# SGD: Street View Synthesis with Gaussian Splatting and Diffusion Prior

Zhongrui Yu[1†], Haoran Wang[2‡], Jinze Yang[3], Hanzhang Wang[4], Zeke Xie[2],
Yunfeng Cai[2], Jiale Cao[5], Zhong Ji[5], and Mingming Sun[2]

[1] ETH Zürich, [2] Baidu Research, [3] University of Chinese Academy of Sciences,
[4] Harbin Institute of Technology, [5] Tianjin University
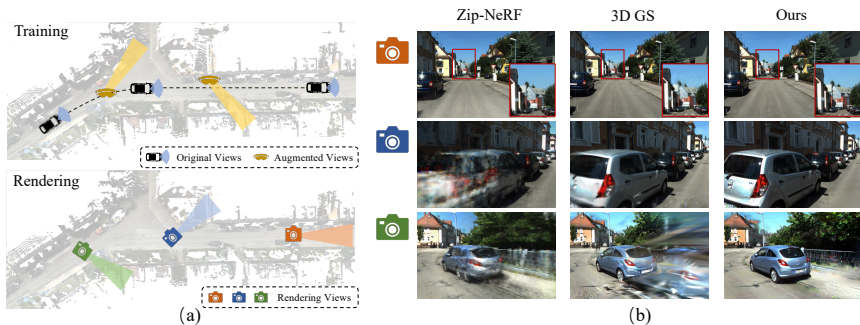zhonyu@ethz.ch, wanghaoran09@baidu.com

**Fig. 1: (a)**. To enable free control of ego-vehicle in autonomous driving simulation with novel view synthesis, we propose a method that leverages the prior from a Diffusion Model to provide 3DGS [12] augmented views during training. **(b)**. Our method preserves photo-realistic rendering quality at viewpoints that are distant from the training views while other approaches [1,12] produce severe artifacts.

**Abstract.** Novel View Synthesis (NVS) for street scenes play a critical role in the autonomous driving simulation. The current mainstream technique to achieve it is neural rendering, such as Neural Radiance Fields (NeRF) and 3D Gaussian Splatting (3DGS). Although thrilling progress has been made, when handling street scenes, current methods struggle to maintain rendering quality at the viewpoint that deviates significantly from the training viewpoints. This issue stems from the sparse training views captured by a fixed camera on a moving vehicle. To tackle this problem, we propose a novel approach that enhances the capacity of 3DGS by leveraging prior from a Diffusion Model along with complementary multi-modal data. Specifically, we first fine-tune a Diffusion Model by adding images from adjacent frames as condition, meanwhile exploiting depth data from LiDAR point clouds to supply additional spatial information. Then we apply the Diffusion Model to regularize the 3DGS at unseen views during training. Experimental results validate the effectiveness of our method compared with current state-of-the-art models, and demonstrate its advance in rendering images from broader views.

---

† Work done when Zhongrui Yu was a Research Intern at Baidu Research.

‡ The corresponding author.

# 1   Introduction

The driving simulation in street scenes holds crucial importance in the development of autonomous driving systems. Through constructing a digital twin of urban streets, we can continually enhance our autonomous driving system with simulated data. Thereby, the dependence of data collection in real scenarios is significantly reduced, making it possible to build a powerful autonomous driving system with lower time and financial cost.

For autonomous driving simulation, the early attempts [7,28,30] deploy Computer Graphics (CG) engines to render the images. It not only requires the time-consuming process to reconstruct virtual scenes, but also yields results with low realism and fidelity. Recently, neural rendering techniques for Novel View Synthesis (NVS), such as Neural Radiance Fields (NeRF) [18] and 3D Gaussian Splatting (3DGS) [12], are introduced for synthesizing photo-realistic street views. Current studies [9,17,20,24,33,37,41,45,53] mainly investigate two challenges faced in street view synthesis: the reconstruction of unbounded scenes and the modeling of dynamic objects. BlockNeRF [33] proposed to split scenes into multiple blocks, aimed at enhancing the model's capacity to present large unbounded street scenes. NSG [20] and some following methods [37,41,43,45,53] separately model the static background and the dynamic foreground to achieve higher background rendering quality while reducing motion blur associated with foreground vehicles.

Although thrilling progress has been made, a critical problem for evaluating the reconstruction quality is not well explored in existing works. It is known that an ideal scene simulation system should have the capacity to achieve free-view rendering with high quality. The current works commonly adopt the views that are sourced from vehicle captures yet unseen in the training stage as **test views**, (such as the red viewpoint in Fig. 1), while neglecting the **novel views** that deviate from the training views (such as the blue and green viewpoints in Fig. 1). When handling these novel views, there is a noticeable reduction in rendering quality with blurring and artifacts for existing works, as shown in Fig. 1. This issue is attributed to the inherently constrained view of the vehicle-collected images. The training images are typically captured along the vehicle's traveling direction and centered around the vehicle's lane. Due to the fast traveling speed of the vehicle, there is limited overlap between frames, thus not allowing for comprehensive multi-view observation of objects in the scene. Therefore, the street view synthesis task for autonomous driving can be comprehended as a reconstruction problem from sparse views.

Previous neural rendering methods proposed to address the challenge of NVS from sparse views can be categorized into two main branches. The first branch [6,32,38,42,48] incorporates scene priors, such as depth [6,25], normal [38], or the features extracted from a deep network [48] to regularize the model training in an explicit manner. Besides, another branch [16,21,29,31,40] attempts to leverage a pre-trained Diffusion Model for NVS. They typically fine-tuned a text-to-image Diffusion Model on some large multi-view datasets [3,5,23,49] to an image-to-image Diffusion Model with relative camera poses as condition, subsequently

apply the Diffusion Model to regularize the training of neural rendering model. However, a significant domain gap persists between multi-view datasets [3,5,23,49] and street scenes. And relying solely on relative camera poses is insufficient to learn the geometric details in more complex street scenes. To resolve this issue in the autonomous driving context, we leverage 3D geometric information obtained from multi-modal data to control the Diffusion model, enabling direct fine-tuning it on autonomous driving datasets and getting rid of the necessity of encoding relative camera poses.

To consolidate the idea, in this paper, we propose a novel NVS approach for street scenes, based on 3D Gaussian Splatting and prior from a fine-tuned Diffusion Model. We start by fine-tuning a Diffusion Model on an autonomous driving dataset [14]. For each input image, we employ its adjacent frames as the condition and leverage the depth information from LiDAR point clouds as control. This fine-tuned Diffusion Model then aids in guiding the 3DGS training by providing prior for the unseen views. Our method demonstrates competitive performance with the state-of-the-art (SOTA) methods [1,12,41] on KITTI [8] and KITTI-360 [14] datasets with dense viewpoint inputs and outperforms them in the sparse-view setting. Remarkably, our approach maintains high rendering qualities even for the viewpoints distant from training views. Moreover, since our approach is only applied during training, it does not compromise the real-time inference capability of 3DGS. Therefore, our model facilitates efficient rendering and versatile viewpoint control within autonomous driving simulation systems.

In Summary, we provide the following contributions:

- We propose a novel framework for Novel View Synthesis in street scenes, enhancing the freedom of view control in the premise of sustaining rendering efficiency for autonomous driving simulations.
- To the best of our knowledge, our method is the first attempt to tackle the street view synthesis task from the perspective of sparse-view-input reconstruction problem, and address this challenge by combining 3D Gaussian Splatting with a customized Diffusion Model.
- A novel strategy for fine-tuning a Diffusion Model on autonomous driving datasets and equipping it with NVS capability is presented, which overcomes the conventional reliance on multi-view datasets and relative camera poses.

## 2   Related Work

**Novel View Synthesis for Street Scenes** The rapid development of the NVS techniques including NeRF [18] and 3DGS [12] has attracted considerable attention within the arena of autonomous driving. A multitude of studies [4,9,15,17,20,24,33,35,37,41,45,47,53] have explored the utilization of these methods for street-view synthesis. Block-NeRF [33] and Mega-NeRF [36] have proposed to segment the scenes into distinct blocks for individual modeling. Urban Randiance Field [24] enhances the NeRF training with geometric information for LiDAR. DNMP [17] utilizes a pre-trained deformable mesh primitive to represent the scene. Streetsurf [9] delimits the scene into close-range,

distant-view and sky, achieving superior reconstruction results for urban street surfaces. For the modeling of dynamic urban scenes, NSG [20] presents the scene as neural graphs. MARS [41] employs distinct networks for modeling background and vehicles, creating an instance-aware simulation framework. With the emergence of 3DGS [12], DrivingGaussian [53] introduces Composite Dynamic Gaussian Graphs and incremental static Gaussians. StreetGaussian [45] optimizes the tracked pose of dynamic Gaussians and introduces 4D SH (spherical harmonics) for varying vehicle appearance in different frames.

In summary, current methods for street view synthesis mainly focus on two challenges: reconstructing the large-scale unbounded scenes and accurately modeling the dynamic vehicles. Yet the sparse-view-input issue within this task has not been adequately addressed.

**Novel View Synthesis with Sparse View Inputs** For NVS methods [12,18], intensive capture of the scene is paramount. Ideally, each part of the scene should be observed from serval perspectives, and a large overlap should exist across frames. The challenge of capturing such data leads to the development of methods that aim at improving rendering quality with sparse inputs [2, 6, 13, 26, 29, 32, 38–40, 42, 44, 46, 48, 54]. Early approaches involve supplementing the training process with scene priors. RefNeRF [38] incorporates a pre-trained model for normal flow to regularize novel viewpoints. DS-NeRF [6] enhances the reconstruction with depth information derived from SfM (Structure from Motion) point clouds. PixelNeRF [48] leverages a CNN encoder to extract features from images. The encoder can be trained across different scenes to acquire diverse priors.

Serval current methods leverage prior from Diffusion Model, which is pretrained on large-scale datasets, to support the synthesis of novel views. Zero-1-to-3 [16] and ZeroNVS [29] fine-tune the Diffusion Model by conditioning it on single image and a relative camera pose. Zero123++ [31] develops various conditioning and training schemes to minimize the effort of fine-tuning. ReconFusion [40] jointly train a PixelNeRF [48] along with fine-tuning the Diffusion Model, with the feature extracted by PixelNeRF serving as conditions for the Diffusion Model. It proves that PixelNeRF's feature provides a more accurate representation of the relative camera pose. However, the fine-tuning process of these methods is typically on large multi-view datasets, including ShapeNet [3], Objaverse [5], CO3D [23], MVImgNet, *etc*. These datasets are object-centric, maintaining a large domain gap from driving scenes. Inspired by the aforementioned methods, we propose a novel fine-tuning approach for Diffusion model tailored to autonomous driving scenarios, leveraging the 3D spatial information provided by multimodal data as conditions.

## 3    Method

The goal of street view synthesis is to render images from any viewpoint $v$, given a set of images and corresponding camera poses $\{I_i, \boldsymbol{p}_i\}_{i=1}^{N}$ captured by a vehicle.
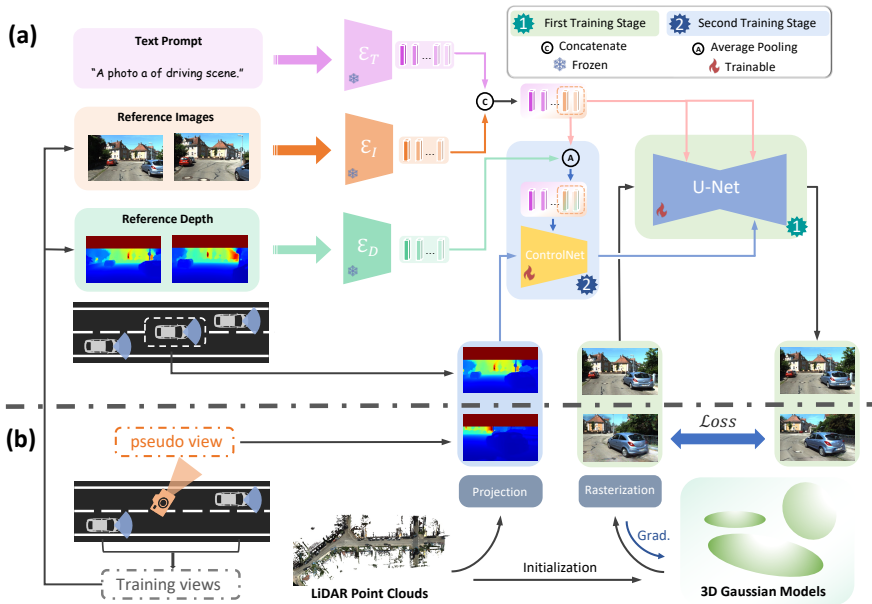
Fig. 2: **Overview of Our Method. (a).** There are two training stages in the Diffusion Model [27] fine-tuning. Firstly, the U-Net is fine-tuned by being injected with the patch-wise CLIP Image features of reference images concatenated with the CLIP text features of a text prompt. Secondly, a ControlNet is trained with the depth of the target image as the control signal. **(b).**The fine-tuned Diffusion Model from (a) guides the 3DGS training by providing regularization in pseudo views. For the sake of simplicity, the VAE encoder and decoder are omitted in the figure.

A big challenge arises from the constrained perspectives of images collected by a moving vehicle, where objects in the scene are often only observed from single viewpoint and appear only in a few images. To address this issue, we propose a novel method that leverages the priors derived from a fine-tuned Diffusion Model and the spatial information from LiDAR, to enhance the 3DGS model's awareness of the unobserved world.

Our method consists of two main components. We first fine-tune a Stable Diffusion Model [27] on a dataset of driving scenes [14] conditioning on reference images from adjacent frames and depth from LiDAR point cloud (detailed in Sec. 3.1). Subsequently, we integrate the fine-tuned Diffusion Model into the 3D Gaussian Splatting pipeline to guide the synthesis of unseen views (detailed in Sec. 3.2).

## 3.1 Fine-tuning Diffusion model

We propose a novel approach to fine-tune a Diffusion Model specifically on driving data. Driving data is collected sequentially, allowing us to easily identify the

closest preceding and succeeding frames from any novel viewpoint. The images from these adjacent frames are taken as reference images as they offer valuable contextual information. Moreover, the 360° LiDAR point clouds allow us to derive depth maps for both reference frames and the novel view, offering a comprehensive understanding of the relative spatial information across viewpoints. Briefly, by fine-tuning the Diffusion model, we guide it to learn about the contents that ought to be present from the context images and the spatial relationships among objects from the depth information.

The structure of our model is illustrated in Fig. 2 (a). During fine-tuning, information from reference images is introduced via cross-attention within the U-Net (in the first training stage), and the depth information is introduced by a ControlNet [50]-like module (in the second training stage).

**Training Stage 1: Image Conditioned Diffusion Model** In the first step, we want the Diffusion Model to learn high-level information about the scene from images of adjacent frames. Unlike other methods, such as ReconFusion [40], which introduces images and poses as conditions at the same time, we opt to not include any pose information at this stage. The reason we conduct this two-step training strategy is that the driving scene is complex, containing a variety of objects including buildings, vehicles, pedestrians, *etc.*, and the objects have mutual occlusion. The Diffusion Model is challenged in fully understanding the 3D scene by only encoding the relative camera pose. Thus, in this step we focus on enabling the Diffusion Model to identify *"what exits within the scene"*, and in next step we utilize the LiDAR point clouds to tell the model *"where each object is located"*.

As depicted in Fig. 2, besides the original structure of the Stable Diffusion Model [27], we introduce an additional pathway into the U-Net to integrate information from reference images encoded by the CLIP Image Encoder [22]. During the fine-tuning, We freeze the VAE encoder, decoder, CLIP Text Encoder and CLIP Image Encoder, and keep the parameters of U-Net $\theta$ trainable. The input image $I$ is encoded as a latent feature map $z_0$ by the VAE encoder. A text prompt $T$ is tokenized and encoded by the CLIP Text Encoder $\mathcal{E}_{text}$ to obtain text embedding,

$$e_T = \mathcal{E}_{\text{text}}(T) = [e_T^0, \ldots, e_T^{L_T}], \quad e_T^i \in \mathbb{R}^D \tag{1}$$

where $L_T$ denotes the length of the text token sequence and $D$ denotes the dimension of the CLIP embeddings. For each input image $I \in \mathbb{R}^{3 \times H \times W}$, a previous and a next frame are selected as the reference images $I_{\text{ref}} = \{I_{\text{pre}}, I_{\text{next}}\}$. Similar to the pathway of text, the reference images $I_{\text{pre}}$ and $I_{\text{next}}$ are encoded by the CLIP Image Encoder [22] $\mathcal{E}_{\text{Image}}$ separately, and the latent patch-wise CLIP Image features $e_{I_{\text{ref}}} = [e_{I_{\text{pre}}}, e_{I_{\text{next}}}]$ are concatenated to the text embedding along the dimension of token length.

$$e_{I_i} = \mathcal{E}_{\text{Image}}(I_i) = [e_{I_i}^0, \ldots, e_{I_i}^{L_I}], \quad i \in \{\text{pre}, \text{next}\} \tag{2}$$

where $L_I$ indicates the number of image patches. Through this process, we conduct a sequence of embeddings including not only the coarse-level information of the scene conveyed by the text prompt, but also the detailed semantic information introduced by the adjacent frames. This embedding is fused with the feature map $z_0$ via the cross-attention layer in the U-Net. The fine-tuning process is optimized by the loss funcion:

$$L(\theta) = \mathbb{E}_{z_0,\epsilon \sim \mathcal{N}(\mathbf{0,1}),t,T,I_{\mathrm{ref}}} ||\epsilon - \epsilon_\theta(z_t, t, e_T, e_{I_{\mathrm{ref}}})||_2^2 \qquad (3)$$

where $\epsilon$ is the random noise, $t$ is the denoising timestep and $z_t$ is the noisy latent at timestep $t$. The efficacy of the first step of fine-tuning is shown in Fig. 5 in the ablation study.

**Training Stage 2: Adding Depth ControlNet** The model acquires a high-level understanding of the scene after the first training step, yet it remains unaware of the spatial relationships between objects within the scene. In the second training step, we intend to utilize 3D information to control the model to achieve more accurate image generation. In autonomous driving scenarios, the 3D structure of the scene is preserved through LiDAR point clouds. Given an arbitrary viewpoint, the depth map can be derived by projecting the point clouds onto the image plane and then completed by an off-the-shelf depth completion model [52].

To enable depth control, we keep the fine-tuned Diffusion Model frozen and add a ControlNet [50] module to it, as depicted in Fig. 2. The ControlNet is initialized as a trainable copy the U-Net's encoder block. From the projection of the LiDAR point cloud, both the depth map for the reference images $D_{I_{\mathrm{ref}}}$ and for the input image $D_I$ can be obtained. $D_I$ serves as the input of the ControlNet. To encode the depth map of the reference images $D_{I_{\mathrm{ref}}}$, a depth encoder $\mathcal{E}_{depth}$ is introduced. The depth encoder is pre-trained to align the CLIP Image Encoder with contrastive learning proposed by [11]. The encoded depth feature of the reference images is fused to their corresponding CLIP Image features via average pooling.

$$e_{D_i} = \mathcal{E}_{\mathrm{Depth}}(I_i) \qquad (4)$$
$$\tilde{e}_{I_i} = \mathrm{AvgPool}(e_{I_i}, e_{D_i}), \quad i \in \{\mathrm{pre, next}\} \qquad (5)$$

The fused feature $\tilde{e}_{I_i}$ concatenated with $e_T$ is injected into the ControlNet via cross-attention. The ControlNet is trained with the loss function:

$$L(\tilde{\theta}) = \mathbb{E}_{z_0,\epsilon \sim \mathcal{N}(\mathbf{0,1}),t,T,\{I_{\mathrm{ref}},D_{\mathrm{ref}}\}} ||\epsilon - \epsilon_\theta(z_t, t, c_D, e_T, \tilde{e}_{I_{\mathrm{ref}}})||_2^2 \qquad (6)$$

where $\tilde{\theta}$ denotes the parameters of the ControlNet. The efficacy of our ControlNet is shown in Fig. 5 in the ablation study.

## 3.2 3D Gaussian Splatting with Diffusion Prior

Once the Diffusion Model is fine-tuned, we use its prior on images from unobserved views to regularize 3D Gaussian Model training.

3D GS [12] represents the scene as a large number of 3D Gaussian Models, each Gaussian Model is parameterized as its mean $\boldsymbol{\mu}$, covariance $\boldsymbol{\Sigma}$, opacity $\alpha$, and spherical harmonics parameters for view-dependent RGB color $\boldsymbol{c}$. For simplicity, we denote $\phi$ as the set of all trainable parameters of a single Gaussian Model. In each training step, a training view $v_o$ is sampled and rendered via differentiable rasterization. The parameters of the Gaussian Models corresponding to view $v_o$, denoted as $\Phi_o = \{\phi_i\}_{i \in N_o}$, are optimized by the loss function in Eq. (7), which is a combination of RGB loss, SSIM loss and depth loss. The depth loss is computed as the $L_1$ loss between the rendered depth map $\tilde{D}_o$ and the completed dense depth map $D_o$ from the projection of LiDAR point clouds.

$$L_{\text{recon}}(\Phi_o) = \mathbb{E}_{v_o}[||I_o, \tilde{I}_o||_1 + \lambda_{\text{SSIM}} L_{\text{SSIM}}(I_o, \tilde{I}_o) + \lambda_{\text{depth}}||D_o, \tilde{D}_o||_1] \quad (7)$$

Besides the training view $v_o$, we also randomly sample a set of pseudo views $\boldsymbol{v}_p = \{v_{p_i}\}_{i=0}^{M}$ in every $k$ training iterations. The position of the pseudo views is interpolated between the current training view and its adjacent views. And the orientation of the pseudo views is rotated from the training views within the range of $[-\delta, \delta]$ along the z-axis (yaw angle).

The loss function to regularize the 3D GS training is similar to the sample loss from [40], as it has been proved to perform better than score distillation sampling (SDS) [21, 40]. The pseudo views are also rendered via differentiable splatting, denoted as $\tilde{\mathbf{I}}_p(\Phi_p, \boldsymbol{v}_p) \in \mathbb{R}^{M,3,H,W}$, where $\Phi_p = \{\phi_i\}_{i \in N_p}$ is the parameters of Gaussian models corresponding to the pseudo views $\boldsymbol{v}_p$. Then the rendered images are fed into our Diffusion Model to generate the guidance image $\mathbf{I}_g$ in an image-to-image manner. The rendered pseudo views are encoded into latents and are added noise with a randomly selected noisy level $t$ between the max noisy level $t_{\max}$ and the min noisy level $t_{\min}$. The noisy latent is denoised from the selected noisy level to $t_{\min}$ and then decoded to obtain the guidance images $\mathbf{I}_g$. The loss function between the guidance images and rendered images from pseudo views is formulated as:

$$L_{\text{pseudo}}(\Phi_p) = \mathbb{E}_{\boldsymbol{v}_p}[||\mathbf{I}_g, \tilde{\mathbf{I}}_p||_1 + \lambda_{\text{p-lpips}} L_{\text{lpips}}(\mathbf{I}_g, \tilde{\mathbf{I}}_p) + \lambda_{\text{p-depth}}||\boldsymbol{D}_p, \tilde{\boldsymbol{D}}_p||_1] \quad (8)$$

The overall loss function is as follows:

$$L(\Phi) = L_{\text{recon}}(\Phi_o) + \lambda_{\text{pseudo}} L_{\text{pseudo}}(\Phi_p), \quad \Phi = \text{Union}\{\Phi_o, \Phi_p\} \quad (9)$$

where $\Phi$ denotes the union of the Gaussian Models' parameters corresponding to the training view and the pseudo views. Because the pseudo views are sampled close to training views, there exists a considerable overlap among their corresponding Gaussian Models. This allows these Models to learn from multi-views simultaneously, preventing them from getting stuck in local minima. Meanwhile, the plausible images from the pseudo views recovered by our fine-tuned Diffusion Model afford the Gaussian Models a more comprehensive observation of the scene, enhancing the capability for free-view rendering.

## 4    Experiments

### 4.1    Implementation Details

**Diffusion Model** Our Diffusion Model is fine-tuned based on Stable Diffusion 1.5 [27]. The additional CLIP Image Encoder is taken from clip-vit-B-32 [22]. It takes images with size $224 \times 224 \times 3$ as input, and its hidden state dimension is $1 \times 50 \times 768$, where 50 is the token length. The first token is the *cls* token and the last 49 tokens are the patch tokens. The Depth Encoder is pre-trained in the mechanism proposed from [11]. Its structure is identical to the CLIP Image Encoder. The Depth ControlNet is initialized as a trainable copy of the fine-tuned U-Net encoder from training stage one. In practice, we fine-tuned the Diffusion Model's U-Net with 625,000 iterations in the first stage and trained the ControlNet with 125,000 steps in the second stage. Both are on 4 32G V100 GPUs with batch size 4.

**3D Gaussian Splatting** . For initializing the Gaussian models of the scene, we only utilize the LiDAR point clouds which is voxelized-downsampled with a voxel size 0.5. No SfM (Structure from Motion) point clouds are used in our approach as we intend to avoid the impact from the prior provided by SfM point clouds and to ensure the exclusive use of data collected by vehicles. We trained the 3DGS model for 50,000 iterations with the learning rate decreasing from $1.6 \times 10^{-4}$ to $1.6 \times 10^{-6}$, which is identical to the original 3DGS.

### 4.2    Experiment Setup

**Datasets** We evaluate our method on two widely-used autonomous driving datasets KITTI [8] and KITTI-360 [14]. Both datasets provide forward-looking images and compact LiDAR point clouds. To prove the generalization ability, we only fine-tuned our Diffusion Model on about 12,000 images randomly selected from KITTI-360 datasets [14]. All images from KITTI datasets [8] are not seen during the fine-tuning. When training the 3DGS model, only monocular images are used.

**Competitors** Since our method is built upon 3DGS framework, we establish it as our baseline for comparison. To ensure a fair comparison, we also implement the depth loss for 3DGS, and keep its hyper-parameters identical to our method. For comparative analysis, we also select Zip-NeRF [1], which is the SOTA method in rendering qualities, and MARS [41], which is the SOTA method in street view synthesis, as our competitors. For Zip-NeRF, we use the PyTorch implementation with nerfstudio [34].

### 4.3    Comparison Results

**Evaluation on Test Views** For evaluating the rendering quality on the test views, we adopt three commonly used metrics PSNR, SSIM and LPIPS [51].

Tables 1 and 2 respectively present the quantitative comparison of our method again 3DGS [12], Zip-NeRF [1] and MARS [41] on KITTI [8] and KITTI-360 [14] datasets. In line with MARS [41], we evaluate our method on KITTI with 75%, 50% and 25% training/testing splits. For KITTI-360 dataset, we follow the setting of its official 50% drop-rate NVS benchmark [14].

Across all conducted experiments, our method substantially outperforms our baseline method 3DGS. Under the KITTI-75% setting, our results are inferior to those of MARS. This is primarily due to the uniform modeling of static background and dynamic objects in our method. This limitation can be similarly observed in Zip-NeRF [1]. However, at the sparse-view input setting in KITTI-25%, our method surpasses all the competitors. This improvement is attributed to the additional information of the unobserved view provided by the Diffusion Model. On KITTI-360 dataset, our method achieve the best performance among all competitors.

**Table 1: Quantitative results on the KITTI dataset [8].** *The results of MARS [41] are taken from their original paper.

| | KITTI - 75% | | | KITTI - 50% | | | KITTI - 25% | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| 3DGS [12] | 22.40 | 0.805 | 0.183 | 21.10 | 0.752 | 0.187 | 19.98 | 0.741 | 0.180 |
| Zip-NeRF [1] | 22.51 | 0.827 | 0.173 | 21.23 | 0.789 | 0.179 | 20.30 | 0.766 | 0.185 |
| MARS* [41] | 24.23 | 0.845 | 0.160 | 24.00 | 0.801 | 0.164 | 23.23 | 0.756 | 0.177 |
| Ours | 23.85 | 0.837 | 0.154 | 23.62 | 0.799 | 0.158 | 23.44 | 0.793 | 0.167 |

**Table 2: Quantitative results on the KITTI-360 dataset [14].** *The results of MARS [41] are taken from the leaderboard of KITTI-360 NVS benchmark [14].

| | KITTI-360 - 50% | | |
| --- | --- | --- | --- |
| | PSNR↑ | SSIM↑ | LPIPS↓ |
| 3DGS [12] | 22.78 | 0.793 | 0.176 |
| Zip-NeRF [1] | 22.86 | 0.802 | 0.167 |
| MARS* [41] | 23.09 | 0.857 | 0.174 |
| Ours | 23.81 | 0.832 | 0.155 |

**Evaluation on Novel Views** To evaluate each model's capability for free-view rendering, we also select novel views that are distant from the training and testing view for evaluation. These novel viewpoints are created by interpolating the position of the training/testing views and adding some perturbations. Their rotation is adjusted by offsetting angle $\pm\delta, \pm2\delta$ along the z-axis, where $\delta$ is randomly chosen from $[15°, 30°]$. Given the absence of corresponding ground-truth

images for these novel views, we adopt a no-reference image quality assessment method BRISQUE [19] to quantitatively measure the image quality, and the FID (Fréchet Inception Distance) score [10] to measure the difference of distribution between the rendered novels views and the training images, while qualitatively comparing rendering cases.

As indicated in Table 3, our method achieves the lowest BRISQUE score among all methods, suggesting that our rendered images in novel views preserve high quality with less noise and blur. Our method also exhibits the lowest FID score, reflecting a better alignment with the original images of street scenes. Figs. 3 and 4 illustrates the qualitative results among the methods on KITTI [8] and KITTI-360 [14]. It can be clearly seen in Fig. 3(a) that 3DGS and Zip-NeRF both produce severe artifacts of the blue vehicle. This is because they are trained predominantly by the rear view of the vehicle without observing it from other directions. Our method successfully mitigates this issue by incorporating the augmented views generated by the Diffusion Model. Fig. 3(b) and Fig. 4 further demonstrate that our method enhances the rendering quality, resulting in smoother road surface, more distinct lane markings, and clearer vehicles in the distance.

**Table 3: Image quality evaluation on novel views.**

|              | BRISQUE↓ | FID↓  |
|--------------|----------|-------|
| GT           | 20.71    | —     |
| 3DGS [12]    | 32.64    | 81.39 |
| Zip-NeRF [1] | 27.02    | 92.72 |
| Ours         | 24.53    | 73.36 |

### 4.4   Ablation Studies

We conduct ablation experiments on the two main processes within our method: the fine-tuning of the Diffusion Model and the 3DGS training process.

When fine-tuning the Diffusion Model, our method utilizes reference images and depth as conditions and conducts a two-phase training strategy. The influences of different condition signals and different training schemes are evaluated, as depicted in Fig. 5. The upper row of Fig. 5(a) and (b) shows the reference image, the target image, and the depth map for both the reference view and target view. The target view in Fig. 5(b) is a novel view, thus its original image is left blank. When solely utilizing a depth ControlNet while not considering reference images, the semantic information of the generated image is governed by the text prompt, resulting in a high diversity and significant deviations from the original images. Conversely, introducing reference images as the sole conditional input without depth guidance allows the Diffusion Model to assimilate scene semantics

(a)



(b)

**Fig. 3: Qualitative comparisons of novel views rendering on the KITTI-360 [14] dataset.** ZipNeRF [1] and 3DGS [12] produce artifacts of the blue vehicle in (a) and blurry lane markings in (b), while our method preserves high rendering quality. Our method also fix the hole on the road surface generated by 3DGS [12].

without understanding the precise locational context of scene objects. Such as in Fig. 5, without the depth information from the reference image, the model recognizes that there is a red car in the scene but misplaces it. Fine-tuning the U-Net and training the ControlNet simultaneously would result in less authentic outputs. This is due to, in such training scheme, the ControlNet has to be initialized by the weights of the original Stable Diffusion 1.5 other than the fine-tuned ones.

We also ablate the loss functions that regularize the training of pseudo views. As illustrated in Tab. 4, each loss validates its efficacy. However, it was observed that introducing $L_1$ loss for pseudo views marginally reduced the PSNR of the test views. This can be attributed to the inherent minor differences between the images produced by the Diffusion Model and the actual scenes, which become apparent upon pixel-level comparison.

**Fig. 4: Qualitative comparisons of novel views rendering on the KITTI [8] dataset.**
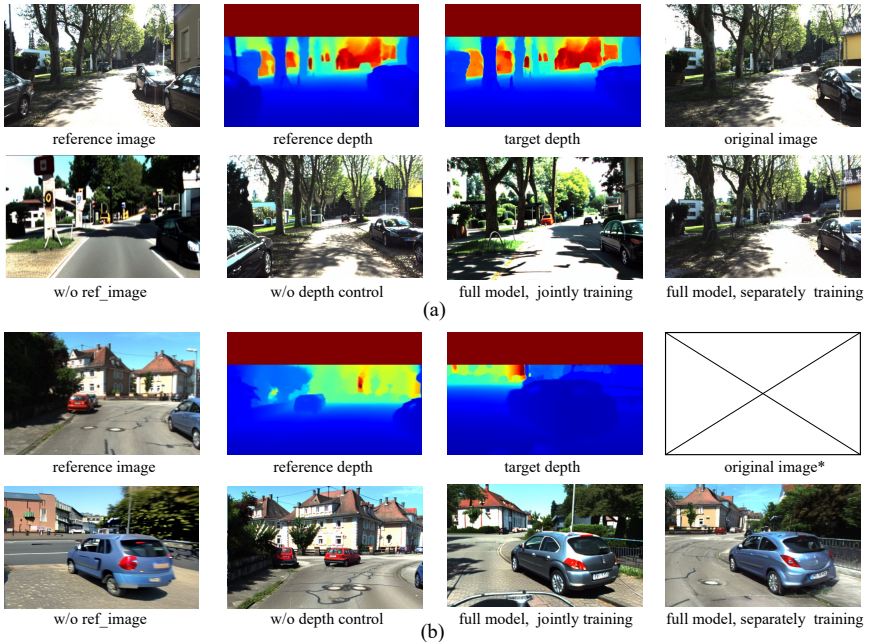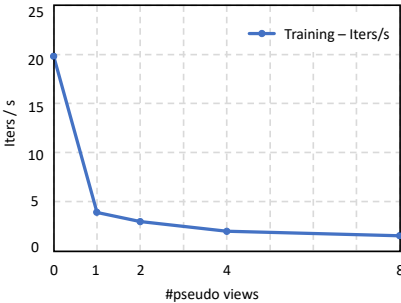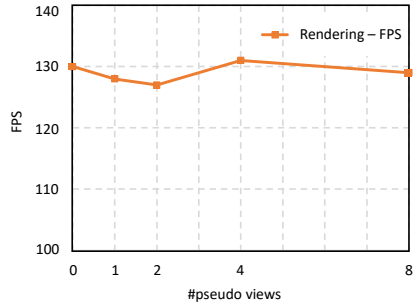


**Fig. 5: Qualitative ablation results on different conditions and different fine-tuning schemes of Diffusion Model [27].** *Target view in (b) is a novel view thus its original image is left blank.

Table 4: Quantitative ablation result on each module in 3DGS training.

|  | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| Complete model | 23.81 | 0.832 | 0.155 |
| w/o $L_{\text{LPIPS\_pseudo}}$ | 23.11 | 0.818 | 0.159 |
| w/o $L_{1\_\text{pseudo}}$ | 23.83 | 0.820 | 0.163 |
| w/o $L_{\text{depth\_pseudo}}$ | 23.29 | 0.801 | 0.169 |
| w/o pseudo views (3DGS) | 22.78 | 0.793 | 0.176 |



Fig. 6: Relationship between the number of sampled pseudo views per iteration and (a). training speed (in iterations per second), (b). rendering speed (in frames per second (FPS)). When pseudo views are introduced the training speed decreases significantly, while the inference speed is not impacted.

## 5    Discussion and Conclusion

**Limitation** The integration of Diffusion Model into 3DGS introduces a notable limitation: longer training time. It is primarily caused by the time-consuming denoising operation of Diffusion Model. Fig. 6a shows the correlation between the number of sample pseudo views and the training speed. A substantial decrease in training speed can be observed with an increment from 0 pseudo views (standard 3DGS) to 1. Since our method does not affect the real-time inference ability of 3DGS, as illustrated in Fig. 6b, and yields proved render quality, we temporarily accept the training time and leave improving the training efficiency as our future work.

In conclusion, we present a method aimed at enhancing the capability of free-viewpoint rendering within autonomous driving scenarios. While certain limitations persist, our method has shown proficiency in maintaining high-quality renderings from novel viewpoints, with considerable efficiency in rendering. This allows our method to offer a broader perspective within autonomous driving simulations, enabling the simulation of potentially hazardous corner cases, and thus enhancing the overall safety and reliability of autonomous driving systems.

# References

1. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Zip-nerf: Anti-aliased grid-based neural radiance fields. ICCV (2023)
2. Chan, E.R., Nagano, K., Chan, M.A., Bergman, A.W., Park, J.J., Levy, A., Aittala, M., De Mello, S., Karras, T., Wetzstein, G.: Generative novel view synthesis with 3d-aware diffusion models. arXiv preprint arXiv:2304.02602 (2023)
3. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012 (2015)
4. Chen, Y., Gu, C., Jiang, J., Zhu, X., Zhang, L.: Periodic vibration gaussian: Dynamic urban scene reconstruction and real-time rendering. arXiv preprint arXiv:2311.18561 (2023)
5. Deitke, M., Schwenk, D., Salvador, J., Weihs, L., Michel, O., VanderBilt, E., Schmidt, L., Ehsani, K., Kembhavi, A., Farhadi, A.: Objaverse: A universe of annotated 3d objects. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13142–13153 (2023)
6. Deng, K., Liu, A., Zhu, J.Y., Ramanan, D.: Depth-supervised nerf: Fewer views and faster training for free. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12882–12891 (2022)
7. Dosovitskiy, A., Ros, G., Codevilla, F., Lopez, A., Koltun, V.: Carla: An open urban driving simulator. In: Conference on robot learning. pp. 1–16. PMLR (2017)
8. Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: The kitti vision benchmark suite. URL http://www. cvlibs. net/datasets/kitti **2**(5) (2015)
9. Guo, J., Deng, N., Li, X., Bai, Y., Shi, B., Wang, C., Ding, C., Wang, D., Li, Y.: Streetsurf: Extending multi-view implicit surface reconstruction to street views. arXiv preprint arXiv:2306.04988 (2023)
10. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems **30** (2017)
11. Huang, T., Dong, B., Yang, Y., Huang, X., Lau, R.W., Ouyang, W., Zuo, W.: Clip2point: Transfer clip to point cloud classification with image-depth pre-training. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 22157–22167 (2023)
12. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. ACM Transactions on Graphics **42**(4) (2023)
13. Kwak, M.S., Song, J., Kim, S.: Geconerf: Few-shot neural radiance fields via geometric consistency. arXiv preprint arXiv:2301.10941 (2023)
14. Liao, Y., Xie, J., Geiger, A.: Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. IEEE Transactions on Pattern Analysis and Machine Intelligence **45**(3), 3292–3310 (2022)
15. Liu, J.Y., Chen, Y., Yang, Z., Wang, J., Manivasagam, S., Urtasun, R.: Real-time neural rasterization for large scenes. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8416–8427 (2023)
16. Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9298–9309 (2023)
17. Lu, F., Xu, Y., Chen, G., Li, H., Lin, K.Y., Jiang, C.: Urban radiance field representation with deformable neural mesh primitives. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 465–476 (2023)

18. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM **65**(1), 99–106 (2021)
19. Mittal, A., Moorthy, A.K., Bovik, A.C.: No-reference image quality assessment in the spatial domain. IEEE Transactions on image processing **21**(12), 4695–4708 (2012)
20. Ost, J., Mannan, F., Thuerey, N., Knodt, J., Heide, F.: Neural scene graphs for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2856–2865 (2021)
21. Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3d using 2d diffusion. In: The Eleventh International Conference on Learning Representations (2022)
22. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021)
23. Reizenstein, J., Shapovalov, R., Henzler, P., Sbordone, L., Labatut, P., Novotny, D.: Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10901–10911 (2021)
24. Rematas, K., Liu, A., Srinivasan, P.P., Barron, J.T., Tagliasacchi, A., Funkhouser, T., Ferrari, V.: Urban radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12932–12942 (2022)
25. Roessle, B., Barron, J.T., Mildenhall, B., Srinivasan, P.P., Nießner, M.: Dense depth priors for neural radiance fields from sparse input views. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12892–12901 (2022)
26. Roessle, B., Müller, N., Porzi, L., Bulò, S.R., Kontschieder, P., Nießner, M.: Ganerf: Leveraging discriminators to optimize neural radiance fields. arXiv preprint arXiv:2306.06044 (2023)
27. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022)
28. Rong, G., Shin, B.H., Tabatabaee, H., Lu, Q., Lemke, S., Možeiko, M., Boise, E., Uhm, G., Gerow, M., Mehta, S., et al.: Lgsvl simulator: A high fidelity simulator for autonomous driving. In: 2020 IEEE 23rd International conference on intelligent transportation systems (ITSC). pp. 1–6. IEEE (2020)
29. Sargent, K., Li, Z., Shah, T., Herrmann, C., Yu, H.X., Zhang, Y., Chan, E.R., Lagun, D., Fei-Fei, L., Sun, D., et al.: Zeronvs: Zero-shot 360-degree view synthesis from a single real image. arXiv preprint arXiv:2310.17994 (2023)
30. Shah, S., Dey, D., Lovett, C., Kapoor, A.: Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In: Field and Service Robotics: Results of the 11th International Conference. pp. 621–635. Springer (2018)
31. Shi, R., Chen, H., Zhang, Z., Liu, M., Xu, C., Wei, X., Chen, L., Zeng, C., Su, H.: Zero123++: a single image to consistent multi-view diffusion base model. arXiv preprint arXiv:2310.15110 (2023)
32. Somraj, N., Soundararajan, R.: ViP-NeRF: Visibility prior for sparse input neural radiance fields (August 2023). https://doi.org/10.1145/3588432.3591539
33. Tancik, M., Casser, V., Yan, X., Pradhan, S., Mildenhall, B., Srinivasan, P.P., Barron, J.T., Kretzschmar, H.: Block-nerf: Scalable large scene neural view synthesis.

In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8248–8258 (2022)

34. Tancik, M., Weber, E., Ng, E., Li, R., Yi, B., Kerr, J., Wang, T., Kristoffersen, A., Austin, J., Salahi, K., Ahuja, A., McAllister, D., Kanazawa, A.: Nerfstudio: A modular framework for neural radiance field development. In: ACM SIGGRAPH 2023 Conference Proceedings. SIGGRAPH '23 (2023)

35. Tonderski, A., Lindström, C., Hess, G., Ljungbergh, W., Svensson, L., Petersson, C.: Neurad: Neural rendering for autonomous driving. arXiv preprint arXiv:2311.15260 (2023)

36. Turki, H., Ramanan, D., Satyanarayanan, M.: Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12922–12931 (2022)

37. Turki, H., Zhang, J.Y., Ferroni, F., Ramanan, D.: Suds: Scalable urban dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12375–12385 (2023)

38. Verbin, D., Hedman, P., Mildenhall, B., Zickler, T., Barron, J.T., Srinivasan, P.P.: Ref-nerf: Structured view-dependent appearance for neural radiance fields. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5481–5490. IEEE (2022)

39. Wang, G., Chen, Z., Loy, C.C., Liu, Z.: Sparsenerf: Distilling depth ranking for few-shot novel view synthesis. arXiv preprint arXiv:2303.16196 (2023)

40. Wu, R., Mildenhall, B., Henzler, P., Park, K., Gao, R., Watson, D., Srinivasan, P.P., Verbin, D., Barron, J.T., Poole, B., Holynski, A.: Reconfusion: 3d reconstruction with diffusion priors. arXiv (2023)

41. Wu, Z., Liu, T., Luo, L., Zhong, Z., Chen, J., Xiao, H., Hou, C., Lou, H., Chen, Y., Yang, R., et al.: Mars: An instance-aware, modular and realistic simulator for autonomous driving. In: CAAI International Conference on Artificial Intelligence. pp. 3–15. Springer (2023)

42. Wynn, J., Turmukhambetov, D.: Diffusionerf: Regularizing neural radiance fields with denoising diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4180–4189 (2023)

43. Xie, Z., Zhang, J., Li, W., Zhang, F., Zhang, L.: S-nerf: Neural radiance fields for street views. arXiv preprint arXiv:2303.00749 (2023)

44. Xiong, H., Muttukuru, S., Upadhyay, R., Chari, P., Kadambi, A.: Sparsegs: Real-time 360 {\deg} sparse view synthesis using gaussian splatting. arXiv preprint arXiv:2312.00206 (2023)

45. Yan, Y., Lin, H., Zhou, C., Wang, W., Sun, H., Zhan, K., Lang, X., Zhou, X., Peng, S.: Street gaussians for modeling dynamic urban scenes. arXiv preprint arXiv:2401.01339 (2024)

46. Yang, J., Pavone, M., Wang, Y.: Freenerf: Improving few-shot neural rendering with free frequency regularization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8254–8263 (2023)

47. Yang, Z., Chen, Y., Wang, J., Manivasagam, S., Ma, W.C., Yang, A.J., Urtasun, R.: Unisim: A neural closed-loop sensor simulator. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1389–1399 (2023)

48. Yu, A., Ye, V., Tancik, M., Kanazawa, A.: pixelnerf: Neural radiance fields from one or few images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4578–4587 (2021)

49. Yu, X., Xu, M., Zhang, Y., Liu, H., Ye, C., Wu, Y., Yan, Z., Zhu, C., Xiong, Z., Liang, T., et al.: Mvimgnet: A large-scale dataset of multi-view images. In: Proceed-

ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9150–9161 (2023)

50. Zhang, L., Rao, A., Agrawala, M.: Adding conditional control to text-to-image diffusion models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3836–3847 (2023)

51. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 586–595 (2018)

52. Zhang, Y., Guo, X., Poggi, M., Zhu, Z., Huang, G., Mattoccia, S.: Completion-former: Depth completion with convolutions and vision transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18527–18536 (2023)

53. Zhou, X., Lin, Z., Shan, X., Wang, Y., Sun, D., Yang, M.H.: Drivinggaussian: Composite gaussian splatting for surrounding dynamic autonomous driving scenes. arXiv preprint arXiv:2312.07920 (2023)

54. Zhou, Z., Tulsiani, S.: Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12588–12597 (2023)

# Appendix

In this appendix, we provide additional details omitted from the main manuscript due to the limited space. First, we present additional figures to underscore the motivation behind our proposed method (Appendix A). Then we present more implementation details on fine-tuning Diffusion Model [27] and training 3DGS [12] (Appendix B). We also explore the influence of the Diffusion Model's prior on the generated results through dedicated experiment (Appendix C). Finally, we showcase more rendering results on the KITTI [8] and KITTI-360 [14] datasets (Appendix D).

# A    Motivation

Novel View Synthesis (NVS) for autonomous driving scenarios is a challenging task. The ideal training images for both NeRF [18] and 3DGS [12] should encompass all possible perspectives of the scene, which exhibit considerable disparities with the data collected by moving vehicles. The viewpoints offered by a vehicle-mounted camera are quite constrained. Take the white car in Fig. 7 as an example, it is only observed from its side rear in the training view, causing the rendering model to overfit these viewpoints. While the current approach, such as Zip-NeRF [1], is able to render the vehicle clearly from a test view close to the training view, it produces unsatisfactory artifacts and deformation when the rendering viewpoint is shifted by a certain distance and rotated by a certain angle.
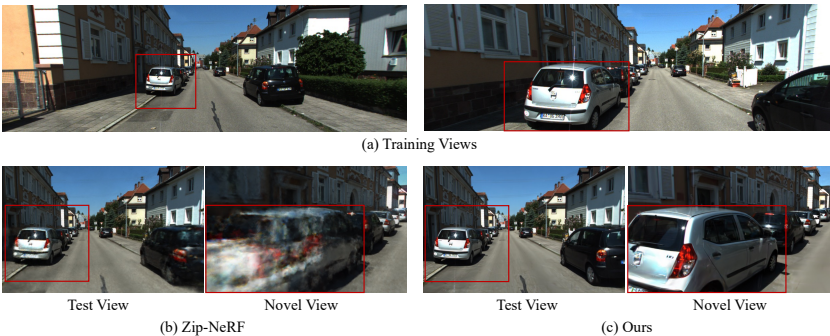


(a) Training Views

Test View          Novel View                          Test View          Novel View
(b) Zip-NeRF                                             (c) Ours

**Fig. 7:** An example of how the current method [1] overfits the training views, while our method overcomes this problem.

# B    More Implementation Details

**Diffusion Model** Our Diffusion Model is adapted from Stable Diffusion 1.5 [27] and is fine-tuned on about 12,000 images with $512 \times 512$ resolution from the KITTI-360 [14] dataset. Considering the original size of KITTI-360 images

is $1408 \times 376$, a preliminary cropping step to $600 \times 376$ is performed before the resizing, to avoid over-distorting the images. We conduct *center-crop* on the training images. For the reference images, we use *random-crop* during the training process, which could ensure a certain perspective gap exists between the reference image and the training image, so as to enhance the robustness of the model. During inference, the reference images are pre-processed with *center-crop*.

When selecting the reference images, we randomly choose one image from the five frames preceding the training image and one from the five frames succeeding the training image separately. During inference for the novel viewpoint, we identify its closest training viewpoint and utilize its adjacent frames as reference images. Regarding the depth maps, due to the limitation of LiDAR point clouds in capturing the scene above a certain height, we apply a mask to the top 80 rows of pixels in the images. In practice, we found that the inpainting capability of the Stable Diffusion Model is effectively able to complete this portion of content. To enable classifier-free guidance in the first training stage, we set both text prompts and reference images to be empty with a 10% probability.

**3D Gaussian Splatting** We only initialize the 3D Gaussian models with Li-DAR point cloud. The detailed procedure involves first projecting LiDAR frame onto its corresponding image frame to assign a color to each LiDAR point. Then these points are re-projected into 3D space, creating colored 3D point clouds. Finally, all frames of point clouds are accumulated and then voxel-downsampled with the voxel-size of 5. We train both our model and the baseline 3DGS model for 50,000 iterations. We first train the model for 500 iterations without sampling pseudo views for adequate warm-up. Subsequently, for every 10 iterations, 4 pseudo views are sampled for training.

## C    Additional Experiment

As described in Sec. 3.2 of the main manuscript, during the training stage of 3DGS, we render some randomly sampled pseudo views, and utilize a fine-tuned Diffusion Model to generate guidance images for these views to regularize the training. Specifically, the pseudo view rendered by 3DGS is passed through the VAE Encoder to obtain a latent feature map, to which noise at level $t$ is added, where $t \sim [t_{\min}, t_{\max}]$. This noised latent feature is denoised by the Diffusion model from level $t$ to $t_{\min}$, and then it is decoded to obtain the generated image. Specifically, we set $t_{\max} = 10$, and employ a hyper-parameter $s$, which indicates strength, to control the noise level $t$, according to $t = s \times t_{\max}$.

In Fig. 8, we show the results of ablation experiments on hyper-parameter $s$. The first column labeled with *original image* refers to the image being fed into the Diffusion Model, while the generated image with hyper-parameter $s$ increasing from 0.2 to 0.8 are exhibited in the other columns. It can be observed that a smaller $s$ makes generated images more similar to the original image, while a large $s$ introduces higher diversity and deviation in details. For novel viewpoints in Fig. 8(c), smaller $s$ makes the generated image preserve noise rendered by 3DGS. As $s$ increases, the image becomes cleaner but loses some details. In
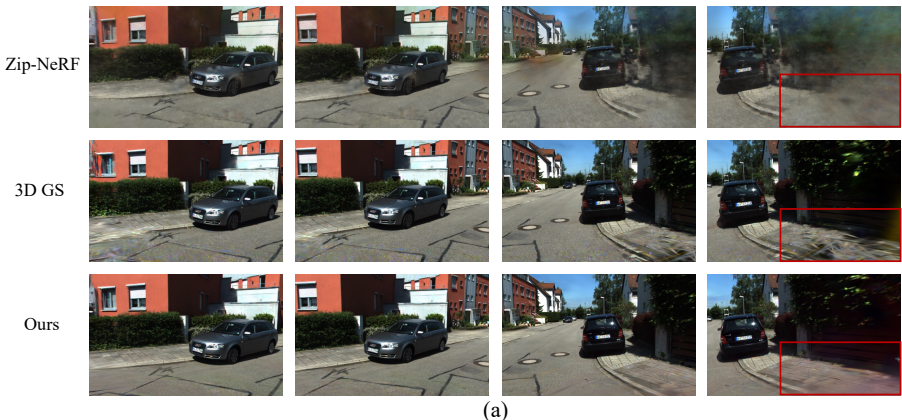
practice, we randomly select $s \sim [s_{\min}, s_{\max}]$ for each sampled pseudo view, where $s_{\min} = 0.2$, $s_{\max}$ starts at 0.6 and decreases to 0.4 over the training process. This strategy guarantees when 3DGS-rendered images are of lower quality in the early stage of training, our model relies more on the guidance from the Diffusion Model's prior. Accompanied by the quality of 3DGS renderings improves with ongoing training, it is necessary to reduce the impact of the Diffusion Model-generated images on the details.



**Fig. 8: The impact of the strength of the Diffusion Model's prior on the generated result.** *(c) is a novel view, its original image is rendered by 3DGS.

## D   More Rendering Results

We provide more novel view rendering results of our method and our competitors [1, 12] on the KITTI [8] and KITTI-360 [14] datasets in Fig. 9.
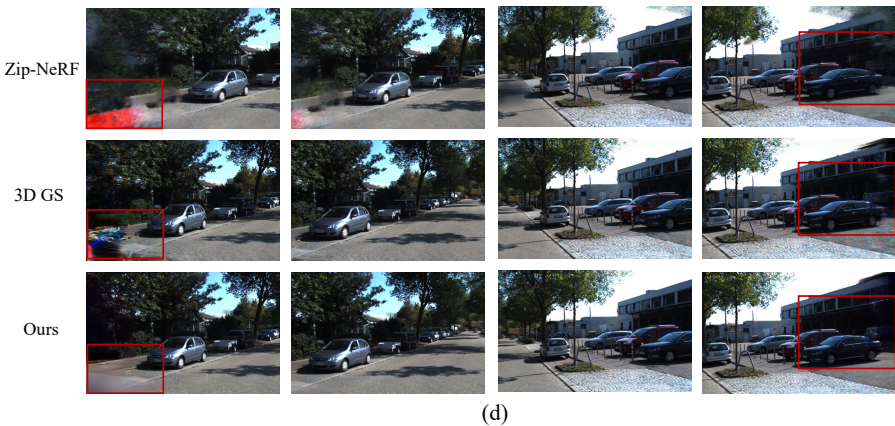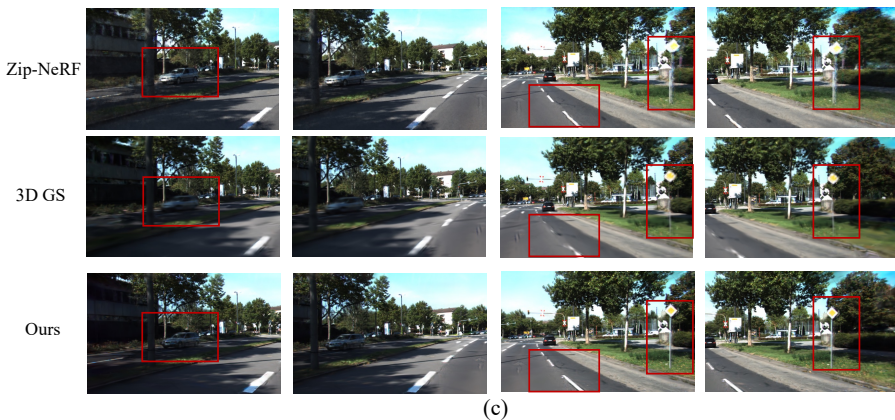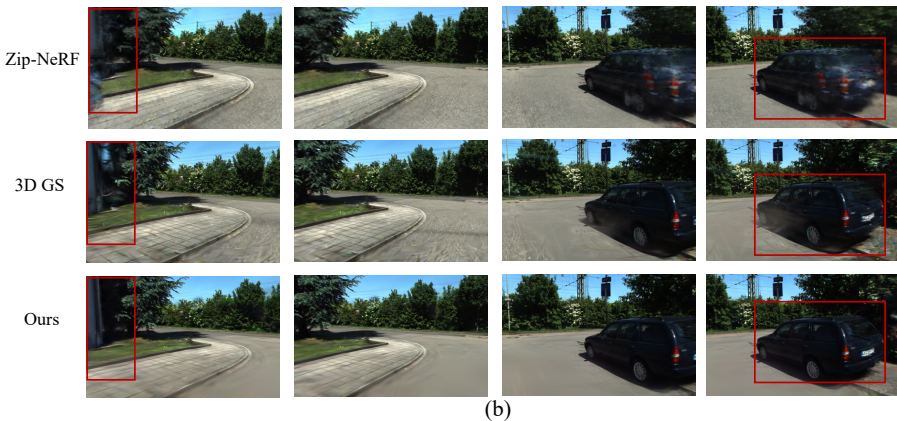


(a)

(b)



(c)



(d)

Fig. 9: **More qualitative results of novel views rendering on the KITTI [8] and KITTI-360 [14] dataset.**