

# SplatFlow: Self-Supervised Dynamic Gaussian Splatting in Neural Motion Flow Field for Autonomous Driving

Su Sun\*<sup>1</sup>, Cheng Zhao\*<sup>2</sup>, Zhuoyang Sun<sup>1</sup>, Yingjie Victor Chen<sup>1</sup>, Mei Chen<sup>2</sup>  
<sup>1</sup>Purdue University, <sup>2</sup>Microsoft

## Abstract

Most existing Dynamic Gaussian Splatting methods for complex dynamic urban scenarios rely on accurate object-level supervision from expensive manual labeling, limiting their scalability in real-world applications. In this paper, we introduce SplatFlow, a Self-Supervised Dynamic Gaussian Splatting within Neural Motion Flow Fields (NMFF) to learn 4D space-time representations without requiring tracked 3D bounding boxes, enabling accurate dynamic scene reconstruction and novel view RGB/depth/flow synthesis. SplatFlow designs a unified framework to seamlessly integrate time-dependent 4D Gaussian representation within NMFF, where NMFF is a set of implicit functions to model temporal motions of both LiDAR points and Gaussians as continuous motion flow fields. Leveraging NMFF, SplatFlow effectively decomposes static background and dynamic objects, representing them with 3D and 4D Gaussian primitives, respectively. NMFF also models the status correspondences of each 4D Gaussian across time, which aggregates temporal features to enhance cross-view consistency of dynamic components. SplatFlow further improves dynamic scene identification by distilling features from 2D foundational models into 4D space-time representation. Comprehensive evaluations conducted on the Waymo Open Dataset and KITTI Dataset validate SplatFlow’s state-of-the-art (SOTA) performance for both image reconstruction and novel view synthesis in dynamic urban scenarios.

## 1. Introduction

As autonomous driving systems increasingly shift toward end-to-end models, there is a growing need for scalable simulation environments without domain gaps where these systems can undergo closed-loop evaluation. A promising approach is real-world closed-loop evaluation, which demands controllable sensor inputs, thereby driving the development of advanced scene reconstruction techniques. In this context, Neural Radiance Fields (NeRFs) [9] and 3D

\*Equally contributed as co-first author.

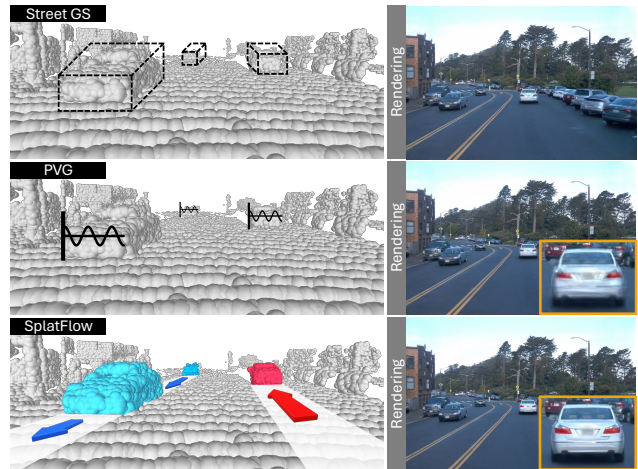


Figure 1. Top: Street GS [23]; Middle: PVG [1]; Bottom: Our SplatFlow. SplatFlow eliminates the need for 3D Bboxes required by Street GS, and enhances rendering quality compared to PVG.

Gaussian Splatting (3DGS) [6] have proven to be effective in creating high-quality 3D scene reconstructions with excellent visual and geometric accuracy. However, accurately and comprehensively reconstructing dynamic driving scenes remains a major challenge, given the complexity of real-world scenarios without dynamic object annotations.

Recent techniques enhance NeRF-based models by using object detection and tracking to animate dynamic objects, enabling photo-realistic view generation of dynamic urban street environments. Approaches [3, 7, 10, 16, 19, 22] create a scene graph, where both dynamic objects and the static background are represented as nodes, reconstructed within their canonical frames. Self-supervised methods [17, 24] model dynamic driving scenes by combining static neural field and time-dependent neural field to handle the static background and moving foreground objects separately. The SUDS [17] uses optical flow to ease the strict requirement for object labeling, while EmerNerf [24] learns temporal attributes and features in a self-supervision manner to minimize reliance on optical flow. However, despite utilizing implicit representations, these NeRF style methods face efficiency challenges in both training and inference, which becomes a major bottleneck for large-scale scene reconstruct-

tion and rendering tasks.

While NeRF has proven effective for driving scenes, 3DGS offers a promising alternative due to faster training and rendering speed with more explicit representation. However, the original 3DGS faces notable challenges in modeling dynamic urban environments due to its limited representation capabilities in the temporal dimension. To address this, DrivingGaussian [27] introduces composite dynamic Gaussian graphs to handle multiple moving objects and incremental static Gaussians for background representation. StreetGaussian [23] optimizes the tracked bounding boxes of dynamic Gaussians along with 4D spherical harmonics to capture changing vehicle appearances. Both methods, however, require accurate object-level supervision, such as 3D object boxes and trackers, to decompose static and dynamic elements, limiting their scalability in real-world applications. In response, PVG [1] proposes a Periodic Vibration Gaussian mechanism to model dynamic driving scenes without relying on manually labeled 3D bounding boxes. However, the scene decomposition in PVG relies on time-dependent Gaussian attributes optimized only through rendering losses, but ignores motion information within the input point cloud for Gaussian initialization, leading to suboptimal performance in challenging scenarios with rapid movements.

To facilitate accurate scene reconstruction and real-time rendering without expensive annotations, we propose SplatFlow, a self-supervised Dynamic Gaussian Splatting method in Neural Motion Flow Fields (NMFF) for dynamic urban scenarios. In contrast to existing methods, as illustrated in Fig. 1, SplatFlow decomposes dynamic objects and the static background in a self-supervised way without requiring expensive 3D bounding box annotations. The key idea of SplatFlow is to seamlessly integrate 4D Gaussian representations within NMFF in a unified framework. The NMFF is a set of implicit functions to model temporal motions of both LiDAR points and Gaussians as continuous motion flow fields. Leveraging NMFF, SplatFlow not only enables the decomposition of dynamic and static elements from 3D LiDAR points, but also enables the status conversion of each 4D Gaussian across time. During Gaussian splatting, we represent dynamic objects using aggregated 4D Gaussians at various viewpoints and timestamps, and the static background using 3D Gaussians, respectively. In SplatFlow, NMFF’s capabilities are enhanced through three components: 1) learning 3D priors by pretraining on 3D LiDAR data, 2) optimizing temporal status transitions of 4D Gaussians on image data, and 3) distilling knowledge from 2D foundation models.

The novel features of SplatFlow are summarized as:

- SplatFlow introduces a unified framework that seamlessly integrates time-varying 4D Gaussian representations into NMFF, enabling self-supervised dynamic scene recon-

struction and rendering.

- NMFF enables scene decomposition, modeling static elements with 3D Gaussians and dynamic elements with 4D Gaussians.
- NMFF models status correspondences of each 4D Gaussian over time, aggregating temporal features to enhance the cross-view consistency of dynamic components.
- SplatFlow enhances the dynamic scene identification by uplifting 2D features distilled from foundational models to 4D space-time through optical flow rendering.

Comprehensive experiments on the Waymo [15] and KITTI [4] benchmarks demonstrate that SplatFlow outperforms the state-of-the-art (SOTA) methods in both image reconstruction and novel view synthesis for dynamic urban scenes. Notably, SplatFlow achieves this by avoiding tracked 3D bounding boxes of dynamic objects, allowing the proposed model to learn from extensive, in-the-wild data sources.

## 2. Related work

Real-world simulation is essential for generating data to support end-to-end autonomous driving solutions. However, current simulation engines like AirSim [14] and CARLA [2] encounter difficulties due to the expensive manual work needed to create virtual environments and the domain gap between real and simulated data. The rapid progress in Novel View Synthesis (NVS) technologies, such as NeRF [9] and 3DGS [6], enables 3D reconstruction and photorealistic image generation, significantly improving real-world simulation for autonomous driving.

For 3D reconstruction and neural rendering in autonomous driving scenarios, sensor fusion solution combining surrounding cameras and LiDAR is commonly used in existing work [5, 8, 11, 13, 25]. Urban Radiance Field [13] improves NeRF training by incorporating 3D information from LiDAR to enhance 3D geometry learning. DNMP [8] employs a pre-trained deformable mesh primitive to represent the 3D scene, enhancing neural rendering quality. NPLF [11] uses explicit 3D reconstructions from LiDAR data to efficiently model the radiance field during rendering. StreetSurf [5] categorizes scenes into close-range, distant-view, and sky sections, achieving superior reconstruction of urban street surfaces. TCLC-GS [25] design a hybrid explicit and implicit 3D representation derived from LiDAR-camera data, to enrich the properties of 3D Gaussians for splatting. While these methods deliver impressive results for static urban scene rendering, they struggle with handling dynamic objects common in driving scenarios.

To model dynamic urban scenes, recent approaches [7, 10, 17, 19, 24] decompose the entire scene into static background and dynamic objects, learning their representations separately. PNF [7] uses monocular 3D bounding box predictions to isolate dynamic objects and fur-



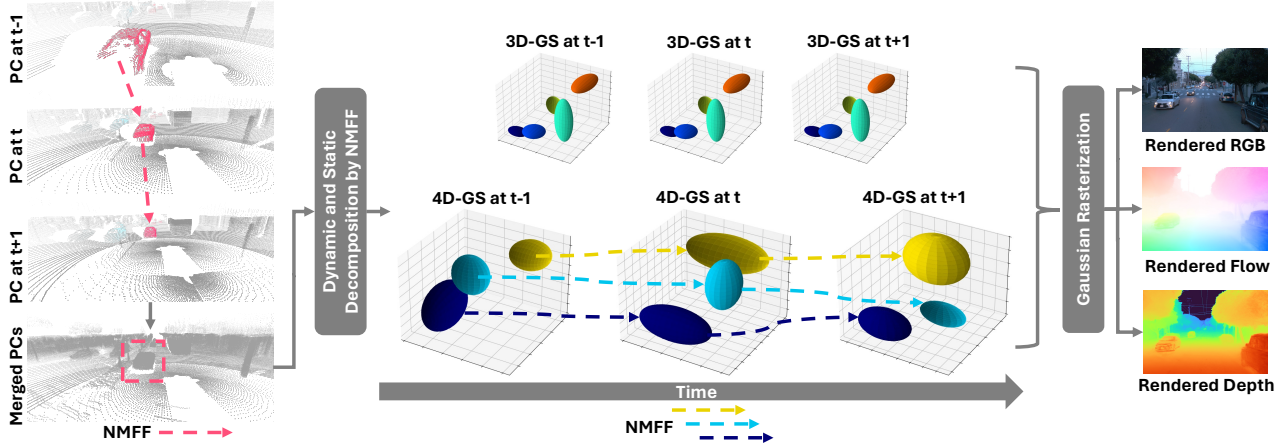


Figure 2. The pipeline of SplatFlow.

they jointly optimize their poses during the reconstruction process. NSG [10] uses neural graphs to represent entire scenes, decomposing dynamic multi-object environments. MARS [19] employs distinct sub-networks to model backgrounds and dynamic objects, creating an instance-aware simulation framework. These methods typically rely on manually annotated or predicted 3D bounding boxes. In order to avoid the need of 3D bounding boxes, SUDS [17] introduces a scalable hash table to represent large-scale dynamic urban scenes, using an off-the-shelf 2D optical flow estimator to track dynamic objects. EmerNerf [24] addresses this challenge by learning scene flow and using it to link corresponding points in the 4D neural field across multi-frame reconstruction, allowing the separation of static and dynamic objects without 3D bounding boxes. However, these methods, which rely on implicit representations, still suffer from inefficiency in both construction and rendering.

Recently, 3DGS [6] introduced a novel explicit 3D scene representation, combining high-quality volume rendering with a fast speed. However, 3DGS is designed for static scenes and fails when modeling dynamic moving objects. To address this, DrivingGaussian [27] and StreetGaussian [23] partition the driving scene into static and dynamic components by separate sets of Gaussians based on vehicles’ 3D bounding boxes. DrivingGaussian [27] uses a composite dynamic Gaussian graph to manage multiple moving objects while incrementally employing incremental static 3D Gaussians to represent the entire background. StreetGaussian [23] models dynamic urban scenes as point clouds with semantic logits and 3D Gaussians, handling foreground vehicles and background separately. HUGS [26] also adopts 3D object boxes to identify dynamic elements from background, optimizing geometry, appearance, semantics, and motion specifically of dynamic Gaussians. However, all these methods rely on expensive annotated or predicted 3D bounding boxes to learn the time-dependent representations of dynamic objects. Most recent

method PVG [1] introduces periodic vibration based Gaussian attributes, optimized through self-supervision without the need of 3D bounding boxes. Each Gaussian models dynamic changes over time through optimizable attributes including vibration directions, life span, and life peak.

### 3. Methodology

#### 3.1. Overview

By collecting images and point clouds with timestamps from surrounding cameras and LiDAR, we aim to learn a space-time 4D Gaussian representation of a dynamic scene without any human annotations, enabling fast and high-quality novel view rendering. As the pipeline of SplatFlow shown in Fig. 2, NMFF models the temporal motions of both 3D points and Gaussians as continuous motion flow fields. It serves two key roles in SplatFlow: 1) decomposing static and dynamic elements from 3D LiDAR points; 2) modeling the temporal status correspondences of each 4D Gaussian across time. We first decompose the time-series LiDAR points into static and dynamic points by NMFF, which are then separately merged to initialize the static and dynamic Gaussians. Dynamic objects are represented by aggregated 4D Gaussians associated with NMFF over time, while the static background is modeled using 3D Gaussians. NMFF learns the status correspondence of each 4D Gaussian in conjunction with optimizing 4D Gaussians’ attributes along time. We also distill optical flow knowledge from a 2D foundational model into 4D space-time representation. Finally, images, depths, and optical flows are rendered from novel viewpoints at different timesteps in the dynamic driving scenario.

#### 3.2. Problem Definition:

In a dynamic urban scenario, we collect time-sequential data from a vehicle equipped with surrounding cameras and a LiDAR sensor. A set of surrounding images  $\mathcal{I}$  is captured

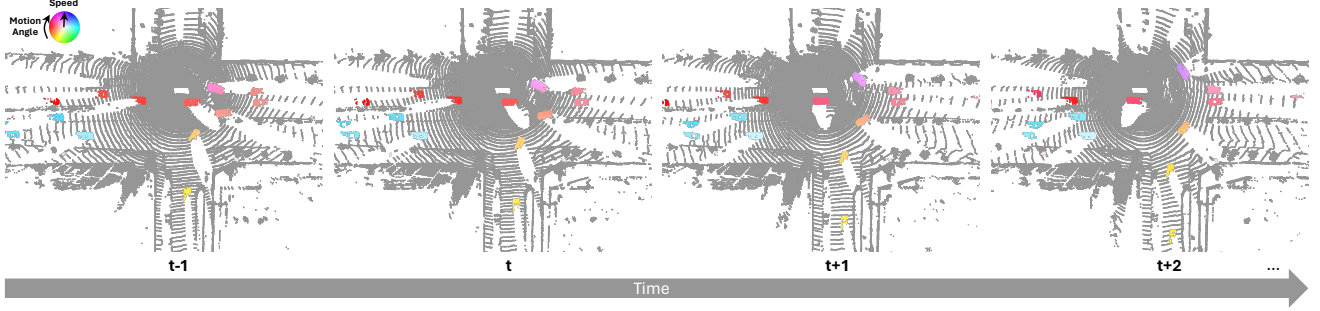


Figure 3. Visualization of 3D LiDAR points within NMFF on Waymo dataset.

by multiple surrounding cameras with corresponding intrinsic matrices  $I$  and extrinsic matrices  $E$ . Meanwhile, a set of 3D points  $\mathcal{P}$  is obtained from LiDAR, along with corresponding extrinsic matrices  $E'$ . We denote the synchronized and calibrated multi-sensor data sequence, including timestamps  $t$ , as  $\{\mathcal{I}_i, \mathcal{P}_i, E_i, I_i, E'_i, t_i | i = 1, 2, 3, \dots, n\}$ , where  $\mathcal{P}_i = \{x_i, y_i, z_i\}$  and  $n$  is the number of frames. The vehicle trajectory is either provided or estimated using multi-sensor-based odometry. Our model  $\mathcal{SF}$  aims to perform accurate 3D reconstruction and synthesize novel viewpoints for any given timestamp  $t$  and camera pose  $[E_t, I_t]$  by rendering  $\hat{\mathcal{I}} = \mathcal{SF}(E_t, I_t, t)$ .

### 3.3. 4D Gaussian Representation

Given images with associated camera poses, 3DGS optimizes a set of anisotropic 3D Gaussians through differentiable rasterization to represent a static 3D scene. We extend 3DGS to 4DGS of spatial-temporal representations for dynamic objects in an urban scene. Each 4D Gaussian primitives  $\mathcal{G}(t)$  is represented by time-varying attributes: 3D center  $\mu(t) = [x(t), y(t), z(t)]^T$  and covariance matrix  $\Sigma(t)$  along with time-invariant attributes: opacity  $\sigma$  and color  $c$ . The density of a Gaussian at point  $(x, t)$  is defined as,

$$\alpha(t) = \sigma \cdot \exp\left(-\frac{1}{2}(x - \mu(t))^T \Sigma(t)^{-1} (x - \mu(t))\right). \quad (1)$$

The covariance matrix  $\Sigma(t) = R(t)SS^T R(t)^T$  is composed of a scaling matrix  $S = \text{diag}(s_x, s_y, s_z)$  and rotation matrix  $R(t) = (q_x(t), q_y(t), q_z(t), q_w(t))$ , constrained to be a positive semi-definite matrix during optimization.

To render an image from a specific viewpoint at timestamp  $t$ , we first transform all the 4D Gaussians from other times to the target time  $t$ , according to the learned correspondence from NMFF. Subsequently, the aggregated 4D Gaussians are splatted onto the image plane, resulting in a collection of 2D Gaussians. The 3D covariance matrix  $\Sigma$  is projected to a 2D covariance matrix  $\Sigma'$  by,

$$\Sigma'(t) = JE\Sigma(t)E^T J^T, \quad (2)$$

where  $E$  refers to the world-to-camera matrix and  $J$  refers to Jacobian of the perspective transformation.

By sorting the Gaussians according to their depth within the camera space, we use  $\alpha$ -blending to estimate the color  $C$  and depth  $D$  at each pixel  $p$  as,

$$C = \sum_{i=1}^N c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), D = \sum_{i=1}^N z_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j). \quad (3)$$

Here  $c_i$  represents the color of Gaussian  $\mathcal{G}_i$ , which is computed using spherical harmonic coefficients.  $z_i$  represents the distance from the image plane to the  $i$ -th Gaussian center.  $\alpha_i$  denotes the density computed from the 2D projection of the  $i$ -th Gaussian and its learned opacity.

### 3.4. Neural Motion Flow Field

To establish the correspondence of dynamic objects, we design NMFF to model 3D motions as a continuous motion flow field, enabling the temporal transition of both 3D points and Gaussians. The NMFF  $\Phi$  contains a temporal sequence of motion flow field, defined as  $\Phi = \{\phi_{t_i:t_{i+1}}\}$ ,  $i = 0, 1, 2, \dots, n-1$ , which predicts 3D motion flow of an arbitrary query point between two consecutive frames as,

$$\phi_{t_1:t_2}(x_{t_1}, y_{t_1}, z_{t_1}) = \Delta x_{t_1:t_2}, \Delta y_{t_1:t_2}, \Delta z_{t_1:t_2}, \Delta R_{t_1:t_2} \quad (4)$$

where  $\Delta R_{t_1:t_2}$  denotes motion angle between two adjacent timestamps. Each field  $\phi$  is eight ReLU-MLP stacks.

**Point Cloud within NMFF:** Although NMFF can be learned during Gaussian Splatting optimization, we observed that the rendering loss alone is insufficient to effectively constrain the motion representation of dynamic scene components. Therefore, we leverage the rich geometric structure inherent in temporally consecutive 3D LiDAR points to derive a robust NMFF prior, as shown in Fig. 3.

We first pre-train  $\Phi$  on a sequence of 3D LiDAR points by forward and backward 3D geometry consistency, enabling NMFF to learn the 3D motion of each query 3D point over time. Given the source point clouds  $\mathcal{P}_{t_1}$  at time  $t_1$  and target point clouds  $\mathcal{P}_{t_2}$  at time  $t_2$ , we minimize the point distance between the source and target point clouds by a bidirectional Chamfer Distance (CD),

$$CD(\mathcal{P}_{t_1}, \mathcal{P}_{t_2}) = \sum_{p_{t_2} \in \mathcal{P}_{t_2}} D(p_{t_2}, \mathcal{P}_{t_1}) + \sum_{p_{t_1} \in \mathcal{P}_{t_1}} D(p_{t_1}, \mathcal{P}_{t_2}), \quad (5)$$

where  $D$  refers to point distance function, where correspondences from source-to-target and target-to-source are searched among the nearest point neighbors.

We generate a 3D dynamic mask to separate static and dynamic points through ego-motion compensation on the

scene flow predicted by NMFF. Static and dynamic point clouds are then merged separately using ego-motion and scene flow. For the static background, we initialize the static Gaussians using the merged 3D static points. For dynamic objects, we aggregate points from different timestamps into a common reference frame using scene flow predictions. This aggregation densifies and completes dynamic point clouds from various viewpoints over time.

**4D Gaussian within NMFF:** The pre-trained  $\Phi$  is subsequently integrated with 4D Gaussian representations of dynamic objects, enabling joint optimization during Gaussian splatting. Given consecutive time  $t_1$  and  $t_2$  where  $t_1 < t_2$  with respective states  $\{\mathcal{G}_i(t_1)\}$  and  $\{\mathcal{G}_i(t_2)\}$ , these states are connected by NMFF for each Gaussian as,

$$\mathcal{G}_i(t_1) = \{\mu(t_1), R(t_1), S, \alpha, c\}, \quad (6)$$

$$\Delta\mu_{t_1:t_2}, \Delta R_{t_1:t_2} = \phi_{t_1:t_2}(\mu(t_1)), \quad (7)$$

$$\hat{\mathcal{G}}_i(t_2) = \{\mu(t_1) + \Delta\mu_{t_1:t_2}, R(t_1) \cdot \Delta R_{t_1:t_2}, S, \alpha, c\}. \quad (8)$$

We employ NMFF  $\Phi$  correspondence to warp the center of all 4D Gaussians across time to the desired timestamp, as  $\{\mathcal{G}_i(t)\} = \{\mathcal{G}(t), \hat{\mathcal{G}}(t_1), \hat{\mathcal{G}}(t_2), \dots, \hat{\mathcal{G}}(t_n)\}$ . For the transition of Gaussians across multiple timestamps, we apply the corresponding NMFF through step-by-step propagation. Here NMFF learns a motion transition of the 4D Gaussian among training frames in a self-supervised manner, promoting a multi-view consistent representation.

For static background, we represent it by static 3D Gaussians  $\mathcal{G}_{static}$  following the strategy [1]. For sky area, we adopt a separate environmental map with sky Gaussians  $\mathcal{G}_{sky}$  following [1, 24]. Given a specific viewpoint  $[E_t, I_t, t]$ , the image  $\hat{\mathcal{I}}$  and depth  $\hat{\mathcal{D}}$  at time  $t$  are rendered by the differentiable rendering using a set of Gaussians,

$$\hat{\mathcal{I}}, \hat{\mathcal{D}} = \text{Render}(\{\mathcal{G}_i(t)\}, \mathcal{G}_{static}, \mathcal{G}_{sky} | E_t, I_t, t). \quad (9)$$

### 3.5. Optical Flow Distillation

The 3D motion flow provided by NMFF enables 2D optical flow rendering. Given two consecutive time  $t_1$  and  $t_2$ , the optical flow  $f_{t_1:t_2}$  of each Gaussian is computed by projecting 3D center  $\mu$  to the image plane using the camera’s intrinsic matrix  $I$  and extrinsic matrix  $E$ ,

$$\mu'(t_1) = I[E]\mu(t_1), \quad \mu'(t_2) = I[E]\mu(t_2), \quad (10)$$

$$f_{t_1:t_2} = \mu'(t_2) - \mu'(t_1). \quad (11)$$

Then we estimate the optical flow  $F$  at pixel  $p$  via by point-based  $\alpha$ -blending as,

$$F = \sum_{i=1}^N f_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j). \quad (12)$$

The optical flow  $\hat{\mathcal{F}}$  at time  $t$  are rendered by the differentiable rendering similar as Equation 9. In order to distill optical flow knowledge from 2D foundation model into 4D

space, we adopt SEA-RAFT[18] to extract optical flow  $\mathcal{F}$  as pseudo ground truth for distillation by optical flow loss,

$$\mathcal{L}_F = \lambda_f \mathcal{L}_1(\mathcal{F}, \hat{\mathcal{F}}) + (1 - \lambda_f) \mathcal{L}_1(\mathcal{I}_{next}, \mathcal{T}(\hat{\mathcal{I}} | \hat{\mathcal{F}})), \quad (13)$$

where  $\lambda_f$  is a scale factor, and  $\mathcal{I}_{next}$  denotes the image at the next time step.  $\mathcal{T}$  is used to warp rendered image to next time step according to optical flow. The optical flow  $\mathcal{F}$  extracted from foundational model is only required during training and is not needed during inference.

### 3.6. Optimization

All Gaussian attributes, along with the parameters of NMFF, are optimized end-to-end in a self-supervised manner. Meanwhile, adaptive densification and pruning strategies are introduced to enhance the fitting of the 3D scene. The overall training loss is given as,

$$\mathcal{L} = \mathcal{L}_I + \lambda_1 \mathcal{L}_D + \lambda_2 \mathcal{L}_F + \lambda_3 \mathcal{L}_{sky} + \lambda_4 \mathcal{L}_{reg}, \quad (14)$$

where  $\lambda_{1,2,3,4}$  are scale factors and  $\mathcal{L}_{reg}$  is the regularization term. The image loss  $\mathcal{L}_I$  combines L1 and SSIM losses between rendered and observed images,

$$\mathcal{L}_I = (1 - \lambda_{ssim}) \mathcal{L}_1(\mathcal{I}, \hat{\mathcal{I}}) + \lambda_{ssim} \mathcal{L}_{ssim}(\mathcal{I}, \hat{\mathcal{I}}), \quad (15)$$

where  $\lambda_{ssim}$  is a scale factor. The depth loss  $\mathcal{L}_D$  is a L1 loss between inverse of rendered depth and generated depth by projecting sparse LiDAR points onto camera plane as,

$$\mathcal{L}_D = \mathcal{L}_1(\mathcal{D}, \hat{\mathcal{D}}). \quad (16)$$

The sky opacity loss  $\mathcal{L}_{sky}$  is a binary cross entropy loss for sky refining using sky mask  $M_{sky}$  from SegFormer [20],

$$\mathcal{L}_{sky} = \mathcal{L}_{bce}(O, 1 - M_{sky}), \quad (17)$$

where the accumulated opacity denotes as  $O = \sum_{i=1}^N \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j)$ .

## 4. Experiments

### 4.1. Datasets, Metrics and Baselines

Following PVG [1], we validate our approach on two widely-used datasets in autonomous driving: Waymo Open Dataset [15] and KITTI Dataset [4]. Both offer multi-sensor data clips including synchronized and calibrated camera and LiDAR data from urban driving scenarios. We train and test our methods, along with all baselines, using full-resolution images:  $1920 \times 1280$  for the Waymo dataset and  $1242 \times 375$  for the KITTI dataset.

Following prior research [1, 17, 23, 24, 27], we evaluate image reconstruction and novel view synthesis using three widely-adopted benchmark metrics: peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM), and learned perceptual image patch similarity (LPIPS).

We compare our method with extensive baselines including S-Nerf [21], StreetSurf [5], 3DGS [6], NSG [10], Mars [19], SUDS [17], EmerNerf [24], PVG [1], StreetGaussian [23] on the Waymo and KITTI benchmarks. More implementation details and results are provided in the supplementary material.



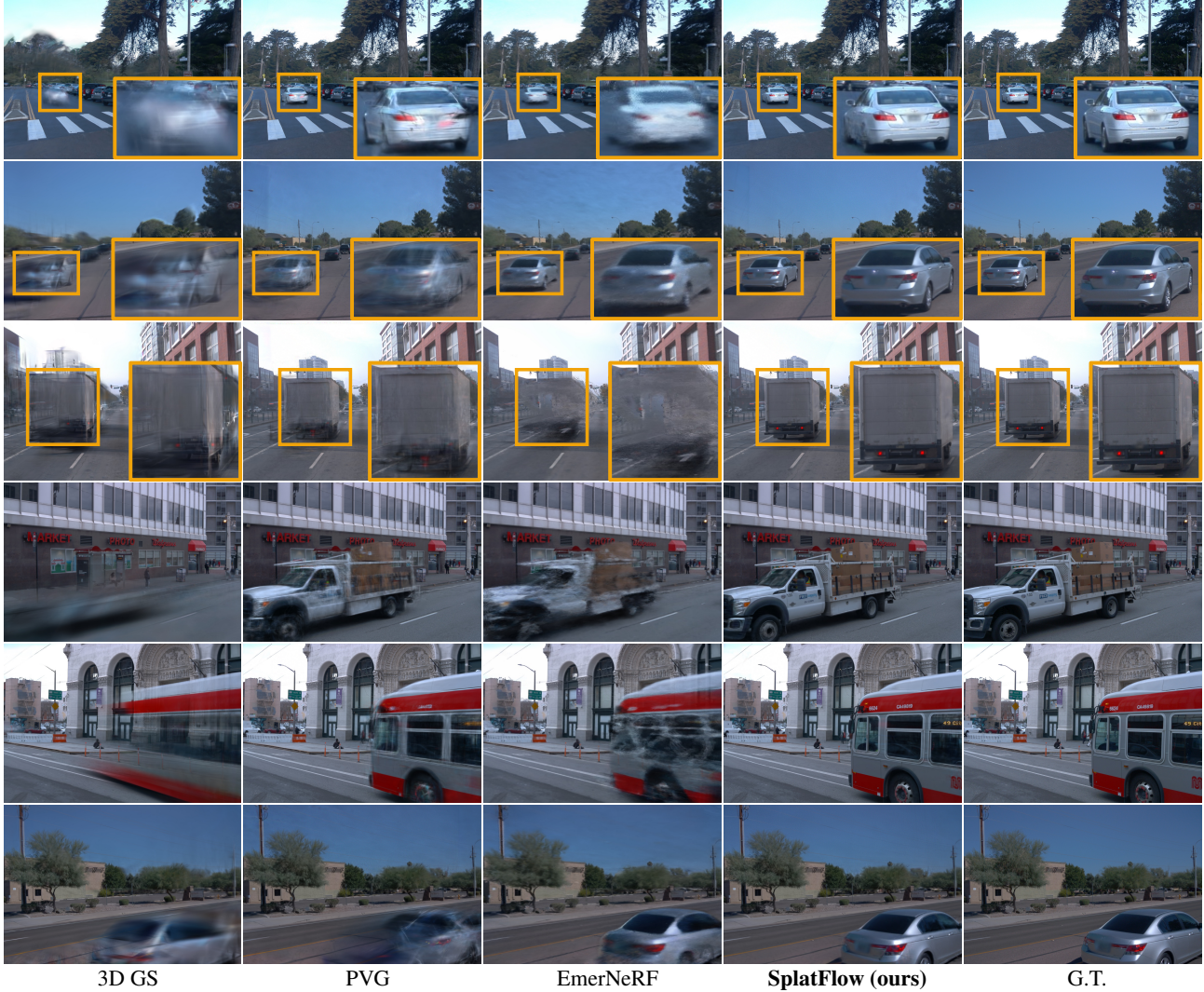


Figure 4. Visual comparison of novel view synthesis on Waymo dataset. Bounding boxes indicate the zoomed-in dynamic areas.

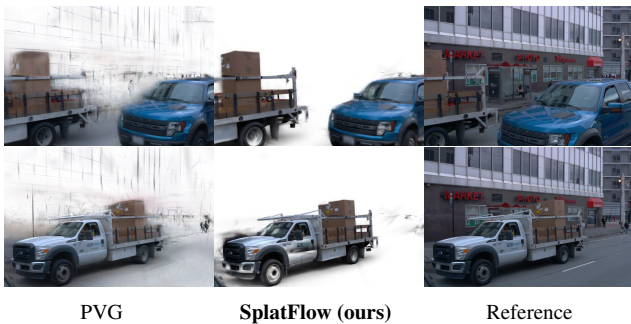


Figure 5. Dynamic object decomposition comparison on Waymo.

## 4.2. Evaluation on Waymo Open Dataset

Following PVG [1], we evaluated our SplatFlow against baselines using the Waymo Open dataset. Table 1 summarizes the average metrics for both image reconstruction and novel view synthesis tasks on selected dynamic scenes. SplatFlow consistently outperforms all baselines across all evaluated metrics for both tasks. In image reconstruction,

	Image Reconstruction			Novel View Synthesis		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
S-NeRF [21]	19.67	0.528	0.387	19.22	0.515	0.400
StreetSurf [5]	26.70	0.846	0.387	23.78	0.822	0.401
3DGS [6]	27.99	0.866	0.3717	25.08	0.822	0.319
NSG [10]	24.08	0.656	0.293	21.01	0.571	0.487
Mars [19]	21.81	0.681	0.441	20.69	0.636	0.453
SUDS [17]	28.83	0.805	0.430	21.83	0.656	0.405
EmerNeRF [24]	28.11	0.786	0.289	25.92	0.763	0.384
PVG [1]	32.46	0.910	0.373	28.11	0.849	0.279
<b>SplatFlow</b>	<b>33.64</b>	<b>0.951</b>	<b>0.198</b>	<b>28.71</b>	<b>0.874</b>	<b>0.239</b>

Table 1. Performance comparison on Waymo dataset.

SplatFlow achieves high image quality with the highest scores: PSNR at 33.64, SSIM at 0.951, and LPIPS at 0.198. For novel view synthesis, SplatFlow produces high-quality renderings of unseen timestamps, reaching PSNR of 28.71, SSIM of 0.874 and LPIPS of 0.239. These improvements are visually validated in Fig. 4 and Fig. 9 (dynamic only), showing exceptional clarity in details for both static and dynamic regions. Compared with baseline methods, which show motion artifacts such as ghosting and blur in regions





Figure 6. Visual comparison of novel view synthesis on KITTI dataset. Bounding boxes indicate the zoomed-in dynamic areas.

	PVG [1]			SplatFlow		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Seg. 1058...	26.86	0.840	0.273	<b>27.46</b>	<b>0.880</b>	<b>0.253</b>
Seg. 2259...	26.45	0.824	0.309	<b>29.30</b>	<b>0.856</b>	<b>0.290</b>
Seg. 7670...	26.58	0.820	0.307	<b>28.89</b>	<b>0.855</b>	<b>0.283</b>
Seg. 5083...	28.35	0.893	0.213	<b>30.29</b>	<b>0.928</b>	<b>0.169</b>
	PSNR* $\uparrow$	SSIM* $\uparrow$	LPIPS* $\downarrow$	PSNR* $\uparrow$	SSIM* $\uparrow$	LPIPS* $\downarrow$
Seg. 1058...	27.42	0.988	0.020	<b>28.50</b>	<b>0.992</b>	<b>0.018</b>
Seg. 2259...	22.55	0.981	0.031	<b>31.61</b>	<b>0.995</b>	<b>0.016</b>
Seg. 7670...	23.56	0.968	0.042	<b>28.76</b>	<b>0.975</b>	<b>0.029</b>
Seg. 5083...	25.58	0.964	0.057	<b>26.71</b>	<b>0.975</b>	<b>0.040</b>

Table 2. Novel view synthesis results on Waymo dataset (\*denotes dynamic elements only).

	Image Reconstruction			Novel View Synthesis		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
S-NeRF [21]	19.23	0.664	0.193	18.71	0.606	0.352
StreetSurf [5]	24.14	0.819	0.257	22.48	0.763	0.304
3DGS [6]	21.02	0.811	0.202	19.54	0.776	0.224
NSG [10]	26.66	0.806	0.186	21.53	0.673	0.254
Mars [19]	27.96	0.900	0.185	24.23	0.845	0.160
SUDS [17]	28.31	0.876	0.185	22.77	0.797	0.171
EmerNeRF [24]	26.95	0.828	0.218	25.24	0.801	0.237
PVG [1]	32.83	0.937	0.070	27.43	0.896	0.114
<b>SplatFlow</b>	<b>33.37</b>	<b>0.943</b>	<b>0.057</b>	<b>28.32</b>	<b>0.932</b>	<b>0.089</b>

Table 3. Performance comparison on KITTI dataset.

of dynamic elements, our method preserves high-quality textures and fine details, producing more accurate renderings at novel viewpoints.

To further evaluate the rendering quality of dynamic objects against the primary baseline PVG, we selected additional four scenes featuring a higher presence of dynamic objects for further comparison. Table 2 displays the metrics for novel view synthesis, evaluated both on entire scenes and exclusively on dynamic objects (denoted as \*). SplatFlow consistently surpasses PVG in all selected scenes, demonstrating higher PSNR and SSIM values and lower

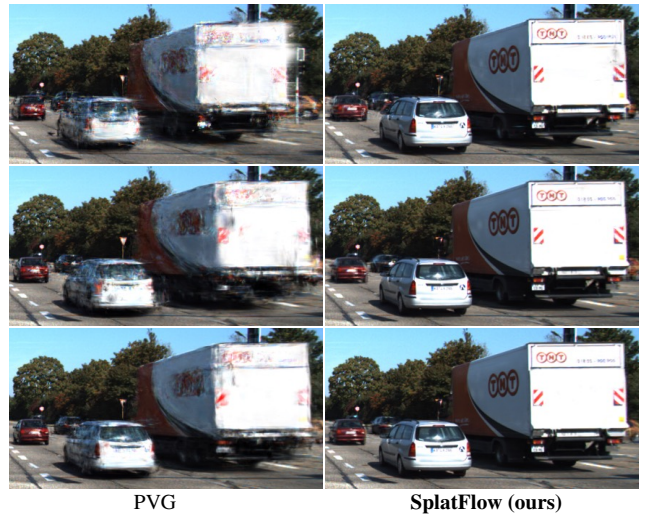


Figure 7. Visual comparison of novel view synthesis on KITTI-25% (row1), 50% (row2), and 75% (row3) dataset.

LPIPS scores, underscoring its effectiveness in handling complex, dynamic content. The dynamic objects masks are obtained from 2D ground truth bounding boxes following [24]. The PSNR\* of dynamic regions on novel view synthesis significantly outperforms PVG both quantitatively and qualitatively, especially in scenes with fast ego-motion and moving objects in Seg. 2259 and Seg. 7670. We also present the visual comparison of dynamic object decomposition with PVG in Fig. 5, which illustrates that SplatFlow renders sharper and clearer image of dynamic objects.

### 4.3. Evaluation on KITTI Dataset

Following PVG [1], we further evaluated our SplatFlow method in comparison to the baselines on KITTI dataset.

	KITTI - 75%			KITTI - 50%			KITTI - 25%		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
NeRF [9]	18.56	0.557	0.554	19.12	0.587	0.497	18.61	0.570	0.510
NSG [10]	21.53	0.673	0.254	21.26	0.659	0.266	20.00	0.632	0.281
SUDS [17]	22.77	0.797	0.171	23.12	0.821	0.135	20.76	0.747	0.198
MARS [19]	24.23	0.845	0.160	24.00	0.801	0.164	23.23	0.756	0.177
3DGS [6]	19.19	0.737	0.172	19.23	0.739	0.174	19.06	0.730	0.180
PVG [1]	27.43	0.896	0.114	25.92	0.882	0.114	22.55	0.833	0.151
StreetGS [23]	25.79	0.844	<b>0.081</b>	25.52	0.841	<b>0.084</b>	24.53	0.824	<b>0.090</b>
<b>SplatFlow</b>	<b>28.32</b>	<b>0.932</b>	0.089	<b>27.90</b>	<b>0.927</b>	0.093	<b>26.10</b>	<b>0.890</b>	0.120

Table 4. Performance comparison of novel view synthesis on KITTI-75%, 50% and 25% dataset.



Figure 8. Visual comparison of ablation study on Waymo dataset.

Table 3 presents the average metrics for both image reconstruction and novel view synthesis tasks on selected dynamic scenes characterized by extensive movement. SplatFlow outperforms all other methods across all metrics in both tasks, with particularly high scores in PSNR, SSIM and LPIPS, achieving 33.37, 0.943 and 0.057 for image reconstruction and 28.32, 0.932 and 0.089 for novel view synthesis. SplatFlow achieves exceptional results in novel view synthesis, as shown in Fig. 6. Our approach generates high-fidelity renderings that preserve fine details while accurately reconstructing dynamic areas where existing approaches struggle to maintain stability.

To evaluate robustness to training data size, we follow the settings of [19, 23] to evaluate our method using different train/test split configurations on the KITTI dataset. Table 4 and Fig. 7 provide the performance comparisons and visualizations. SplatFlow consistently outperforms all other methods across all dataset subsets, achieving the highest scores. Even with reduced training data, SplatFlow robustly delivers high-quality novel view synthesis across varying levels of data availability, highlighting its efficient use of available information and strong generalization capabilities.

#### 4.4. Ablation Study

To demonstrate the effectiveness of each component in SplatFlow, we performed ablation studies on the Waymo dataset. Our main contributions, NMFF prior, NMFF opti-

	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
w/o NMFF prior	27.69	0.863	0.282
w/o NMFF optimization	28.14	0.874	0.269
w/o optical flow distillation	28.28	0.877	0.252
<b>Full</b>	<b>28.99</b>	<b>0.880</b>	<b>0.249</b>
	PSNR* $\uparrow$	SSIM* $\uparrow$	LPIPS* $\downarrow$
w/o NMFF prior	27.28	0.963	0.056
w/o NMFF optimization	27.97	0.975	0.036
w/o optical flow distillation	28.51	0.980	0.028
<b>Full</b>	<b>28.90</b>	<b>0.984</b>	<b>0.026</b>

Table 5. Ablation study on Waymo dataset (\*denotes dynamic elements only).



Figure 9. Detail comparison of novel view synthesis on Waymo. EmerNeRF, PVG, SplatFlow, and G.T. optimization, and optical flow distillation are analyzed in this study to assess their individual impact. We trained four different variations: 1) SplatFlow without pretraining on LiDAR data; 2) SplatFlow without optimization with 4D Gaussians on image data; 3) SplatFlow without optical flow distillation from foundational model; 4) SplatFlow full method. Table 5 presents the average evaluation metrics across test scenes, showing results for both entire scenes and dynamic objects (indicated as \*). Ablation study validates the contribution of NMFF prior, NMFF optimization, optical flow distillation to overall performance.

## 5. Conclusion

In this paper, we propose SplatFlow, a novel self-supervised Dynamic Gaussian Splatting within Neural Motion Flow Fields (NMFF) for accurate reconstruction and real-time rendering in dynamic urban scenarios. The core idea of SplatFlow is to seamlessly integrate time-dependent 4D Gaussian representations within NMFF in a unified framework, where NMFF implicitly models the motions of dynamic components across time through self-supervision. Experimental evaluations show that SplatFlow outperforms state-of-the-art across all standard metrics on the Waymo Open and KITTI datasets, without requiring expensive annotations of dynamic object detection and tracking.



# SplatFlow: Self-Supervised Dynamic Gaussian Splatting in Neural Motion Flow Field for Autonomous Driving

## Supplementary Material

In this supplementary material, we first provide more implementation details of SplatFlow in Appendix 6. Second, we provide more visualization of 3D LiDAR points within NMFF in Appendix 7. Third, we present more detailed comparison visualizations of dynamic object synthesis from novel views in Appendix 8. Fourth, we include visualizations of rendered RGB image, optical flow and depth in Appendix 9. Fifth, runtime performance comparisons are provided in Appendix 10. Finally, video demonstrations are included in Appendix 11.

### 6. Implementation Details

Each field in the Neural Motion Flow Field (NMFF) consists of eight ReLU-MLP stacks. All MLPs are followed by a ReLU activation, except for the final prediction layer, where the middle hidden dimensions are configured as 128.

For NMFF pre-training, we follow the approach in [12] to generate pseudo scene flow labels, excluding ground points from the Waymo and KITTI datasets. The raw 3D LiDAR points are utilized without cropping to a smaller range. We use a learning rate of  $8e-3$  with the Adam optimizer, optimizing each scene for up to 4000 iterations with early stopping. Additionally, point cloud densification is performed by accumulating point clouds through Euler integration, using per-pair scene flow estimations.

During the 4D Gaussian with NMFF optimization, we configure the position learning rate to a range from  $1.6e-5$  to  $1.6e-6$ , the opacity learning rate to 0.05, the scale learning rate to 0.005, the feature learning rate to  $2.5e-3$ , and the rotation learning rate to 0.001. The intervals for densification and opacity reset are set to 500 and 3000, respectively. We set the densify gradient threshold for decomposed static and dynamic Gaussians as  $1.7e-4$  and  $1e-4$ , respectively. The Spherical Harmonics degree for each Gaussian is set to 3. For NMFF optimization, we set the learning rate to  $1e-4$ . For training losses, we use coefficients  $\lambda_1 = 0.1$ ,  $\lambda_2 = 0.005$ ,  $\lambda_3 = 0.05$ ,  $\lambda_4 = 0.001$ ,  $\lambda_{sim} = 0.2$  and  $\lambda_f = 0.8$ .

### 7. Visualization of LiDAR Points in NMFF

We provide more visualization of 3D LiDAR points within NMFF on the Waymo dataset in Fig. 10. The color wheel in the center represents the flow magnitude through color intensity and the flow direction via angle. As illustrated, NMFF accurately predicts the 3D motion flow of 3D LiDAR points for dynamic objects in driving scenarios.

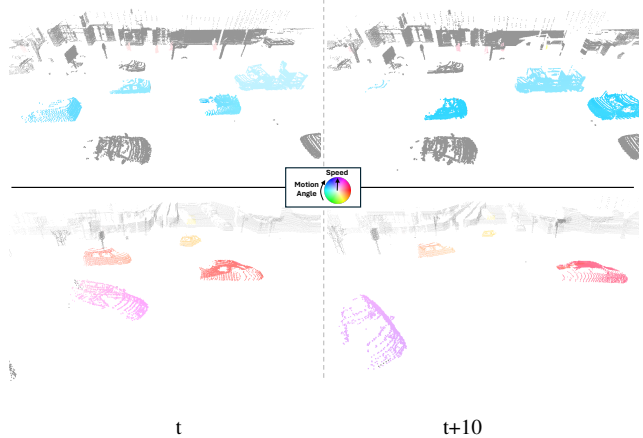


Figure 10. Visualization of 3D LiDAR points within NMFF on Waymo dataset.

	Waymo FPS	KITTI FPS
S-NeRF [21]	0.0014	0.0075
StreetSurf [5]	0.097	0.37
NSG [10]	0.032	0.19
Mars [19]	0.030	0.31
SUDS [17]	0.008	0.04
EmerNerf [24]	0.053	0.28
3DGS [6]	<b>63</b>	<b>125</b>
PVG [1]	50	59
<b>SplatFlow</b>	40	44

Table 6. The comparison running-time analysis on Waymo and KITTI datasets.

### 8. Detailed Visual Comparison

We present more visual comparison details of dynamic object synthesis from novel views, showcasing results on the Waymo dataset in Fig. 11 and on the KITTI dataset in Fig. 12. As shown, SplatFlow generates sharper images with fewer blurred artifacts, particularly for high-speed vehicles, compared to the baselines.

### 9. Rendered Depth and Flow Visualization

We also present visualization of rendered RGB image, optical flow, and depth on the Waymo dataset in Fig. 13 and 14, and on the KITTI dataset in Fig. 15. In these Figures, the first row displays the rendered RGB images, the second row shows the rendered optical flow, and the third row presents the rendered depth, all generated by SplatFlow. As shown, our SplatFlow renders sharp, clear, and dense optical flow

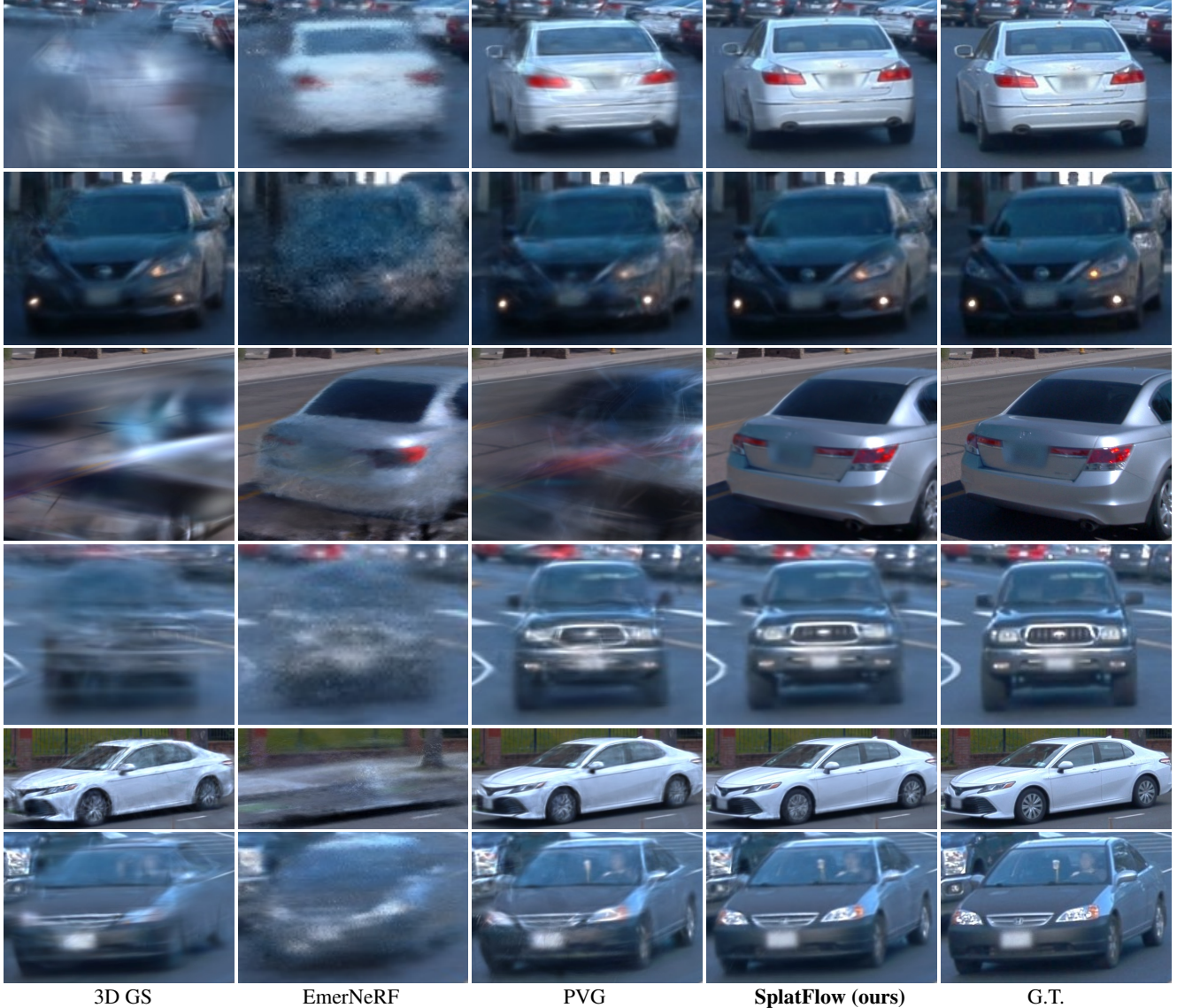


Figure 11. Detailed comparison of dynamic object synthesis from novel views on Waymo dataset.

and depth images in dynamic driving scenarios.

## 10. Running-time Analysis

We compare the runtime performance of SplatFlow with various baseline methods on the Waymo and KITTI datasets, as summarized in Table 6. Utilizing a single NVIDIA GeForce A6000, SplatFlow achieves real-time rendering speeds for high-resolution images after quantization and pruning optimization, delivering approximately 40 FPS at  $1920 \times 1280$  resolution on the Waymo dataset and around 44 FPS at  $1242 \times 375$  resolution on the KITTI dataset. Compared to NeRF-based methods such as S-NeRF [21], StreetSurf [5], NSG [10], SUDS [17], EmerNerf [24], SplatFlow significantly surpasses the speed of

these methods, delivering real-time performance. Compared to GS-based methods such as 3DGS [6] and PVG [1], SplatFlow achieves higher accuracy in dynamic object rendering while maintaining efficient performance.

## 11. Video Demos

We include five video demos in the supplementary material according to size limitation.

Demos 1 and 2 showcase the results of image synthesis from novel views produced by SplatFlow, alongside baseline methods and ground truth (G.T.) data, in dynamic driving scenarios from the Waymo and KITTI datasets respectively. In these videos, the first and second rows display the surrounding images rendered by the baselines: 3D-GS [6],



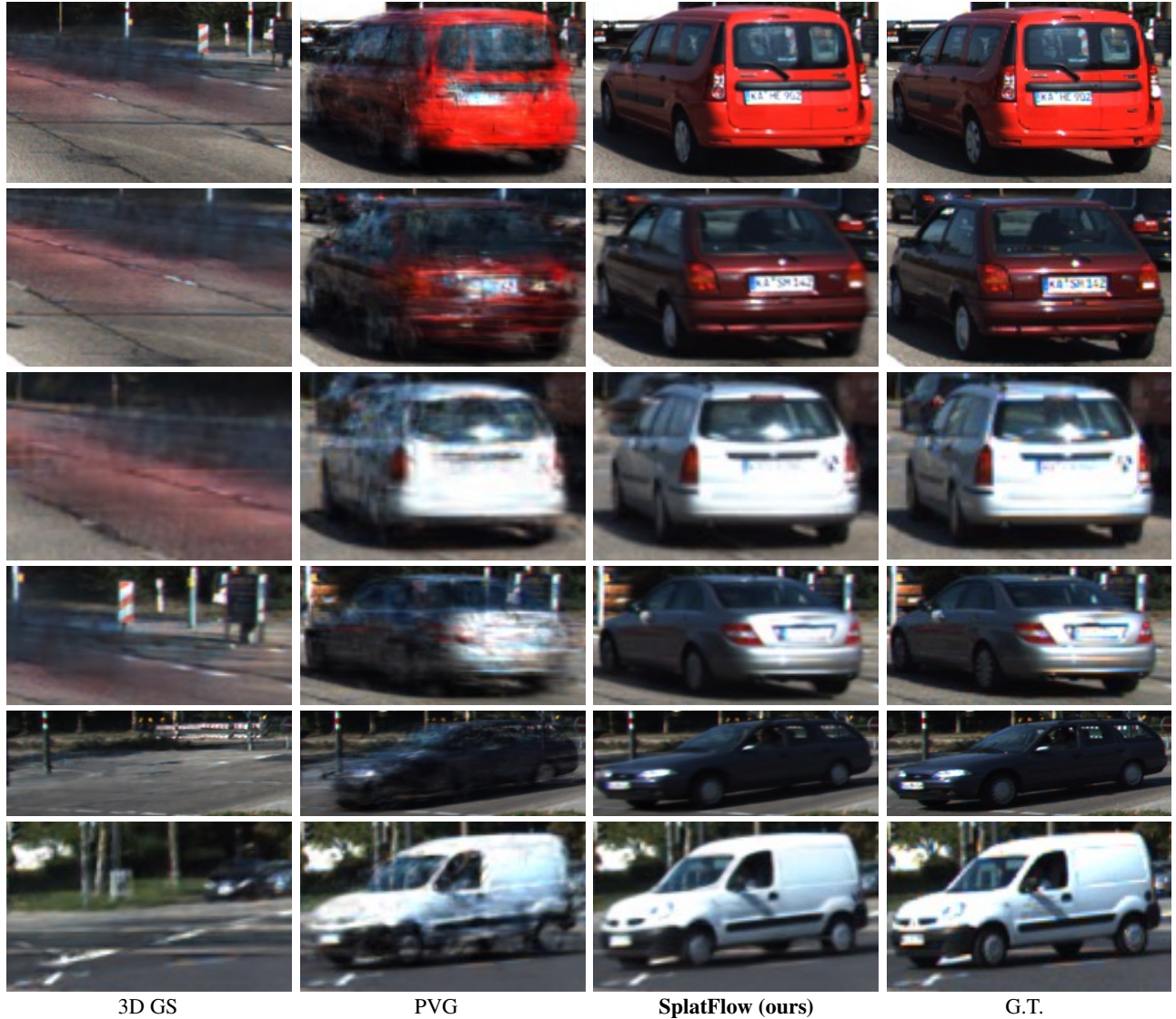


Figure 12. Detailed comparison of dynamic object synthesis from novel views on KITTI dataset.

PVG [1], or EmerNeRF [24]. The third row presents the surrounding images rendered by our SplatFlow, while the final row shows the G.T. surrounding images. To provide a clearer comparison in the visualization, video demos 1 and 2 are played at 0.1x speed.

Demos 3 and 4 showcase the rendered images, optical flow, and depth produced by our SplatFlow in dynamic driving scenarios from the Waymo and KITTI datasets. In these videos, the first row shows the G.T RGB images. The second row displays the rendered RGB images, the third row shows the rendered optical flow, and the fourth row presents the rendered depth, all generated by our SplatFlow. For a clearer visualization, video demos 3 and 4 are played at 0.1x speed.

Demo 5 showcases the 3D motion prediction of LiDAR

points within the NMFF. The color wheel in the top right corner represents the flow magnitude through color intensity and flow direction via the angle. For clearer visualization, video Demo 5 is played at 0.1x speed.



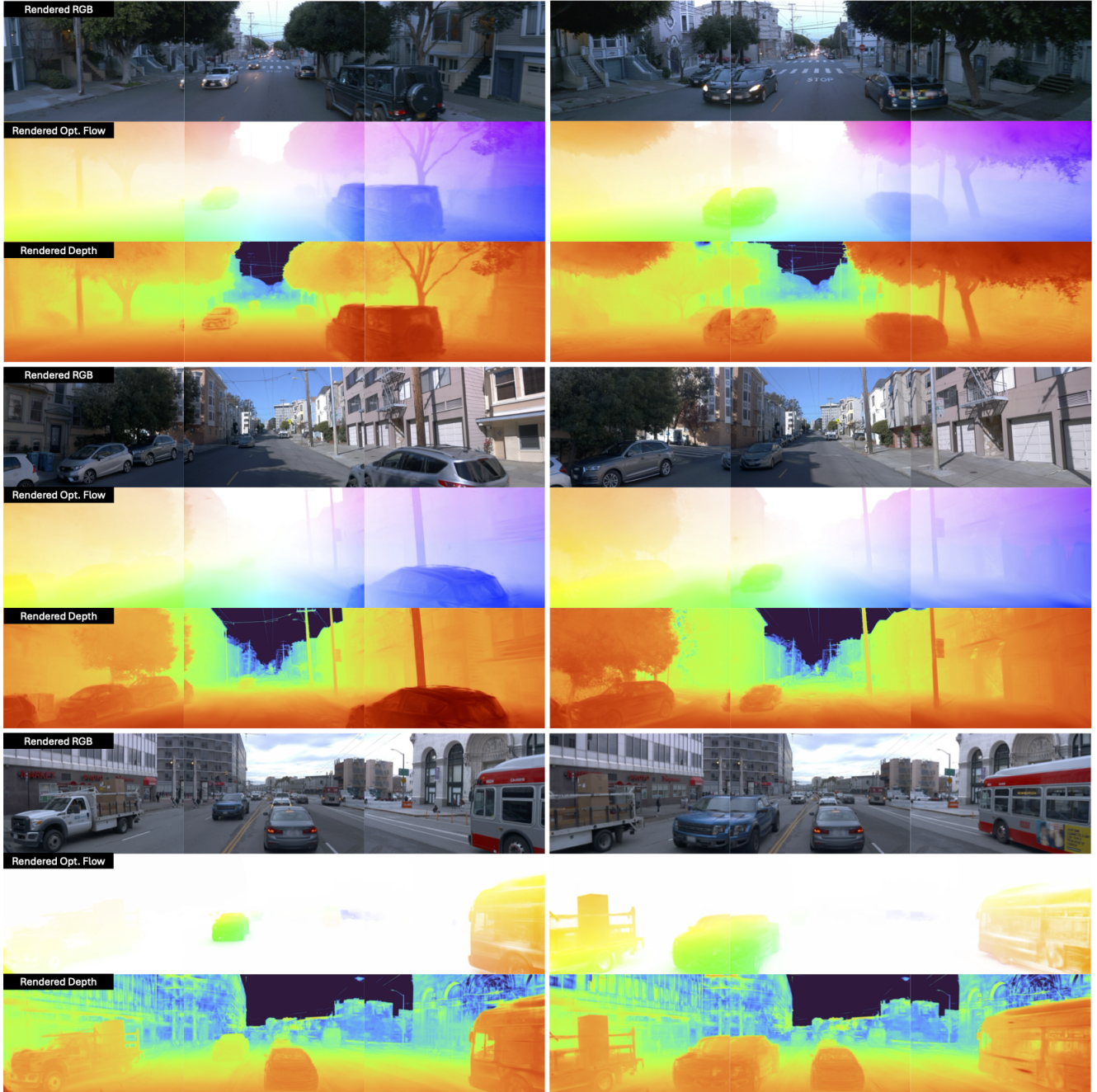


Figure 13. Visualization of rendered RGB image, optical flow, and depth by SplatFlow on Waymo dataset.

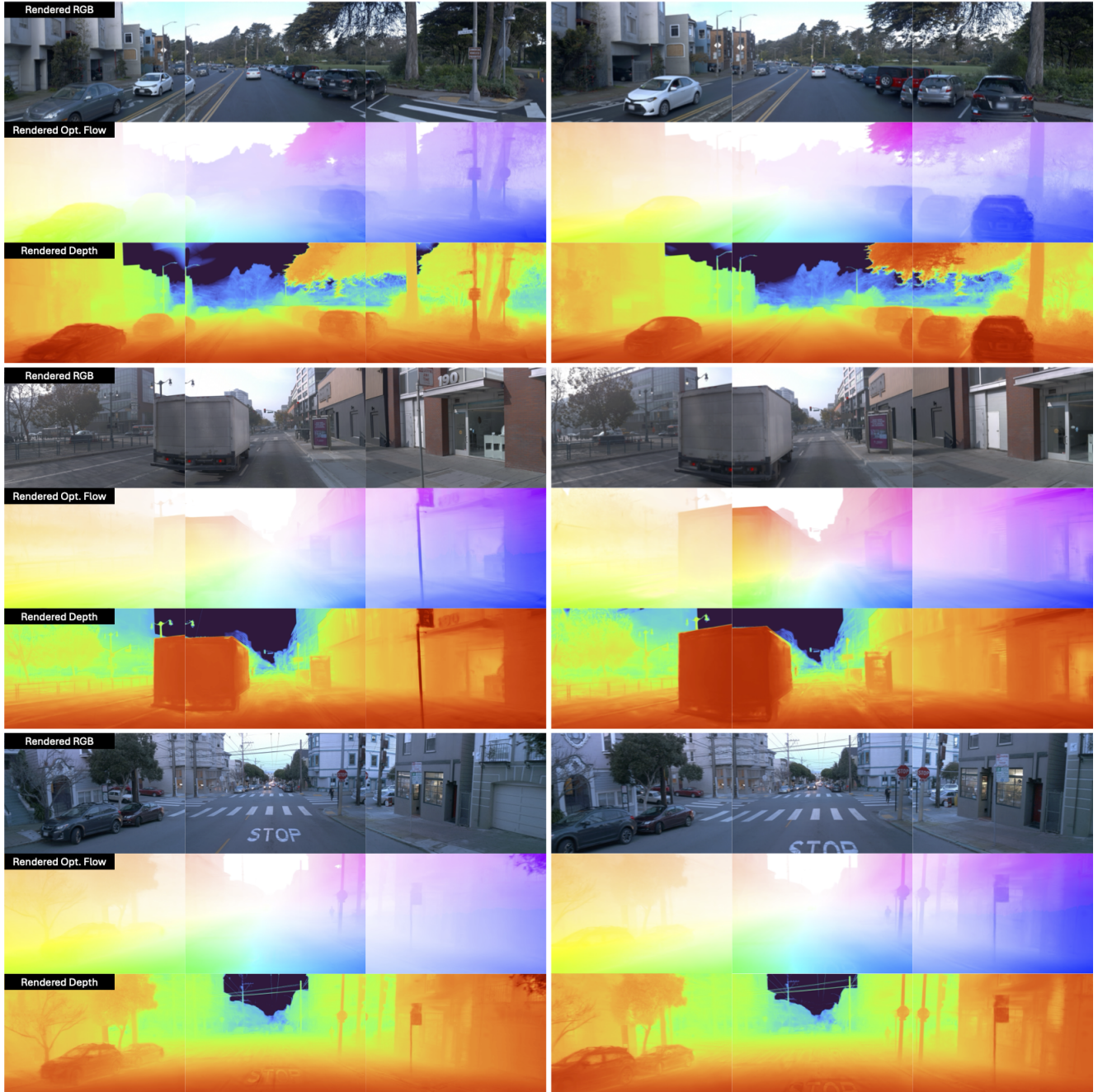


Figure 14. Visualization of rendered RGB image, optical flow, and depth by SplatFlow on Waymo dataset.



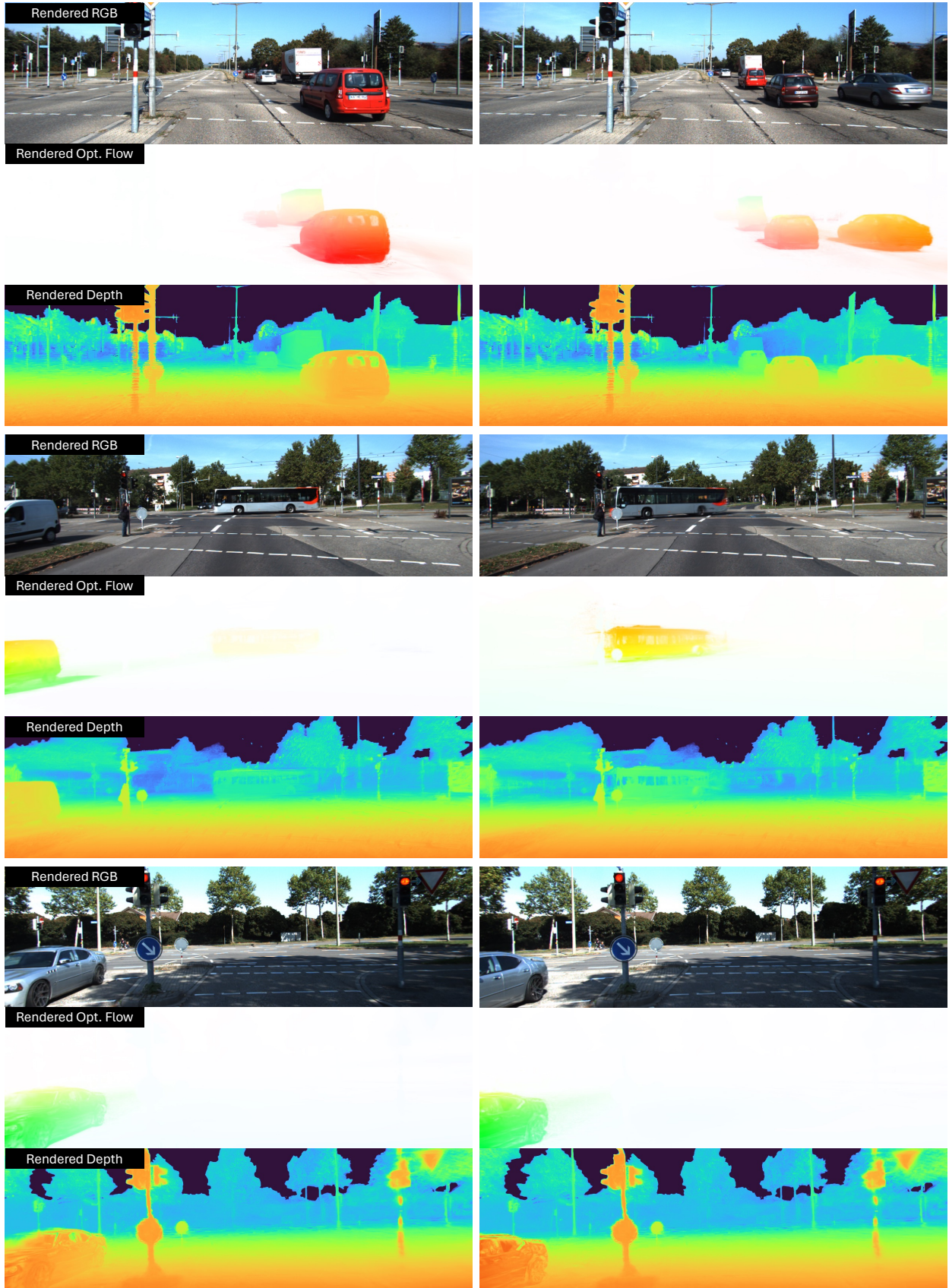


Figure 15. Visualization of rendered RGB image, optical flow, and depth by SplatFlow on KITTI dataset.



## References

- [1] Yurui Chen, Chun Gu, Junzhe Jiang, Xiatian Zhu, and Li Zhang. Periodic vibration gaussian: Dynamic urban scene reconstruction and real-time rendering. *arXiv preprint arXiv:2311.18561*, 20. 1, 2, 3, 5, 6, 7, 8
- [2] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. In *Conference on robot learning*, pages 1–16. PMLR, 2017. 2
- [3] Tobias Fischer, Lorenzo Porzi, Samuel Rota Bulo, Marc Pollefeys, and Peter Kotschieder. Multi-level neural scene graphs for dynamic urban environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21125–21135, 2024. 1
- [4] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012. 2, 5
- [5] Jianfei Guo, Nianchen Deng, Xinyang Li, Yeqi Bai, Botian Shi, Chiyu Wang, Chenjing Ding, Dongliang Wang, and Yikang Li. Streetsurf: Extending multi-view implicit surface reconstruction to street views. *arXiv preprint arXiv:2306.04988*, 2023. 2, 5, 6, 7, 1
- [6] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. 1, 2, 3, 5, 6, 7, 8
- [7] Abhijit Kundu, Kyle Genova, Xiaoqi Yin, Alireza Fathi, Caroline Pantofaru, Leonidas J Guibas, Andrea Tagliasacchi, Frank Dellaert, and Thomas Funkhouser. Panoptic neural fields: A semantic object-aware neural scene representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12871–12881, 2022. 1, 2
- [8] Fan Lu, Yan Xu, Guang Chen, Hongsheng Li, Kwan-Yee Lin, and Changjun Jiang. Urban radiance field representation with deformable neural mesh primitives. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 465–476, 2023. 2
- [9] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2, 8
- [10] 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. 1, 2, 3, 5, 6, 7, 8
- [11] Julian Ost, Issam Laradji, Alejandro Newell, Yuval Bahat, and Felix Heide. Neural point light fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18419–18429, 2022. 2
- [12] Jhony Kaesemodel Pontes, James Hays, and Simon Lucey. Scene flow from point clouds with or without learning. In *2020 international conference on 3D vision (3DV)*, pages 261–270. IEEE, 2020. 1
- [13] Konstantinos Rematas, Andrew Liu, Pratul P Srinivasan, Jonathan T Barron, Andrea Tagliasacchi, Thomas Funkhouser, and Vittorio Ferrari. Urban radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12932–12942, 2022. 2
- [14] 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. 2
- [15] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2446–2454, 2020. 2, 5
- [16] Adam Tonderski, Carl Lindström, Georg Hess, William Ljungbergh, Lennart Svensson, and Christoffer Petersson. Neurad: Neural rendering for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14895–14904, 2024. 1
- [17] 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. 1, 2, 3, 5, 6, 7, 8
- [18] Yihan Wang, Lahav Lipson, and Jia Deng. Sea-raft: Simple, efficient, accurate raft for optical flow. In *European Conference on Computer Vision*, pages 36–54. Springer, 2025. 5
- [19] 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. 1, 2, 3, 5, 6, 7, 8
- [20] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *Neural Information Processing Systems (NeurIPS)*, 2021. 5
- [21] Ziyang Xie, Junge Zhang, Wenye Li, Feihu Zhang, and Li Zhang. S-nerf: Neural radiance fields for street views. *arXiv preprint arXiv:2303.00749*, 2023. 5, 6, 7, 1, 2
- [22] Ziyang Xie, Junge Zhang, Wenye Li, Feihu Zhang, and Li Zhang. S-nerf: Neural radiance fields for street views. In *International Conference on Learning Representations (ICLR)*, 2023. 1
- [23] Yunzhi Yan, Haotong Lin, Chenxu Zhou, Weijie Wang, Haiyang Sun, Kun Zhan, Xianpeng Lang, Xiaowei Zhou, and Sida Peng. Street gaussians for modeling dynamic urban scenes. *arXiv preprint arXiv:2401.01339*, 2024. 1, 2, 3, 5, 8
- [24] Jiawei Yang, Boris Ivanovic, Or Litany, Xinshuo Weng, Seung Wook Kim, Boyi Li, Tong Che, Danfei Xu, Sanja Fidler, Marco Pavone, et al. Emernerf: Emergent spatial-temporal scene decomposition via self-supervision. *arXiv preprint arXiv:2311.02077*, 2023. 1, 2, 3, 5, 6, 7
- [25] Cheng Zhao, Su Sun, Ruoyu Wang, Yuliang Guo, Jun-Jun Wan, Zhou Huang, Xinyu Huang, Yingjie Victor Chen, and

- Liu Ren. Tlc-gs: Tightly coupled lidar-camera gaussian splatting for surrounding autonomous driving scenes. In *European Conference on Computer Vision*. Springer, 2024. 2
- [26] Hongyu Zhou, Jiahao Shao, Lu Xu, Dongfeng Bai, Weichao Qiu, Bingbing Liu, Yue Wang, Andreas Geiger, and Yiyi Liao. Hugs: Holistic urban 3d scene understanding via gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21336–21345, 2024. 3
- [27] Xiaoyu Zhou, Zhiwei Lin, Xiaojun Shan, Yongtao Wang, Deqing Sun, and Ming-Hsuan Yang. Drivinggaussian: Composite gaussian splatting for surrounding dynamic autonomous driving scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21634–21643, 2024. 2, 3, 5