-time few-shot view synthesis using gaussian splatting. In *European Conference on Computer Vision*, pages 145–163. Springer, 2024. 3, 5, 6, 7