# Modeling uncertainty for Gaussian Splatting

Luca Savant, *Student Member, IEEE*, Diego Valsesia, *Member, IEEE*, Enrico Magli, *Fellow, IEEE*

*Abstract*—We present Stochastic Gaussian Splatting (SGS): the first framework for uncertainty estimation using Gaussian Splatting (GS). GS recently advanced the novel-view synthesis field by achieving impressive reconstruction quality at a fraction of the computational cost of Neural Radiance Fields (NeRF). However, contrary to the latter, it still lacks the ability to provide information about the confidence associated with their outputs. To address this limitation, in this paper, we introduce a Variational Inference-based approach that seamlessly integrates uncertainty prediction into the common rendering pipeline of GS. Additionally, we introduce the Area Under Sparsification Error (AUSE) as a new term in the loss function, enabling optimization of uncertainty estimation alongside image reconstruction. Experimental results on the LLFF dataset demonstrate that our method outperforms existing approaches in terms of both image rendering quality and uncertainty estimation accuracy. Overall, our framework equips practitioners with valuable insights into the reliability of synthesized views, facilitating safer decision-making in real-world applications.

## I. INTRODUCTION

Novel-view synthesis, the task of generating images of a scene from viewpoints not observed during data collection, is a fundamental problem in computer vision with numerous applications, including virtual reality, augmented reality, and robotics. Traditionally, this task has been addressed using methods such as Structure from Motion [1], also called multi-view stereo, which rely on geometric reconstruction techniques. However, recent advances in deep learning, particularly with the introduction of Neural Radiance Fields (NeRF), have revolutionized the field by enabling high-fidelity synthesis of novel views directly from the underlying scene representation.

NeRF [2] has recently enjoyed great success by representing a scene as a continuous volumetric function that maps 3D spatial coordinates and viewing directions to radiance values. By learning this function from a set of posed images, NeRF can generate photorealistic images from novel viewpoints. However, while NeRF achieves impressive results, its computational complexity and memory requirements limit its practicality for real-time applications. This is the focus of the emerging Gaussian Splatting (GS) technique [3] which offers a more computationally efficient alternative to NeRF while maintaining high-quality novel-view synthesis. GS learns to approximate the radiance field by using a set of Gaussian kernels, enabling real-time rendering with competitive visual fidelity.

At the same time, research in novel view synthesis has started addressing the problem of estimating the epistemic uncertainty in order to understand the reliability of the generated views. Indeed, any practical downstream task that involves taking actions in the real world (such as robotics and autonomous systems) must consider not only the newly synthesized views but also their corresponding uncertainties, in order to potentially discard too uncertain yet promising actions. This scenario was first addressed in the seminal work of Shen et al. [4], where they proposed a deep architecture, based on NeRF, called S-NeRF to also estimate meaningful uncertainty maps for each generated view.

At the moment, GS lacks a mechanism for estimating uncertainty in the synthesized views. In this paper, we seek to address this limitation by proposing a novel framework for uncertainty estimation in GS. We extend the traditional deterministic GS framework to incorporate stochasticity, allowing us to predict uncertainty alongside synthesized views. Our approach leverages Variational Inference (VI) to learn the parameters of the GS radiance field in a Bayesian framework, enabling us to accurately estimate uncertainty without sacrificing computational efficiency.

We can summarize our novel contributions as follows:

- we introduce a novel framework for uncertainty estimation in GS, called Stochastic Gaussian Splatting (SGS), enabling real-time synthesis of high-quality images with accurate uncertainty predictions;
- we propose a VI-based approach to learn the parameters of the GS radiance field, allowing us to incorporate uncertainty prediction seamlessly into the rendering pipeline. Moreover, we innovate this learning process by augmenting Empirical Bayes with a loss function dependent on the area under the sparsification curve;
- we demonstrate the effectiveness of our approach through experiments on the challenging LLFF dataset, showing significant improvements in both rendering quality and uncertainty estimation metrics compared to state-of-the-art methods.

## II. BACKGROUND

### A. NeRF and Gaussian Splatting

In recent years, the Structure from Motion and novel view synthesis tasks have been revolutionized by novel techniques based on Neural Radiance Fields (NeRF), introduced in [2] and based on deep neural network architectures. Following the paradigm of implicit neural representation [5], the NeRF network learns a mapping $(\mathbf{x} \in \mathbb{R}^3, \mathbf{d} \in \mathbb{S}^2) \rightarrow (c(\mathbf{x}, \mathbf{d}), \sigma(\mathbf{x}))$ where $\mathbf{x}$ represents a 3D point, $\mathbf{d}$ a view direction, $c$ the color field, and $\sigma$ the density field. Using the Direct Volume Rendering Integral [6], [7], the physical radiance field $\mathbf{I}$ can

be obtained from $c$ and $\sigma$, hence an image can be formed. This learning process can be supervised with the pixels from a number of available views, so that the network correctly resynthesizes the original images. After successful training, the network is also able to coherently synthesize novel views.

In an ever-growing literature around NeRF, GS [3] distinguishes itself by achieving state-of-the-art visual quality while maintaining competitive training times and, importantly, allowing high-quality and real-time novel-view synthesis ($\geq$ 30fps at 1080p resolution). In order to achieve such high inference speed, GS replaces the deep neural network with an elliptical basis function approximation. By carefully choosing the basis functions as elliptical Gaussian kernels, the EWA Volume Splatting technique [8] can be used to achieve real-time view synthesis.

### B. Structure from Motion uncertainty estimation

An emerging research direction for the Structure from Motion field is the quantification of the uncertainty of the synthesized novel-views or of the 3D radiance field itself. This task was first addressed by [4], where it is cast as a Bayesian learning problem and solved with the Variational Inference (VI) framework [9], applied to NeRF. Subsequently, the same authors proposed an evolution of their method [10], which differs from the previous work by dropping the independence assumptions between the color and opacity fields, hence recovering a more complex stochastic dependency graph using the Conditional Normalizing Flows framework [11]. Following these works, [12] drops all independence assumptions by using a generative Flow-GAN model [13]. Another approach is addressed in [14] where the Laplace Approximation framework [9] is used. Finally, two similar works, [15] and [16], focus on developing methods that actively estimate high uncertainty in spatial areas that are not covered by the input views, using VI and Ensemble Learning frameworks, respectively. It is worth noting that all these methods rely on the use use of NeRF, while, to the best of our knowledge, there is no work currently addressing the GS framework.

### C. Radiance Field and Direct Volume Rendering

In this section, we recall the notation for volume rendering, which will be useful in the remainder of the paper. Following the optical derivation of [6], let's denote the physical light intensity field with $I$, which increases or decreases along a ray segment, parameterized by $t$, by interacting with particles in a particle-filled volume. These interactions are described using the density field $\sigma$ (also called the extinction coefficient) and the color field $c$ (also called the emission term), with the following Cauchy problem:

$$\begin{cases} \frac{dI}{dt} = -c(t)\sigma(t) + \sigma(t)I(t) & \forall t \in [0, T] \\ I(T) = I_T \end{cases} \quad (1)$$

whose solution value at $t = 0$ is:

$$I(0) = I_T e^{-\int_0^T \sigma(t)dt} + \int_0^T e^{-\int_0^T \sigma(z)dz} c(t)\sigma(t)dt \quad (2)$$

Reparametrizing the segment with $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$, setting $I_T = 0$ (in the case of a dark background), and assuming an isotropic density field, the rendering integral is obtained:

$$I(\mathbf{o}, \mathbf{d}) = \int_0^T e^{-\int_0^t \sigma(\mathbf{r}(s))ds} c(\mathbf{r}(t), \mathbf{d}) \sigma(\mathbf{r}(t)) dt \quad (3)$$

From a given set of posed views, a set of corresponding posed pixels and colors can be extracted as $\mathcal{D} = \{((\mathbf{o}_i, \mathbf{d}_i), y_i)\}_i$, where $i$ is an index that runs over all the pixels of all the images in the training set, $\mathbf{o}_i$ is the camera center of the $i$-th pixel, $\mathbf{d}_i$ is the 3d spatial direction of the $i$-th pixel from its camera center and $y_i$ is the color of the $i$-th pixel, so that a reconstruction loss can be used to regress the color and density fields:

$$\mathcal{L}_{\text{rec}} = \sum_i \ell(y_i, I(\mathbf{o}_i, \mathbf{d}_i)) \quad (4)$$

To synthesize a novel view, the integral in Eq. (3) must be evaluated by placing $\mathbf{o}$ at the camera origin and letting $\mathbf{d}$ vary for each pixel in the new image plane.

### III. METHOD

In this section, we present the proposed method, called Stochastic Gaussian Splatting (SGS), to enable uncertainty quantification in the Gaussian Splatting framework.

### A. Stochastic Gaussian Splatting

The main difference between NeRF and GS lies in the fact that in the original NeRF formulation [2], the color and density fields are regressed with a multilayer perceptron neural network. Instead, in GS [3], the two fields are learned with an elliptical basis function approximation. Let $\mathbf{x}$ be any point in 3D space, and let $\mathbf{d}$ identify a view direction, then the color and density fields are obtained as:

$$\begin{cases} \sigma(\mathbf{x}) \approx \sum_{k=1}^K \alpha_k \phi(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ c(\mathbf{x}, \mathbf{d}) \approx \sum_{k=1}^K c_k(\mathbf{d}) \phi(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{cases} ; \quad (5)$$

where:
- $\phi$ is the elliptical 3D Gaussian kernel, with learnable parameters $\boldsymbol{\mu}_k \in \mathbb{R}^3$ and $\boldsymbol{\Sigma}_k \in \mathbb{R}^{3 \times 3}$ to control the shape and position of the $k$-th kernel:

$$\phi(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right); \quad (6)$$

- $\alpha_k \in [0, 1]$ is a learnable scalar to control the opacity of the $k$-th kernel;
- $c_k(\mathbf{d})$ is a learned linear combination of spherical harmonics $Y_{lm}$:

$$c_k(\mathbf{d}) = \sum_{l=0}^L \sum_{m=0}^l c_{klm} Y_{lm}(\mathbf{d}), \quad (7)$$

where $l, m$ are degree and order of the used spherical harmonics.

In GS, the learnable parameters for each Gaussian kernel are $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, c_{klm}$, and $\alpha_k$.
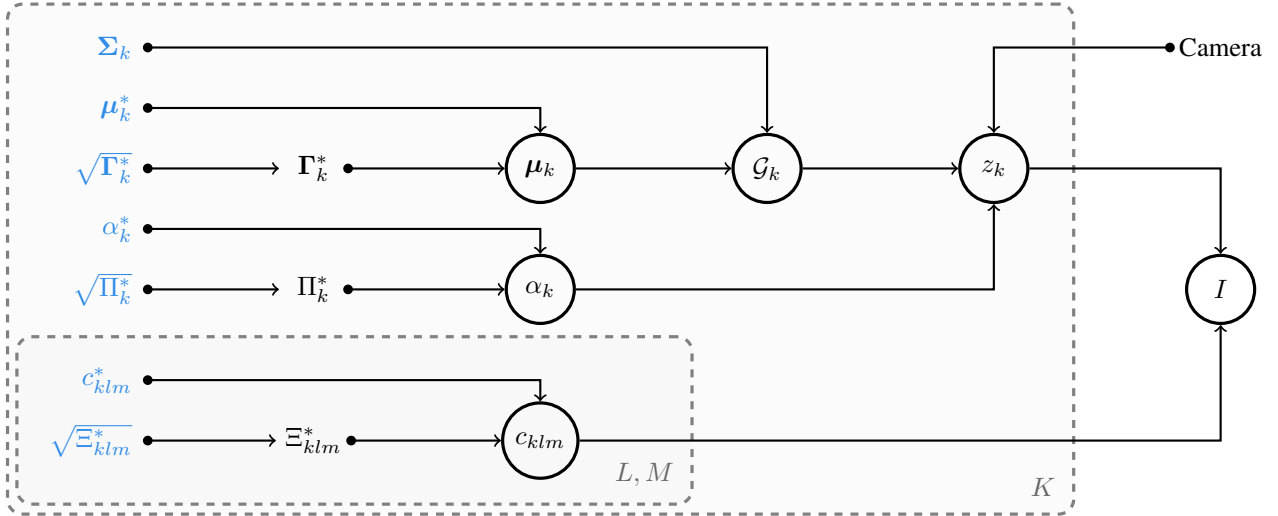
Fig. 1. Bayesian Network Graphical Model of SGS. Learnable variables are depicted in blue, while stochastic variables are circled. The "Camera" node represents both the spatial coordinate of the pixel $(\mathbf{o}, \mathbf{d})$ and the corresponding camera intrinsic and extrinsic parameters. Gray dashed rectangles are used for the plate notation, i.e. variables repetitions.

The specific definition of $\phi$ enables the EWA Volume Splatting technique [8], so that the integral in Eq. (3) can be reduced to the simpler alpha blending technique:

$$I(\mathbf{o}, \mathbf{d}) = \sum_{k=1}^{K} \alpha_k e^{z_k(\mathbf{o},\mathbf{d})} c_k(\mathbf{d}) \prod_{j=1}^{k-1} (1 - \alpha_j e^{z_j(\mathbf{o},\mathbf{d})}) \quad (8)$$

where $z_k(\mathbf{o}, \mathbf{d})$ is the splatting coefficient [8].

Aligning with previous works in the field, the proposed SGS method predicts an uncertainty value for each pixel in the novel-synthesized view. It is natural to define this uncertainty as the standard deviation of the predicted pixel color, requiring the predicted colors to be random variables. Hence, we need to inject stochasticity into the otherwise deterministic process of Eq. (8). We propose to use a Monte Carlo method to approximate the variance of pixel colors, so that the expression (8) should be evaluated multiple times, each time by sampling some random variables.

Following the Variational Inference framework, we directly expand upon GS by imposing a prior distribution on each of the following parameters:

- $\boldsymbol{\mu}_k \sim \mathcal{N}(\boldsymbol{\mu}_k^*, \boldsymbol{\Gamma}_k^*)$
- $\ln \frac{\alpha_k}{1-\alpha_k} \sim \mathcal{N}(\alpha_k^*, \Pi_k^*)$
- $c_{klm} \sim \mathcal{N}(c_{klm}^*, \Xi_{klm}^*)$

so that the original GS parameters are no longer learned but are sampled from the above distributions, whose parameters are the new learning variables.

Thanks to the reparameterization trick [17], gradients can flow from the loss in Eq. (4), which is a function of the samples, to the new distribution parameters $\boldsymbol{\mu}_k^*, \boldsymbol{\Gamma}_k^*, \alpha_k^*, \Pi_k^*, c_{klm}^*$ and $\Xi_{klm}^*$, enabling learning with standard backpropagation techniques.

### B. Learning with Variational Inference

The Variational framework introduced in [9] and used in [4] and [10] is required for optimization, since direct optimization

of just the loss would otherwise incur in underestimation of the pixel variance as the model could reduce to a deterministic one.

The VI framework stems from an approximation of Bayes' Theorem. Let $\mathcal{G}$ represent a Gaussian Splatting Radiance Field, composed of $K$ Gaussian kernels $\mathcal{G}_k$, and $\mathcal{D}$ represent the pixels dataset in Eq. (4). Then the Bayes' Theorem reads as:

$$\mathbb{P}(\mathcal{G}|\mathcal{D}) = \frac{\mathbb{P}(\mathcal{D}|\mathcal{G})\,\mathbb{P}(\mathcal{G})}{\mathbb{P}(\mathcal{D})}. \quad (9)$$

Since this is generally intractable, the VI framework prescribes to approximate the true posterior $\mathbb{P}(\mathcal{G}|\mathcal{D})$ by introducing a parametric distribution $q_\theta$ over all GS radiance fields $\mathcal{G}$ and to learn these parameters $\theta$ in order to minimize the Kullback-Leibler (KL) divergence between the approximate posterior and the true one:

$$\min_\theta \mathrm{KL}(q_\theta(\mathcal{G}) \,||\, \mathbb{P}(\mathcal{G}|\mathcal{D})) \quad (10)$$

Now, this problem is further manipulated to get a tractable expression. Let us start with the following manipulation, where only properties of the logarithm and linearity on the expectation are used:

$$\mathrm{KL}(q_\theta(\mathcal{G}) \,||\, \mathbb{P}(\mathcal{G}|\mathcal{D})) =$$
$$= \mathbb{E}_{q_\theta}\left[\log \frac{q_\theta(\mathcal{G})}{\mathbb{P}(\mathcal{G}|\mathcal{D})}\right] = \mathbb{E}_{q_\theta}\left[\log \frac{q_\theta(\mathcal{G})\mathbb{P}(\mathcal{D})}{\mathbb{P}(\mathcal{D}|\mathcal{G})\,\mathbb{P}(\mathcal{G})}\right]$$
$$= -\mathbb{E}_{q_\theta}[\log \mathbb{P}(\mathcal{D}|\mathcal{G})] + \mathbb{E}_{q_\theta}\left[\log \frac{q_\theta(\mathcal{G})}{\mathbb{P}(\mathcal{G})}\right] + \mathbb{E}_{q_\theta}[\log \mathbb{P}(\mathcal{D})]$$
$$= -\mathbb{E}_{q_\theta}[\log \mathbb{P}(\mathcal{D}|\mathcal{G})] + \mathrm{KL}(q_\theta(\mathcal{G}) \,||\, \mathbb{P}(\mathcal{G})) + \log \mathbb{P}(\mathcal{D})$$
$$(11)$$

Now, remembering that $\mathcal{D} = \{((\mathbf{o}_i, \mathbf{d}_i), y_i)\}_i$, and that the dataset samples are independently sampled, the first term in the last expression is equivalent to:

$$\mathbb{E}_{q_\theta}[\log \mathbb{P}(\mathcal{D}|\mathcal{G})] = \tag{12}$$

$$= \mathbb{E}_{q_\theta}\left[\log \prod_i \mathbb{P}(\mathcal{D}_i|\mathcal{G})\right] \tag{13}$$

$$= \sum_i \mathbb{E}_{q_\theta}[\log \mathbb{P}(\mathcal{D}_i|\mathcal{G})] \tag{14}$$

$$= \sum_i \mathbb{E}_{q_\theta}[\log \mathbb{P}(((\mathbf{o}_i, \mathbf{d}_i), y_i)|\mathcal{G})] \tag{15}$$

$$= \sum_i \mathbb{E}_{q_\theta}[\log \mathbb{P}(y_i|(\mathbf{o}_i, \mathbf{d}_i), \mathcal{G})\,\mathbb{P}((\mathbf{o}_i, \mathbf{d}_i)|\mathcal{G})] \tag{16}$$

$$= \sum_i \mathbb{E}_{q_\theta}[\log \mathbb{P}(y_i|(\mathbf{o}_i, \mathbf{d}_i), \mathcal{G})\,\mathbb{P}(\mathbf{o}_i, \mathbf{d}_i)] \tag{17}$$

$$= \sum_i \mathbb{E}_{q_\theta}[\log \mathbb{P}(y_i|(\mathbf{o}_i, \mathbf{d}_i), \mathcal{G})] + \mathbb{E}_{q_\theta}[\log \mathbb{P}(\mathbf{o}_i, \mathbf{d}_i)] \tag{18}$$

$$= \sum_i \mathbb{E}_{q_\theta}[\log \mathbb{P}(y_i|(\mathbf{o}_i, \mathbf{d}_i), \mathcal{G})] + \log \mathbb{P}(\mathbf{o}_i, \mathbf{d}_i) \tag{19}$$

Plugging (12) into (11), we get:

$$\begin{aligned}
\mathrm{KL}(q_\theta(\mathcal{G}) \,||\, \mathbb{P}(\mathcal{G}|\mathcal{D})) &= \\
&= -\sum_i \mathbb{E}_{q_\theta}[\log \mathbb{P}(y_i|(\mathbf{o}_i, \mathbf{d}_i), \mathcal{G})]+ \\
&\quad -\sum_i \log \mathbb{P}(\mathbf{o}_i, \mathbf{d}_i)+ \\
&\quad + \mathrm{KL}(q_\theta(\mathcal{G}) \,||\, \mathbb{P}(\mathcal{G})) + \log \mathbb{P}(\mathcal{D})
\end{aligned} \tag{20}$$

As the optimization variable is the parameters $\theta$, all the terms that do not depend on $\theta$ in the latter equation can be discarded. Finally, we have the optimization problem:

$$\min_\theta -\sum_i \mathbb{E}_{q_\theta}[\log \mathbb{P}(y_i|(\mathbf{o}_i, \mathbf{d}_i), \mathcal{G})] + \mathrm{KL}(q_\theta(\mathcal{G}) \,||\, \mathbb{P}(\mathcal{G})) \tag{21}$$

The first term in the loss function of problem (21) is the expected negative log-likelihood. This term forces $\theta$ to maximize the expected log-likelihood by matching the observations in the dataset $\mathcal{D}$. So, in spirit, it replaces the standard loss (4). It is estimated with the Monte Carlo method and, for simplicity, by defining the conditional probability distribution $\mathbb{P}(y_i|(\mathbf{o}_i, \mathbf{d}_i), \mathcal{G})$ to be a normal distribution that is pixel-wise independent, i.e.:

$$\log \mathbb{P}(y_i|(\mathbf{o}_i, \mathbf{d}_i), \mathcal{G}) \propto (y_i - I(\mathbf{o}_i, \mathbf{d}_i))^2 \tag{22}$$

Note that this requirement is not too strict as the independence is required in the conditional probability distribution and not for the unconditional probability distribution.

The second term in Eq. (21), instead, limits $\theta$ from moving too far away from the prior distribution of the GS radiance field $\mathcal{G}$. For it to be efficiently tractable, we suppose independence among the $K$ Gaussian kernels, for both the prior and the approximate posterior distributions, i.e.:

$$q_\theta(\mathcal{G}) = \prod_{k=1}^K q_{\theta_k}(\mathcal{G}_k) \tag{23}$$

$$\mathbb{P}(\mathcal{G}) = \prod_{k=1}^K \mathbb{P}(\mathcal{G}_k) \tag{24}$$

This independence assumption between Gaussian kernels is more general with respect to previous work. For instance, Shen et al. [4] prescribed independence between the values of the opacity fields for every pair of points in the radiance field, even if they are very close to each other. Meanwhile, in the proposed method, the independence is prescribed at the Gaussian kernel level and not at the infinitesimal 3D point level.

Thanks to the assumptions in Eqs. (23) and (24), the second term in Eq. (21) can be expanded to become:

$$\mathrm{KL}(q_\theta(\mathcal{G}) \,||\, \mathbb{P}(\mathcal{G})) = \sum_{k=1}^K \mathrm{KL}(q_{\theta_k}(\mathcal{G}_k) \,||\, \mathbb{P}(\mathcal{G}_k)) \tag{25}$$

If we suppose that the prior distributions are themselves multivariate Normal, each KL term in the right-hand side of Eq. (25) has the following general closed form expression:

$$\begin{aligned}
\mathrm{KL}(\mathcal{N}_{\mu_0, \Sigma_0} \,||\, \mathcal{N}_{\mu_1, \Sigma_1}) = &-\frac{3}{2} + \frac{1}{2}\mathrm{Tr}\left(\Sigma_1^{-1}\Sigma_0\right) \\
&+ \frac{1}{2}(\mu_1 - \mu_0)^T \Sigma_1^{-1}(\mu_1 - \mu_0) \\
&+ \frac{1}{2}\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right)
\end{aligned}$$

Hence, we define the total KL contribution to the training loss $\mathcal{L}_{\mathrm{KL}}$ as the sum over all Gaussians and over all the learnable parameters of the KL divergence between the prior (hat variables) and the posterior (starred variables) of that parameter:

$$\begin{aligned}
\mathcal{L}_{\mathrm{KL}} = \sum_{k=1}^K \Big[ &\mathrm{KL}(\mathcal{N}_{\hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Gamma}}_k} \,||\, \mathcal{N}_{\boldsymbol{\mu}_k^*, \boldsymbol{\Gamma}_k^*}) \\
&+ \mathrm{KL}(\mathcal{N}_{\hat{\alpha}_k, \hat{\Pi}_k} \,||\, \mathcal{N}_{\alpha_k^*, \Pi_k^*}) \\
&+ \mathrm{KL}(\mathcal{N}_{\hat{c}_{klm}, \hat{\Xi}_{klm}} \,||\, \mathcal{N}_{c_{klm}^*, \Xi_{klm}^*}) \Big]
\end{aligned} \tag{26}$$

### C. Learning with AUSE

To assess the accuracy of uncertainty estimation, a quantitative approach involves examining its correlation with the true error map using the Sparsification Curve [18]. First, the predicted values are sorted based on decreasing predicted uncertainty and then progressively removed, starting from those with high predicted uncertainty. By keeping track of a quality metric applied to the remaining values, the Sparsification Curve is generated. The area under this curve is called Area Under Sparsification Curve (AUSC). The Area Under the Sparsification Error (AUSE) metric is defined as the difference between the AUSC of the method and the AUSC

of the oracle, i.e., the curve obtained by sorting the predicted values according to the true error.

If the uncertainty prediction was random, the percolation process would also be random, resulting in a flat curve and a high AUSE. Otherwise, if the uncertainty prediction was positively correlated with the prediction error, improvements in the tracked quality metric would be observed. As the GS technique has significantly lower memory requirements compared to the NeRF used in previous works, it is capable of sampling the whole view multiple times in a single forward pass. This enables us to directly compute the AUSE metric applied to all the pixels in a view, taking the standard deviation of the samples of each pixel as the uncertainty map.

As we will show in the experiments section, using only the standard VI framework is suboptimal, as VI would not fully exploit the fact that the GS technique is orders of magnitude faster than neural-network-based NeRFs. Hence, in this work, we exploit the efficiency and low memory requirements of GS by augmenting the VI loss of Eq. (10) with the AUSE metric.

### D. End-to-end SGS Training

We now have all the necessary ingredients to define the overall loss used to train the proposed SGS method. In particular, the overall loss is the following combination:

$$\mathcal{L} = \mathcal{L}_{\text{rec}} + \lambda_{\text{SSIM}}\mathcal{L}_{\text{ssim}} + \lambda_{\text{KL}}\mathcal{L}_{\text{KL}} + \lambda_{\text{AUSE}}\mathcal{L}_{\text{AUSE}} \qquad (27)$$

where $\mathcal{L}_{\text{rec}}$ is defined as in (4) using the $\ell_1$ norm, $\mathcal{L}_{\text{ssim}}$ is the SSIM perceptual loss, $\mathcal{L}_{\text{KL}}$ is the Kullback–Leibler divergence with the prior from Eq. (26), and $\mathcal{L}_{\text{AUSE}}$ is the loss induced by the AUSE RMSE metric. The $\ell_1$, SSIM and AUSE losses augment the conventional KL loss in order to more explicitly enforce the training tradeoff between distortion, perceptual quality and uncertainty estimation.

Finally, one more aspect in which SGS training differs from previous works is the approach to learning the distribution of the priors. In previous works, stochasticity was introduced in the weights of neural networks, which are typically randomly initialized [19]. However, in GS, the parameters have a more direct physical meaning in the 3D space. For example, minimizing the KL-divergence in GS would tend to fix the center of a Gaussian kernel fixed in a randomly initialized position in 3D space. Instead, we tackle this convergence issue, by taking inspiration from Empirical Bayes. We thus first learn an informative prior with some iterations of classic GS, which serves as an initialization before switching to the SGS formulation.

## IV. EXPERIMENTAL RESULTS

### A. Experimental Setting

This work addresses the task of synthesizing novel views and associated uncertainty maps from multiple posed views of a static scene. These views consist of RGB photos captured from various camera positions and orientations. The spatial positions of the cameras are referred to as extrinsic parameters, while the intrinsic parameters include the camera's focal length and central point. Similar to prior research, we assume that both extrinsic and intrinsic camera parameters are known

for all input views. The objective is to perform standard novel view synthesis using the provided views while also generating an uncertainty map for each synthesized view. Consequently, for every pixel in the novel view, both its color and uncertainty are predicted. It is crucial that the predicted uncertainty correlates with the true error.

We remark that this is the first paper addressing uncertainty estimation for the GS setting, while current literature focuses on NeRF models. This makes direct comparisons challenging, as both absolute image quality as well as metrics for measuring the effectiveness of uncertainty estimation need to be reported. In particular, we remark that the AUSE metric benchmarks uncertainty estimation against the oracle method, so it is relatively insensitive to the absolute image quality generated by each method.

We compare our proposed method with the current literature on NeRF: the state-of-the-art CF-NeRF [10], the pioneering work of S-NeRF [4], and also with NeRF-W [20], Deep-Ensembles (D.E.) [21] and MC-Dropout [22], as done in previous works. We remark that some methods [4], [10] use extra information in the form of depth maps while we do not, resulting in a setting that is slightly unfair towards SGS. Nevertheless, we report improvements over such methods.

As common practice, all the experiments are conducted on the LLFF dataset from the original NeRF paper [2], which is composed of eight scenes (*fern*, *flower*, *fortress*, *horns*, *leaves*, *orchids*, *room* and *trex*), using the standard train-test split (e.g., the test split for *fern* is {*IMG4026*, *IMG4034*, *IMG4042*}). All the experiments are performed at $\frac{1}{8}$ of the original resolution, so that all synthesized images and uncertainty maps are composed by $504 \times 378$ pixels.

The hyperparameters in the final loss function are: $\lambda_{\text{KL}} = 10^{-3}$, $\lambda_{\text{AUSE}} = 5$, $\lambda_{\text{ssim}} = 0.2$. For the first 2500 iterations, we use the standard GS method with the following hyperparameters: the highest spherical harmonics degree is 1, the densification step is applied until iteration 1000, and the learning rate of the gaussians centers is fixed to $10^{-2}$.

At iteration 2500, the current learned GS is fixed and taken as the prior, and the Bayesian regime is introduced. All the prior covariance matrices are initialized as $10^{-2}\mathbb{I}$, where $\mathbb{I}$ is the identity matrix of the correct dimension. The learning rate for all the posterior learnable parameters is set to $10^{-4}$. Then we continue the training until iteration 10000. During both training and testing, Monte Carlo sampling from the posterior is performed for 8 times.

### B. Main Results

Table I reports our results on the LLFF dataset, evaluating the quality of the rendered images, as well as the reliability of the associated uncertainty maps. The rendered images are quantitatively evaluated with three metrics: PSNR as a distortion metric and LPIPS and SSIM as perceptual metrics. It can be noticed that our method improves by a large margin all these metrics, hence our method has a much lower test prediction bias with respect to previous work. The uncertainty maps are quantitatively evaluated with two metrics: AUSE RMSE and AUSE MAE. Both metrics are obtained as the

TABLE I
RESULTS ON THE LLFF DATASET. * DENOTES EXTRA DEPTH INFORMATION.

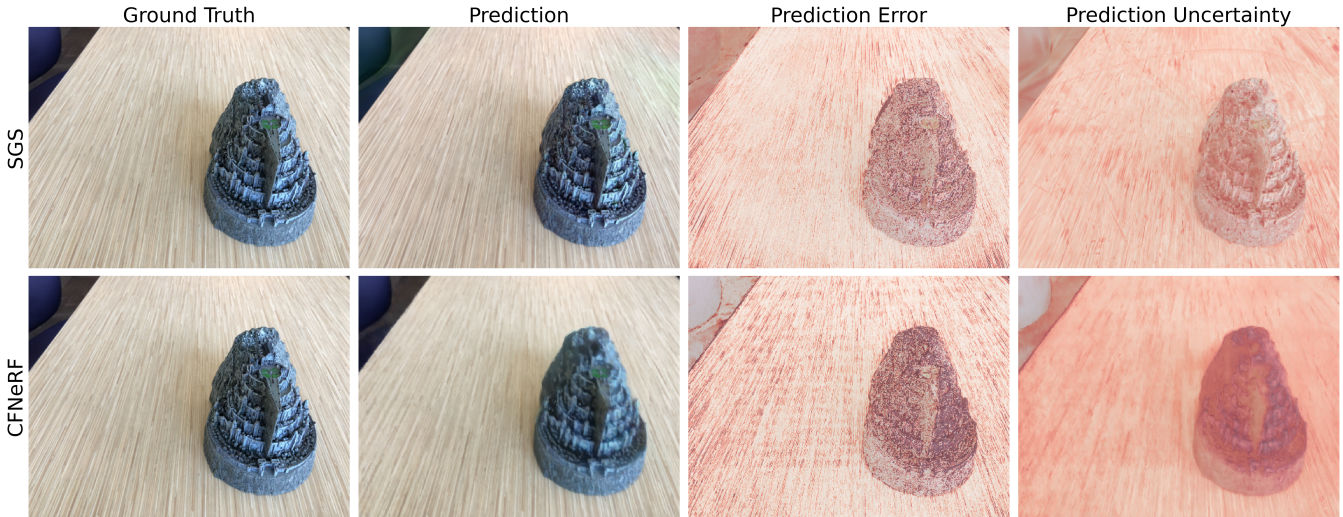| | Rendering Metrics | | | Uncertainty Metrics | |
|---|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | AUSE RMSE ↓ | AUSE MAE ↓ |
| D.E. (Lakshminarayanan, Pritzel, and Blundell 2017) | <u>22.32</u> | 0.788 | 0.236 | 0.0254 | 0.0122 |
| Drop. (Gal and Ghahramani 2016) | 21.90 | 0.758 | 0.248 | 0.0316 | 0.0162 |
| NeRF-W (Martin-Brualla et al. 2021) | 20.19 | 0.706 | 0.291 | 0.0268 | 0.0113 |
| S-NeRF* (Shen et al. 2021) | 20.27 | 0.738 | 0.229 | 0.0248 | 0.0101 |
| CF-NeRF* (Shen et al. 2022) | 21.96 | <u>0.790</u> | <u>0.201</u> | <u>0.0177</u> | **0.0078** |
| SGS (Ours) | **24.20** | **0.842** | **0.121** | **0.0147** | <u>0.0092</u> |



Fig. 2. A qualitative example of our method SGS with CF-NeRF [10]. The last column is a visualization of the predicted uncertainty map.

TABLE II
AUSE LOSS TERM ablation.

| | PSNR (dB) ↑ | SSIM ↑ | LPIPS ↓ | AUSE RMSE ↓ |
|---|---|---|---|---|
| No AUSE Loss | **26.65** | **0.869** | **0.082** | 0.0291 |
| **SGS (Ours)** | 24.20 | 0.842 | 0.121 | **0.0147** |

area under the sparsification curve, but considering as the error metric the Root Mean Square Error and the Mean Absolute Error, respectively. As shown in Table I, our method improves the AUSE RMSE metric, while keeping an AUSE MAE metric comparable with the state of the art, which however also exploits depth information. Fig. 2 shows a qualitative result for a novel view generated by SGS and CF-NeRF together with the predicted uncertainty maps. We can notice that SGS is capable of producing a sharp view as well as prediction uncertainty which correlates well with the true rendering error.

### C. AUSE Loss Ablation

Table II compares our SGS method with an ablated version, where the proposed AUSE loss term $\mathcal{L}_{\text{AUSE}}$ in equation (27) is removed in order to verify its effectiveness. As reported in the table, the removal of this loss term improves the photometric reconstruction (measured by the three quality metrics: PSNR, SSIM, and LPIPS), while deteriorating the model's ability to predict accurate uncertainty maps. Hence, this ablation study proves that one of our key contribution, that is to incorporate the AUSE loss term, improves the quality of the predicted uncertainty maps. Moreover, the hyperparameter $\lambda_{\text{AUSE}}$ provides a natural way to control the impact of this loss term, so that an application-specific trade-off between reconstruction quality and accurate uncertainty prediction can be found for downstream tasks.

## V. CONCLUSIONS

In this paper, we proposed a novel approach for uncertainty estimation in GS-based novel-view synthesis tasks. Leveraging the efficiency and real-time capabilities of GS, we introduced a stochastic extension to the traditional deterministic GS framework. Our method incorporates uncertainty prediction through a Bayesian framework, specifically using VI to learn the parameters of the GS radiance field. The training is further augmented with direct optimization of the AUSE metric to control the tradeoff between reconstruction quality and accuracy of uncertainty estimation.

Experimental results on the LLFF dataset demonstrated the effectiveness of our approach. We outperformed state-of-the-art methods in terms of rendering quality metrics, while also improving upon uncertainty estimation metrics. Notably, our work advances the state of the art by being the first to introduce uncertainty estimation for GS-based novel-view synthesis tasks.

# REFERENCES

[1] S. Ullman, "The interpretation of structure from motion," *Proceedings of the Royal Society of London. Series B. Biological Sciences*, vol. 203, no. 1153, pp. 405–426, 1979.

[2] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[3] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics (ToG)*, vol. 42, no. 4, pp. 1–14, 2023.

[4] J. Shen, A. Ruiz, A. Agudo, and F. Moreno-Noguer, "Stochastic neural radiance fields: Quantifying uncertainty in implicit 3d representations," in *2021 International Conference on 3D Vision (3DV)*. IEEE, 2021, pp. 972–981.

[5] Z. Chen and H. Zhang, "Learning implicit fields for generative shape modeling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5939–5948.

[6] N. Max, "Optical models for direct volume rendering," *IEEE Transactions on Visualization and Computer Graphics*, vol. 1, no. 2, pp. 99–108, 1995.

[7] S. Chandrasekhar, *Radiative transfer*. Courier Corporation, 2013.

[8] M. Zwicker, H. Pfister, J. Van Baar, and M. Gross, "Ewa volume splatting," in *Visualization, 2001. VIS 01. Proceedings*. IEEE, 2001, pp. 29–538.

[9] L. V. Jospin, H. Laga, F. Boussaid, W. Buntine, and M. Bennamoun, "Hands-on bayesian neural networks—a tutorial for deep learning users," *IEEE Computational Intelligence Magazine*, vol. 17, no. 2, pp. 29–48, 2022.

[10] J. Shen, A. Agudo, F. Moreno-Noguer, and A. Ruiz, "Conditional-flow nerf: Accurate 3d modelling with reliable uncertainty quantification," in *European Conference on Computer Vision*. Springer, 2022, pp. 540–557.

[11] G. Yang, X. Huang, Z. Hao, M.-Y. Liu, S. Belongie, and B. Hariharan, "Pointflow: 3d point cloud generation with continuous normalizing flows," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 4541–4550.

[12] S. Wei, J. Zhang, Y. Wang, F. Xiang, H. Su, and H. Wang, "Fg-nerf: Flow-gan based probabilistic neural radiance field for independence-assumption-free uncertainty estimation," *arXiv preprint arXiv:2309.16364*, 2023.

[13] A. Grover, M. Dhar, and S. Ermon, "Flow-gan: Combining maximum likelihood and adversarial learning in generative models," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, 04 2018.

[14] L. Goli, C. Reading, S. Selllán, A. Jacobson, and A. Tagliasacchi, "Bayes' rays: Uncertainty quantification for neural radiance fields," *arXiv preprint arXiv:2309.03185*, 2023.

[15] J. Shen, R. Ren, A. Ruiz, and F. Moreno-Noguer, "Estimating 3d uncertainty field: Quantifying uncertainty for neural radiance fields," *arXiv preprint arXiv:2311.01815*, 2023.

[16] N. Sünderhauf, J. Abou-Chakra, and D. Miller, "Density-aware nerf ensembles: Quantifying predictive uncertainty in neural radiance fields," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 9370–9376.

[17] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[18] F. K. Gustafsson, M. Danelljan, and T. B. Schon, "Evaluating scalable bayesian deep learning methods for robust computer vision," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 318–319.

[19] M. V. Narkhede, P. P. Bartakke, and M. S. Sutaone, "A review on weight initialization strategies for neural networks," *Artificial intelligence review*, vol. 55, no. 1, pp. 291–322, 2022.

[20] R. Martin-Brualla, N. Radwan, M. S. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, "Nerf in the wild: Neural radiance fields for unconstrained photo collections," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7210–7219.

[21] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," *Advances in neural information processing systems*, vol. 30, 2017.

[22] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *international conference on machine learning*. PMLR, 2016, pp. 1050–1059.