

Drone-assisted Road Gaussian Splatting with Cross-view Uncertainty

Gaining Zhang^{1,2*}

SAINING001@e.ntu.edu.sg

Baijun Ye^{1,3*}

yebaijun52@gmail.com

Xiaoxue Chen¹

chenxx21@mails.tsinghua.edu.cn

Yuantao Chen¹

yuantaochen973@gmail.com

Zongzheng Zhang¹

zzongzheng0918@gmail.com

Cheng Peng^{1,4}

120211642@bit.edu.cn

Yongliang Shi¹

shiyongliang@air.tsinghua.edu.cn

Hao Zhao^{1†}

zhaohao@air.tsinghua.edu.cn

¹ Institute for AI Industry Research (AIR),
Tsinghua University,
Beijing, China

² College of Computing and Data
Science,
Nanyang Technological University,
Singapore

³ IIIS,
Tsinghua University,
Beijing, China

⁴ School of Computer Science and
Technology,
Beijing Institute of Technology,
Beijing, China

Abstract

Robust and realistic rendering for large-scale road scenes is essential in autonomous driving simulation. Recently, 3D Gaussian Splatting (3D-GS) has made groundbreaking progress in neural rendering, but the general fidelity of large-scale road scene renderings is often limited by the input imagery, which usually has a narrow field of view and focuses mainly on the street-level local area. Intuitively, the data from the drone’s perspective can provide a complementary viewpoint for the data from the ground vehicle’s perspective, enhancing the completeness of scene reconstruction and rendering. However, training naively with aerial and ground images, which exhibit large view disparity, poses a significant convergence challenge for 3D-GS, and does not demonstrate remarkable improvements in performance on road views. In order to enhance the novel view synthesis of road views and to effectively use the aerial information, we design an uncertainty-aware training method that allows aerial images to assist in the synthesis of areas where ground images have poor learning outcomes instead of weighting all pixels equally in 3D-GS training like prior work did. We are the first to introduce the cross-view uncertainty to 3D-GS by matching the car-view ensemble-based rendering uncertainty to aerial images, weighting the contribution of each pixel to the training process. Additionally, to systematically quantify evaluation metrics, we assemble a high-quality synthesized dataset comprising both aerial and ground images for road scenes. Through comprehensive results,

*Equal contribution

†Corresponding author

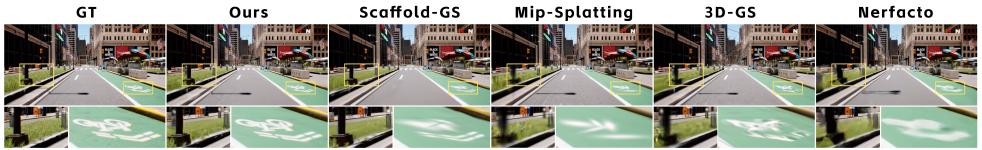


Figure 1: Qualitative results of our Drone-assisted Road Gaussian Splatting with Cross-view Uncertainty and several baseline methods. The dataset is 1.6m test set of New York City. The quality improvement is highlighted by boxes.

we show that: (1) Jointly training aerial and ground images helps improve representation ability of 3D-GS when test views are shifted and rotated, but performs poorly on held-out road view test. (2) Our method reduces the weakness of the joint training, and out-performs other baselines quantitatively on both held-out tests and scenes involving view shifting and rotation on our datasets. (3) Qualitatively, our method shows great improvements in the rendering of road scene details, as shown in Fig. 1. The code and data for this work will be released at <https://github.com/SainingZhang/UC-GS>.

1 Introduction

Autonomous driving simulation serves as a critical platform for scaling up to real-world deployment. Realistic rendering and novel view synthesis (NVS) for large-scale road scenes have become increasingly important in autonomous driving simulation, as they enable the synthesis of high-quality training and testing data at a significantly lower cost compared to using real-world data. This capability also benefits a wide range of applications, including digital cities [14, 47], virtual reality [43], and embodied AI [27].

Recently, NeRF [21, 25, 38, 40, 41, 52] has greatly enhanced fidelity of NVS by parameterizing the 3D scene as implicit neural fields but suffers from slow rendering due to exhaustive per-pixel ray sampling process especially in large-scale road scenes. 3D Gaussian Splatting (3D-GS) [17] achieves real-time rendering by rasterizing the learnable Gaussian primitives.

However, the rendering quality for both NeRF and 3D-GS is highly dependent on the input views. On the contrary, current road scene datasets, such as KITTI [9] or nuScenes [4], only contain car-view images, which are limited by the field of view and focus on the street-level local areas. The most related work may be MatrixCity [15], which offers both aerial and street-level city views, but the aerial imagery’s high altitude limits its ability to capture fine-grained road details. This mismatch in scale and granularity between global aerial views and local street views makes MatrixCity unsuitable for drone-assisted road scene synthesis.

To address the aforementioned issues, we dedicate to establish a new paradigm for **drone-assisted road scene synthesis**, aiming to overcome the limitations of input views by integrating an aerial perspective to provide a comprehensive global view of road scenes, in contrast to the localized perspective obtained from ground-level vehicle cameras. The aerial perspective can be captured using devices such as drones.

Since it is difficult to capture well-aligned view synthesis ground truth dataset for evaluation in real world, we first create a synthesized dataset (Fig. 2) comprising aerial-ground imagery with viewpoints that have similar levels of information granularity across large-scale road scenes. For the aerial perspective, we simulate drone flight trajectories and behavior patterns using AirSim [29]. For the ground-view images, we sample them to simulate the perspective and field of view from onboard cameras on vehicles.

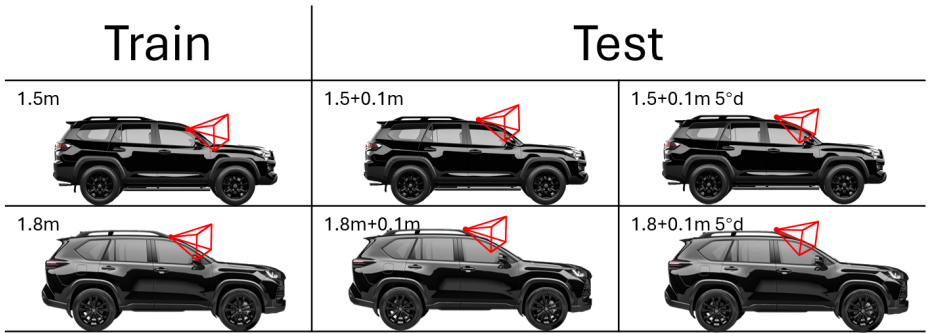


Figure 2: General view of the synthesized dataset. 5°d means 5 degrees downward.

On the other hand, the training process of 3D-GS leads it to an imperfect modeling when it comes to significant different viewpoints, such as aerial and ground views. This limited generalization capability is exacerbated in our specific setting, making simultaneous fitting of aerial and ground data challenging. For instance, the inherent ambiguity in aerial-ground data and the information outside co-view regions dilute useful information. Naively incorporating aerial data into the joint training with ground data adversely affects the convergence of the Gaussian field and compromises rendering quality for NVS when shifting and rotating views in autonomous driving scenes.

To overcome this problem, we introduce a novel uncertainty-aware training paradigm that enables more effective use of aerial imagery and guides 3D-GS to focus on challenging areas where ground data alone may struggle. By excluding irrelevant portions in aerial perspectives, such as the upper floors of buildings, which are less related to street scenes, we successfully mitigate the ambiguity and improve the fidelity of NVS such as view shifting and rotation on street.

The uncertainty is first computed within the ground-view image space through an ensemble-based method, and then projected to the aerial space to assist the training of 3D-GS, which is named as cross-view uncertainty as a new concept. Extensive experiments demonstrate that our uncertainty-aware training paradigm for drone-assisted road scene synthesis outperforms the naive joint training with aerial and ground images, as well as training with only ground images, both quantitatively and qualitatively. Our method offers significant benefits for applications like autonomous driving simulation. To summarize, the contributions of our work include:

- We formalize the problem of drone-assisted road scene synthesis and craft a high-quality and appropriate dataset for this new and important problem;
- We propose an uncertainty-aware training strategy and are the first to demonstrate that cross-view uncertainty can facilitate a pixel-weighted training paradigm while prior works use all pixels of images equally for 3D Gaussians' training;
- Through extensive experiments and evaluations, we demonstrate notably improved performance on both held-out tests and scenarios involving view shifting and rotation on road scene synthesis.

2 Related Works

2.1 3D Scene Representation

As the foundation for 3D computer vision, various 3D representations have been proposed to depict real-world scenes such as point cloud-based representation [22, 45], voxel-based representation [19, 50], or implicit representation [21, 26]. Among the implicit representation, NeRF [21] stands out as a groundbreaking neural rendering method that represents 3D scenes as continuous radiance fields parameterized by neural networks, taking coordinates and viewing directions as inputs. With the rise of NeRF, many efforts have been made to enhance its quality and efficiency [0, 1, 3, 5, 7, 17, 22, 42, 48, 51]. Recently, 3D Gaussian splatting (3D-GS) [12] has been proposed as a novel 3D scene representation, utilizing a set of 3D positions, opacity, anisotropic covariance, and spherical harmonic (SH) coefficients to represent a 3D scene. Compared with NeRF, 3D-GS based methods [6, 18, 35, 46, 49], shows superior performance in rendering speed, fidelity, and training time. In this work, we also leverage 3D-GS as the scene representation to resolve the problem of drone-assisted road scene synthesis.

2.2 Uncertainty Modeling

Modeling uncertainty has been a long-standing problem in deep learning. Early works usually resolve uncertainty estimation through Bayesian Neural Network (BNN) [23, 24]. However, these methods can be computationally expensive and challenging to implement. Later, dropout-based methods [8, 20, 37] have emerged as a computationally efficient alternative that adds dropout during inference to estimate uncertainty. Besides, ensemble-based methods [13, 16] have been proposed to model uncertainty by merging the prediction from multiple independently trained neural networks. As for the field of 3D scene representation, a series of works [10, 31, 32, 33] have focused on quantifying uncertainty in the prediction of NeRF. For example, S-NeRF [31] employ a probabilistic model to learn a simple distribution over radiance fields, while CF-NeRF [32] learns a distribution over possible radiance fields with latent variable modeling and conditional normalizing flows. With the emergence of 3D-GS, SGS [28] first addresses uncertainty modeling in 3D-GS, and integrates a variational inference-based method with the rendering pipeline of 3D-GS. CG-SLAM [11] also introduce the uncertainty-aware 3D-GS to SLAM. In this work, we introduce a novel cross-view uncertainty training paradigm to facilitate the training of 3D Gaussians on road scenes.

3 Method

Fig. 3 depicts the overview of our method. In Sec. 3.1, we briefly introduce the basic principles of original 3D-GS. Next, we construct the first drone-assisted road scene dataset in Sec. 3.2. Then, Sec. 3.3 illustrates how we model cross-view uncertainty through an ensemble-based rendering paradigm and an uncertainty projection module. Finally, by incorporating the uncertainty map into the loss function, we can build an uncertainty-aware training module, which facilitates the training of 3D-GS (Sec. 3.4).

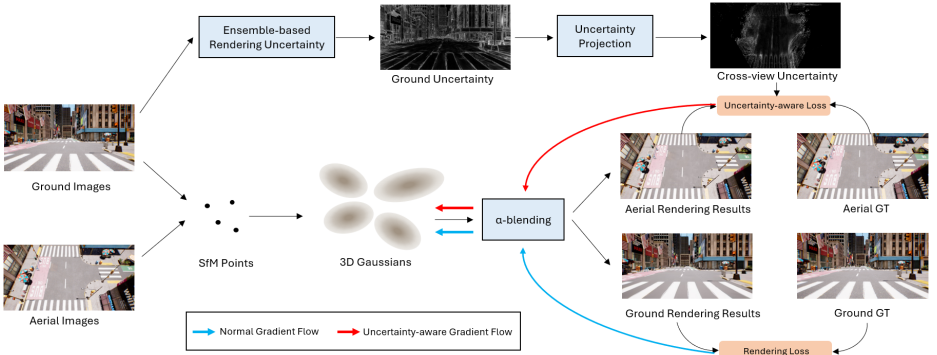


Figure 3: Overview of Drone-assisted Road Gaussian Splatting with Cross-view Uncertainty. We first adopt an ensemble-based rendering uncertainty to quantify the learning outcomes of 3D Gaussians on ground images. Next, the ground uncertainty is projected to the air to build the cross-view uncertainty. Subsequently, we introduce the cross-view uncertainty to the training of 3D Gaussians as weight for each pixel of aerial images in the loss function, together with the original rendering loss of 3D-GS for ground images.

3.1 Preliminaries

3D-GS [1] represents a 3D scene by a set of differentiable 3D Gaussians, which could be efficiently rendered to images through tile-based rasterization.

Specifically, initialized by a bunch of Structure-from-Motion (SfM) points, each 3D Gaussian is defined as:

$$G(x) = e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}, \quad (1)$$

where $x \in \mathbb{R}^{3 \times 1}$ is a random 3D position in the scene, $\mu \in \mathbb{R}^{3 \times 1}$ stands for the mean vector of the 3D Gaussian, and $\Sigma \in \mathbb{R}^{3 \times 3}$ refers to its covariance matrix. In order to maintain its positive semi-definite, Σ is further formulated as $\Sigma = RSS^T R^T$, where $R \in \mathbb{R}^{3 \times 3}$ is the rotation matrix and $S \in \mathbb{R}^{3 \times 3}$ is the scaling matrix.

To render the Gaussians into the image space, each pixel p is colored by α -blending N sorted Gaussians overlapping p as:

$$c(p) = \sum_{i=1}^N c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (2)$$

where α_j is calculated by multiplying the 2D Gaussian projected from 3D Gaussian G in p with the opacity of G , and c_i is the color of G . Through the differentiable tile-based rasterizer technique, all attributes of the Gaussians could be learnable and optimized end-to-end via training view reconstruction.

In this work, we utilize Scaffold-GS [18] as our baseline, as it represents the SOTA among 3D-GS based methods in road scene synthesis tasks. However, we posit that our proposed strategy holds promise for application across other 3D-GS based methods as well.

3.2 Drone-assisted Road Scene Synthesized Dataset

Crucial for autonomous driving simulation, high-fidelity view synthesis for large road scenes is often hindered by poor road rendering due to reliance on limited car-view imagery. To

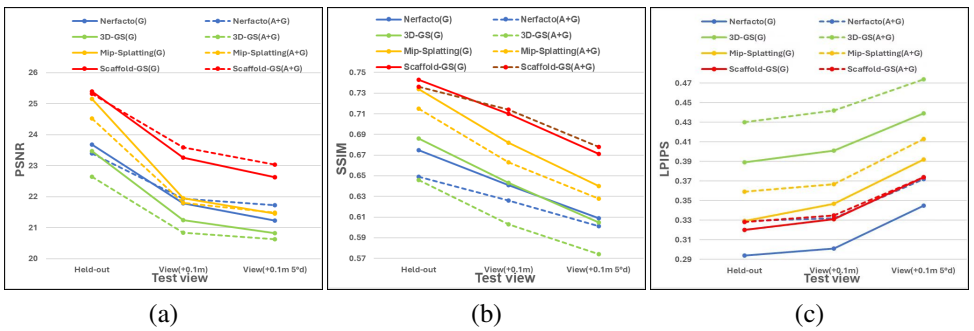


Figure 4: Results for training with ground or aerial and ground images on various models. (G), (A+G) are training with ground or aerial and ground images.

address this, we introduce drone-assisted road scene synthesis, leveraging aerial images as an additional input for a better scene reconstruction.

We present a new benchmark for assessing aerial images in large-scale road scene synthesis, featuring both aerial and ground posed images. Using Unreal Engine, we create two high-fidelity scenes to simulate real-world road imagery. AirSim [24] controls drones and vehicles for precise trajectory generation, simulating real scenarios. With the trajectories of drones, we employ AirSim to simulate the camera perspectives and render corresponding image data through the Unreal Engine. For ground-view imagery, we utilize vehicle trajectories to generate forward-facing images. As shown in Fig. 2, to replicate real-world driving conditions, we capture front-view ground images at heights of 1.5m and 1.8m, while aerial-view images are collected at the height 20m with the angle of 60° downward from the front view (based on some tests). Additional test data at 1.6m and 1.9m heights evaluate perspective impact. Each scene includes a training set of 315 ground and 351 aerial images, and a test set of 36 ground images, aiming to simulate diverse driving scenarios for a more representative benchmark dataset.

3.3 Cross-view Uncertainty Modeling

Preliminary Experiments and Motivation. From the comparison between the dotted and solid lines of different colors in Fig. 4 (a/b/c), it is clearly that jointly training aerial and ground images mitigates the decline in metrics during road view shifting and rotation compared with merely training with ground images. However, aerial images do not enhance the result on the held-out test of road scene synthesis, as shown by the point in the held-out column of Fig. 4 (a/b/c). Weighting all pixels from aerial and ground images equally while training will let aerial images have same synthesis priority as road views do for 3D Gaussians. The areas that are non-overlapped with road scene and the areas that have little contribution to the road scene synthesis in the aerial images not only fail to enhance the effectiveness of road reconstruction but also pose more challenges to 3D Gaussians’ convergence. This leads to poor rendering quality in the ground perspective when jointly training aerial and ground images.

Baseline and Implementation. In order to enhance the rendering result of road views, we attempt to quantify the contribution of each pixel in the aerial image to the road scene synthesis. However, undertaking such a task is very challenging. We decide to approach from a different angle by quantifying the quality of the learning outcomes of ground images to in-

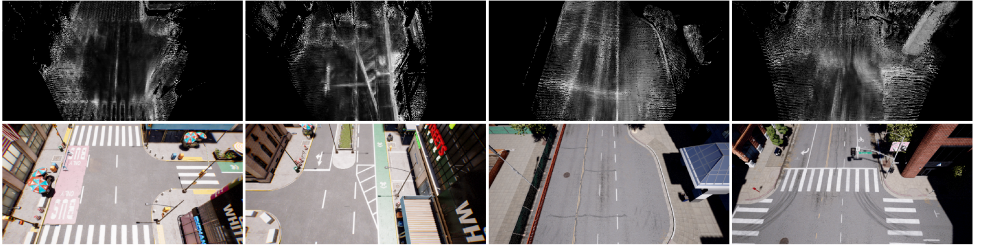


Figure 5: The first row shows the visualization of the cross-view uncertainty, and the second row shows the corresponding aerial data.

for the weight to each pixel in the aerial image during Gaussians’ training. To conveniently and plausibly compute the learning outcomes, we adopt an ensemble-based rendering uncertainty [43] paradigm to quantify the uncertainty of each pixel in the ground imagery. To be more specific, we train an *ensemble* of M gaussian splatting(GS)s initialised from the structure from motion (SfM) generated by ground imagery. By interpreting the ensemble as a uniformly-weighted mixture model, the members’ predictions are combined through averaging, and the predictive uncertainty is expressed as the variance over the individual member predictions. With an ensemble of GSs, the expected color of pixel p in a scene is

$$\mu_{\text{RGB}}(p) = \frac{1}{M} \sum_{k=1}^M c_k(p). \quad (3)$$

The predictive uncertainty can be expressed as the variance over the individual member predictions:

$$\sigma_{\text{RGB}}^2(p) = \frac{1}{M} \sum_{k=1}^M (\mu(p) - c_k(p))^2. \quad (4)$$

μ_{RGB} and σ_{RGB}^2 can be calculated very easily by rendering the M individual RGB images and calculating the mean and variance directly in pixel space. Both will be 3-vectors over the RGB colour channels.

We combine the variances from the colour channels into a single uncertainty value by:

$$u(p) = \frac{1}{3} \sum_{c \in \{\text{RGB}\}} \log(\sigma_{\text{RGB},(c)}^2(p) + 1), \quad (5)$$

where $\sigma_{\text{RGB},(c)}^2(p)$ indicates the variance associated with colour channel c , \log is the logarithmic transformation to smooth and tighten the values for further normalization process.

Cross-view Uncertainty Projection. To project the uncertainty map from ground-view to aerial-view, we test several methods. Neural field-based methods like NeRF and 3D-GS are prone to overfitting, so neural fields trained with ground uncertainty maps are unable to render high-quality uncertainty in the air. Besides, recently appeared end-to-end dense stereo model—DUST3R [49] has set SoTAs on many 3D tasks, which could be used as a 2D-2D pixel matcher between aerial and ground images. In this way, uncertainty maps from ground are projected to air through matches between ground and aerial images, and by averaging the uncertainties at pixels with multiple matches, we build reasonable cross-view uncertainty maps for training. The visualization of the cross-view uncertainty map is shown in Fig. 5.

3.4 Uncertainty-aware Training

In this section, we elaborate on how we introduce the cross-view uncertainty map to the training process. $U_k(x)$ presents the uncertainty value on the pixel position x of the k -th aerial image. First of all, we normalize all the uncertainty to the range (0, 1) like this:

$$U'_k = \left(\frac{U_k}{\max(U_1 \dots U_M) - \min(U_1 \dots U_M)} \right)^{\frac{1}{n}}, \quad (6)$$

where n refers to a hyperparameter for taking the n -th root, with the purpose of enhancing the impact of non-zero values.

Then, we introduce the uncertainty map to the color and SSIM loss as a weight map:

$$\mathcal{L}_{\text{color}} = \frac{1}{HW} \sum_{x=1}^{HW} U'(x) |\hat{C}(x) - C(x)|, \quad (7)$$

$$\mathcal{L}_{\text{SSIM}} = \text{mean}(U'(1.0 - \text{SSIM_MAP}(\hat{C}, C))), \quad (8)$$

where $\hat{C}(x)$ and \hat{C} represent ground-truth color, H and W stand for height and width of the image, and SSIM_MAP is the structural similarity of the inference and the ground-truth. The final loss function is given by:

$$\mathcal{L} = (1.0 - \lambda_{\text{SSIM}})\mathcal{L}_{\text{color}} + \lambda_{\text{SSIM}}\mathcal{L}_{\text{SSIM}} + \lambda_{\text{vol}}\mathcal{L}_{\text{vol}}, \quad (9)$$

where \mathcal{L}_{vol} is the volume regularization used in [18] to encourages the neural Gaussians to be small with minimal overlapping.

The loss function achieves the cross-view uncertainty-aware training which weights the effects of each pixel of aerial images to better assist in the road scene synthesis.

4 Experiments

4.1 Experimental Setup

Dataset and Metrics. In order to ensure the authenticity of the simulation data, we use two realistic city scene model, hereafter as New York City (NYC) and San Francisco (SF), from Kyrlo Sibiriakov [52] and Tav Shande [50] to collect the data as mentioned in Sec. 3.2. In addition to the 960×480 aerial and ground images, we also collected 1280×720 (HD) aerial images for experiments. All models are trained on 1.5m and 1.8m ground images, respectively, and tested at the viewpoint of the front and 5° downward at 0.1 meter above the ground level. All results are measured by three metrics: PSNR, SSIM and LPIPS.

Baseline and Implementation. Through preliminary experiments among several methods, Scaffold-GS [18] is selected as the baseline since its outstanding performance. All methods are trained for 900k iterations. Furthermore, we record the results of other SoTA methods in NVS like Nerfacto [56], 3D-GS [17] and Mip-Splatting [49].

For hyperparameters in eq. (9), $\lambda_{\text{SSIM}} = 0.2$ and $\lambda_{\text{vol}} = 0.001$ as in [18]. For n in eq. (6), the default value is 6 and we will discuss it in Sec. 4.3.

4.2 Main Results

From preliminary experiments (Fig. 4), it could be easily concluded that when testing the view shifting and rotation, the metrics of all methods decline. The inclusion of aerial images helps to slow down this trend compared to training the ground data along, indicating that aerial images can provide more perspective-rich information to maintain the rendering ability of the neural field during the view shifting and rotation. However, in the held-out viewpoints testing, training with aerial images performs no better than merely training on ground images, thus failing to demonstrate the superiority of aerial images in terms of view shifting and rotation.

Tab. 1 reports comprehensive results of road view synthesis on our datasets. After the implementation of the cross-view uncertainty, the paradigm makes a great progress. The average growth of PSNR on the held-out test set is 0.68 (NYC), and 0.41 (SF) compared with SoTA methods training with ground images. The SSIM and LPIPS also make a significant improvement. When shifting views, the PSNR is about 0.90 (NYC) and 0.80 (SF) more than training with ground images. The SSIM and LPIPS exhibit similar advancement trends. All our results out-performs previous SOTA solutions and Fig. 1 shows the improvement of our methods on certain details of road scene. In a word, our method not only enhances the representation of GS from ground perspectives but also improves the quality of road scenes synthesis during the view shifting and rotation.

Test set Method/Metrics	Held-out			View(+0.1m)			View(+0.1m 5°down)		
	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓
Nerfacto(G)	23.54	0.719	0.245	20.19	0.663	0.258	19.69	0.632	0.300
3D-GS(G)	23.71	0.706	0.363	20.54	0.688	0.346	20.01	0.646	0.387
Mip-Splatting(G)	25.35	0.779	0.302	20.51	0.710	0.302	20.03	0.668	0.350
Scaffold-GS(G)	25.64	0.790	0.265	22.19	0.746	0.281	21.55	0.705	0.326
Scaffold-GS(A+G)	25.66	0.782	0.273	22.56	0.744	0.286	22.10	0.709	0.328
Scaffold-GS(A*+G)	25.68	0.784	0.274	22.91	0.751	0.284	22.38	0.715	0.326
Ours	26.32	0.802	0.244	23.11	0.766	0.258	22.49	0.725	0.303

(a)

Test set Method/Metrics	Held-out			View(+0.1m)			View(+0.1m 5°down)		
	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓
Nerfacto(G)	23.82	0.631	0.344	23.38	0.618	0.344	22.77	0.587	0.421
3D-GS(G)	23.22	0.630	0.449	21.96	0.597	0.457	21.64	0.564	0.492
Mip-Splatting(G)	24.95	0.690	0.381	23.38	0.654	0.391	22.90	0.613	0.434
Scaffold-GS(G)	25.16	0.697	0.375	24.37	0.675	0.386	23.69	0.637	0.423
Scaffold-GS(A+G)	24.98	0.691	0.383	24.63	0.683	0.384	23.97	0.647	0.421
Scaffold-GS(A*+G)	24.99	0.689	0.386	24.70	0.684	0.385	23.88	0.649	0.422
Ours	25.57	0.723	0.337	25.18	0.715	0.338	24.55	0.678	0.376

(b)

Table 1: Results on NYC (a) and SF (b). A* is HD aerial images. (G), (A+G) are training with ground or aerial and ground images.

4.3 Ablation studies

Efficacy of Cross-view Uncertainty. Compared with equally training all aerial and ground images Tab. 1, the cross-view uncertainty-aware training achieves a 0.66 (NYC) and 0.59 (SF) increase in PSNR on held-out test set, and about 0.47 (NYC) and 0.57 (SF) when the view shifting and rotation. Moreover, our method also reverses the adverse effects of the joint training on SSIM and LPIPS, resulting in improvements in both metrics. It is also very impressive that our method performs even better than using HD aerial data when the view shifting and rotation. This reflects the effective utilization of aerial data in road scene

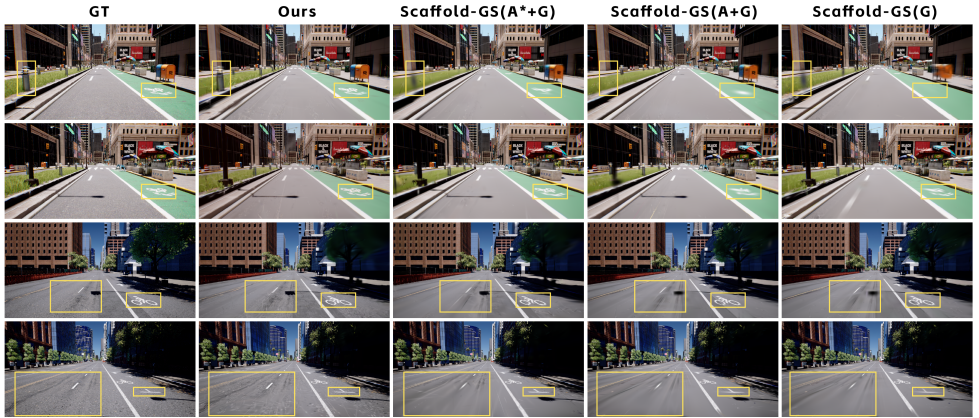


Figure 6: Rendering results for the ablation study of cross-view uncertainty. NYC: 1.6m (row 1), 1.6m 5° down (row2); SF: 1.9m (row 3), 1.9m 5° down (row4). A* is HD aerial images. (G), (A+G) are training with ground or aerial and ground images.

synthesis. From Fig. 6, it is clear that our method not only contributes to the rendering effect of road textures but also enhances the clarity of roadside obstacles, lane markings and ground signs, which will greatly aid in autonomous driving simulation.

Efficacy of Hyperparameter n . Tab. 2 presents the experimental results for different values of n in eq. (6). When n is set to 1 (i.e., no n), there is no significant improvement in the metrics compared to equally training on all aerial and ground images. However, as n increases to 2 or greater, the metrics improve with the increment of n , and the results become stable when $n \geq 6$. This indicates that when $n \geq 6$, the potential of aerial imagery is fully realized.

Test set n /Metrics	Held-out			View(+0.1m)			View(+0.1m 5°down)		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
1	25.64	0.748	0.310	23.71	0.723	0.320	23.11	0.685	0.360
2	25.72	0.754	0.301	23.83	0.729	0.309	23.24	0.690	0.353
3	25.80	0.758	0.298	23.93	0.733	0.305	23.39	0.694	0.347
4	25.85	0.759	0.296	24.04	0.735	0.303	23.47	0.698	0.344
6	25.94	0.762	0.291	24.14	0.741	0.298	23.52	0.701	0.339
8	25.89	0.761	0.291	24.13	0.739	0.298	23.47	0.700	0.339
10	25.91	0.763	0.290	24.20	0.742	0.296	23.55	0.704	0.337

Table 2: The results for the ablation study of hyperparameter n . n is for taking the n -th root to the uncertainty map. Metrics are averaged over testing on two datasets.

5 Conclusion

In this work, we propose a novel drone assisted road Gaussian Splatting with cross-view uncertainty. To use the global information from images in drones’ view to assist ground-view training, we are the first to introduce the cross-view uncertainty into the 3D-GS based model for weighting pixels in aerial images during training. This method reduces the impact of superfluous aerial information and effectively utilizes aerial images for road scene synthesis. From the experimental results, we achieve SoTA on two high-fidelity synthesized datasets. Our method enhances various metrics for held-out ground view synthesis while maintaining the robustness of aerial-ground training during the view shifting and rotation. The superiority of the method shows a great potential for the improvement of autonomous driving simulation in the near future.

Acknowledgement

This research is sponsored by Tsinghua-Toyota Joint Research Fund (20223930097).

References

- [1] 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.
- [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.
- [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.
- [4] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020.
- [5] Xiaoxue Chen, Junchen Liu, Hao Zhao, Guyue Zhou, and Ya-Qin Zhang. Nerf: 3d reconstruction and view synthesis for transparent and specular objects with neural refractive-reflective fields. *arXiv preprint arXiv:2309.13039*, 2023.
- [6] 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.
- [7] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022.
- [8] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [9] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11): 1231–1237, 2013.
- [10] Lily Goli, Cody Reading, Silvia Sellán, Alec Jacobson, and Andrea Tagliasacchi. Bayes’ rays: Uncertainty quantification for neural radiance fields. *arXiv preprint arXiv:2309.03185*, 2023.

- [11] Jiarui Hu, Xianhao Chen, Boyin Feng, Guanglin Li, Liangjing Yang, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. Cg-slam: Efficient dense rgb-d slam in a consistent uncertainty-aware 3d gaussian field. *arXiv preprint arXiv:2403.16095*, 2024.
- [12] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023.
- [13] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- [14] Ruilong Li, Sanja Fidler, Angjoo Kanazawa, and Francis Williams. Nerf-xl: Scaling nerfs with multiple gpus. *arXiv preprint arXiv:2404.16221*, 2024.
- [15] Yixuan Li, Lihan Jiang, Linning Xu, Yuanbo Xiangli, Zhenzhi Wang, Dahua Lin, and Bo Dai. Matrixcity: A large-scale city dataset for city-scale neural rendering and beyond. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3205–3215, 2023.
- [16] Jeremiah Liu, John Paisley, Marianthi-Anna Kioumourtoglou, and Brent Coull. Accurate uncertainty estimation and decomposition in ensemble learning. *Advances in neural information processing systems*, 32, 2019.
- [17] Junchen Liu, Wenbo Hu, Zhuo Yang, Jianteng Chen, Guoliang Wang, Xiaoxue Chen, Yantong Cai, Huan-ang Gao, and Hao Zhao. Rip-nerf: Anti-aliasing radiance fields with ripmap-encoded platonic solids. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024.
- [18] Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. *arXiv preprint arXiv:2312.00109*, 2023.
- [19] Jiageng Mao, Yujing Xue, Minzhe Niu, Haoyue Bai, Jiashi Feng, Xiaodan Liang, Hang Xu, and Chunjing Xu. Voxel transformer for 3d object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3164–3173, 2021.
- [20] Daily Milanés-Hermosilla, Rafael Trujillo Codorníu, René López-Baracaldo, Roberto Sagaró-Zamora, Denis Delisle-Rodríguez, John Jairo Villarejo-Mayor, and José Ricardo Núñez-Álvarez. Monte carlo dropout for uncertainty estimation and motor imagery classification. *Sensors*, 21(21):7241, 2021.
- [21] 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.
- [22] Raúl Mur-Artal and Juan D. Tardós. ORB-SLAM2: an open-source SLAM system for monocular, stereo and RGB-D cameras. *IEEE Transactions on Robotics*, 33(5): 1255–1262, 2017. doi: 10.1109/TRO.2017.2705103.
- [23] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.

- [24] Richard E Neapolitan et al. *Learning bayesian networks*, volume 38. Pearson Prentice Hall Upper Saddle River, 2004.
- [25] Julian Ost, Fahim Mannan, Nils Thuerey, Julian Knodt, and Felix Heide. Neural scene graphs for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2856–2865, 2021.
- [26] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019.
- [27] Ri-Zhao Qiu, Yafei Hu, Ge Yang, Yuchen Song, Yang Fu, Jianglong Ye, Jiteng Mu, Ruihan Yang, Nikolay Atanasov, Sebastian Scherer, and Xiaolong Wang. Learning generalizable feature fields for mobile manipulation, 2024.
- [28] Luca Savant, Diego Valsesia, and Enrico Magli. Modeling uncertainty for gaussian splatting. *arXiv preprint arXiv:2403.18476*, 2024.
- [29] Shital Shah, Debadeepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In *Field and Service Robotics: Results of the 11th International Conference*, pages 621–635. Springer, 2018.
- [30] Tav Shande. Artstation page <https://www.artstation.com/tavshande>, 2022.
- [31] Jianxiong Shen, Adria Ruiz, Antonio Agudo, and Francesc Moreno-Noguer. Stochastic neural radiance fields: Quantifying uncertainty in implicit 3d representations. In *2021 International Conference on 3D Vision (3DV)*, pages 972–981. IEEE, 2021.
- [32] Jianxiong Shen, Antonio Agudo, Francesc Moreno-Noguer, and Adria Ruiz. Conditional-flow nerf: Accurate 3d modelling with reliable uncertainty quantification. In *European Conference on Computer Vision*, pages 540–557. Springer, 2022.
- [33] Jianxiong Shen, Ruijie Ren, Adria Ruiz, and Francesc Moreno-Noguer. Estimating 3d uncertainty field: Quantifying uncertainty for neural radiance fields. *arXiv preprint arXiv:2311.01815*, 2023.
- [34] Kirill Sibiriakov. Artstation page <https://www.artstation.com/vegaart>, 2022.
- [35] Xiaowei Song, Jv Zheng, Shiran Yuan, Huan-ang Gao, Jingwei Zhao, Xiang He, Weihao Gu, and Hao Zhao. Sa-gs: Scale-adaptive gaussian splatting for training-free anti-aliasing. *arXiv preprint arXiv:2403.19615*, 2024.
- [36] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, et al. Nerfstudio: A modular framework for neural radiance field development. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–12, 2023.
- [37] Beiwen Tian, Liyi Luo, Hao Zhao, and Guyue Zhou. Vibus: Data-efficient 3d scene parsing with viewpoint bottleneck and uncertainty-spectrum modeling. *ISPRS Journal of Photogrammetry and Remote Sensing*, 194:302–318, 2022.

- [38] Haithem Turki, Jason Y Zhang, Francesco Ferroni, and Deva Ramanan. Suds: Scalable urban dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12375–12385, 2023.
- [39] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. *arXiv preprint arXiv:2312.14132*, 2023.
- [40] Yuxi Wei, Zi Wang, Yifan Lu, Chenxin Xu, Changxing Liu, Hao Zhao, Siheng Chen, and Yanfeng Wang. Editable scene simulation for autonomous driving via collaborative llm-agents. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15077–15087, 2024.
- [41] Zirui Wu, Tianyu Liu, Liyi Luo, Zhide Zhong, Jianteng Chen, Hongmin Xiao, Chao Hou, Haozhe Lou, Yuantao Chen, Runyi Yang, et al. Mars: An instance-aware, modular and realistic simulator for autonomous driving. In *CAAI International Conference on Artificial Intelligence*, pages 3–15. Springer, 2023.
- [42] DeJia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Humphrey Shi, and Zhangyang Wang. Sinnerf: Training neural radiance fields on complex scenes from a single image. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII*, pages 736–753. Springer, 2022.
- [43] Linning Xu, Vasu Agrawal, William Laney, Tony Garcia, Aayush Bansal, Changil Kim, Samuel Rota Bulò, Lorenzo Porzi, Peter Kotschieder, Aljaž Božič, et al. Vr-nerf: High-fidelity virtualized walkable spaces. In *SIGGRAPH Asia 2023 Conference Papers*, pages 1–12, 2023.
- [44] Shiyao Xu, Caiyun Liu, Yuantao Chen, Zhenxin Zhu, Zike Yan, Yongliang Shi, Hao Zhao, and Guyue Zhou. Camera relocalization in shadow-free neural radiance fields. *arXiv preprint arXiv:2405.14824*, 2024.
- [45] Wei Xu, Yixi Cai, Dongjiao He, Jiarong Lin, and Fu Zhang. Fast-lio2: Fast direct lidar-inertial odometry. *IEEE Transactions on Robotics*, 38(4):2053–2073, 2022.
- [46] Runyi Yang, Zhenxin Zhu, Zhou Jiang, Baijun Ye, Xiaoxue Chen, Yifei Zhang, Yuantao Chen, Jian Zhao, and Hao Zhao. Spectrally pruned gaussian fields with neural compensation. *arXiv preprint arXiv:2405.00676*, 2024.
- [47] Baijun Ye, Caiyun Liu, Xiaoyu Ye, Yuantao Chen, Yuhai Wang, Zike Yan, Yongliang Shi, Hao Zhao, and Guyue Zhou. Blending distributed nerfs with tri-stage robust pose optimization. *arXiv preprint arXiv:2405.02880*, 2024.
- [48] 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.
- [49] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d gaussian splatting. *arXiv preprint arXiv:2311.16493*, 2023.
- [50] Chongjian Yuan, Wei Xu, Xiyuan Liu, Xiaoping Hong, and Fu Zhang. Efficient and probabilistic adaptive voxel mapping for accurate online lidar odometry. *IEEE Robotics and Automation Letters*, 7(3):8518–8525, 2022.

-
- [51] Shiran Yuan and Hao Zhao. Slimmerf: Slimmable radiance fields. In *2024 International Conference on 3D Vision (3DV)*, pages 64–74. IEEE, 2024.
- [52] Tianyuan Yuan, Yucheng Mao, Jiawei Yang, Yicheng Liu, Yue Wang, and Hang Zhao. Presight: Enhancing autonomous vehicle perception with city-scale nerf priors. *arXiv preprint arXiv:2403.09079*, 2024.