

# EMD: Explicit Motion Modeling for High-Quality Street Gaussian Splatting

Xiaobao Wei<sup>1,2,\*†</sup>, Qingpo Wuwu<sup>1,2,\*†</sup>, Zhongyu Zhao<sup>1,2,†</sup>, Zhuangzhe Wu<sup>1</sup>,  
Nan Huang<sup>1</sup>, Ming Lu<sup>1</sup>, Ningning Ma<sup>2</sup>, Shanghang Zhang<sup>1,‡</sup>

<sup>1</sup>Peking University <sup>2</sup>Autonomous Driving Development, NIO  
weixiaobao0210@gmail.com

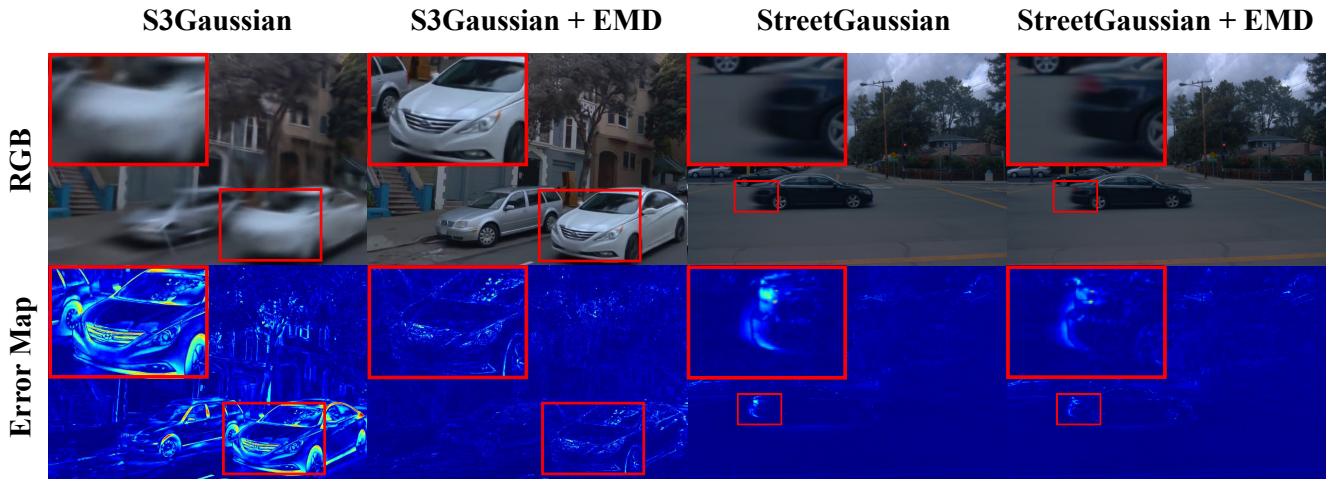


Figure 1: Previous street Gaussian splatting methods find it challenging to accurately model the motion of dynamic objects, which leads to blurry reconstructions. With the introduction of the proposed Explicit Motion Decomposition (EMD), which improves the decomposition of dynamic object motion, the current method achieves significantly better reconstruction quality.

## Abstract

Photorealistic reconstruction of street scenes is essential for developing real-world simulators in autonomous driving. While recent methods based on 3D/4D Gaussian Splatting (GS) have demonstrated promising results, they still encounter challenges in complex street scenes due to the unpredictable motion of dynamic objects. Current methods typically decompose street scenes into static and dynamic objects, learning the Gaussians in either a supervised manner (e.g., w/ 3D bounding-box) or a self-supervised manner (e.g., w/o 3D bounding-box). However, these approaches do not effectively model the motions of dynamic objects (e.g., the motion speed of pedestrians is clearly different from that of vehicles), resulting in sub-optimal scene decomposition. To address this, we propose Explicit Motion Decomposition (EMD), which models the motions of dynamic objects by introducing learnable motion embeddings to the Gaussians, enhancing the decomposition in street scenes. The proposed EMD is a plug-and-play approach

applicable to various baseline methods. We also propose tailored training strategies to apply EMD to both supervised and self-supervised baselines. Through comprehensive experimentation, we illustrate the effectiveness of our approach with various established baselines. The code will be released at: <https://qingpowuwu.github.io/emdgaussian.github.io/>.

## 1 Introduction

Novel view synthesis for dynamic scenes is essential in autonomous driving, enabling a variety of applications such as simulation, testing, and validation of perception systems. Traditional simulators like CARLA (Dosovitskiy et al. 2017) and AirSim (Shah et al. 2018), while providing controlled environments, suffer from limited realism and require substantial manual effort in creating virtual scenes. Recent advances in neural rendering, particularly Neural Radiance Fields (NeRF) (Mildenhall et al. 2020) and 3D Gaussian Splattting (3DGS) (Kerbl et al. 2023), have emerged as promising alternatives for photorealistic scene reconstruction. These pioneering approaches excel in capturing complex geometries and appearances through implicit or explicit neural representations, and

\*Equal contribution.

†Work done during internship at NIO.

‡Corresponding Author.

have been extended to dynamic scenes by incorporating an additional time dimension into their representations and learning deformation networks to model non-rigid motions (Park et al. 2021a; Pumarola et al. 2020; Li et al. 2021; Tretschk et al. 2020; Park et al. 2021b; Xian et al. 2020; Xu, Alldieck, and Sminchisescu 2021; Wang et al. 2022; Cleac'h et al. 2022; Cao and Johnson 2023; Rabich, Stotko, and Klein 2023). Despite advancements, reconstructing autonomous driving scenes remains difficult due to complex multi-object dynamics, expansive environments, and varied motion patterns.

To tackle this challenge, existing works usually build upon dynamic NeRF and 3DGS frameworks, further separating autonomous driving scenes into static and dynamic components through two primary paradigms: Supervised methods utilize pre-trained models to acquire auxiliary conditions, such as segmentation masks from SAM (Kirillov et al. 2023), depth maps from DepthAnything (Yang et al. 2024a), and both depth and optical flow from Dynamo (Sun and Hariharan 2023), or 3D bounding boxes from various datasets. Representative works like StreetGaussian (Yan et al. 2024) demonstrate the effectiveness of using supervision signals for trajectory optimization and appearance modeling. In contrast, self-supervised methods achieve static-dynamic separation without explicit supervision, as shown by S3Gaussian (Huang et al. 2024), which leverages inherent motion cues. While both paradigms have shown promising results, their binary classification of scene elements as either static or dynamic overlooks the continuous spectrum of motion inherent in real-world street scenes.

To better address the different motion patterns in street scenes, we propose an Explicit Motion Decomposition (EMD) module that can be easily integrated into existing supervised and self-supervised frameworks (Fig. 2). EMD improves scene decomposition by incorporating motion-aware feature encoding and dual-scale deformation modeling. Specifically, we enhance each Gaussian primitive with learnable motion embeddings to capture its motion characteristics and design a hierarchical deformation framework that separately manages fast, global motions and slow, local deformations. This design allows for more efficient analysis of complex street scenes with varying motion speeds.

To showcase the versatility of our approach, we conduct extensive experiments on the Waymo-NOTR dataset by integrating EMD with representative methods from both paradigms: StreetGaussian for supervised and S3Gaussian for self-supervised settings. Our main contributions include:

- We propose EMD, a new plug-and-play module that effectively addresses varying motion speeds in street scenes through explicit motion modeling.
- We validate our method in both supervised and self-supervised settings, demonstrating consistent performance improvements across various evaluation protocols.
- We carried out a thorough evaluation of the Waymo-NOTR dataset, showcasing substantial improvements across various settings, achieving better reconstruction quality in both full scenes (+1.81 PSNR) and vehicle-specific regions (+2.81 PSNR) compared to respective

baselines.

## 2 Related Work

**Dynamic Scene Representation.** Neural Radiance Fields (NeRF) (Mildenhall et al. 2020) revolutionized novel view synthesis by introducing volumetric scene representation using MLPs (Taud and Mas 2018), and has been significantly improved through various extensions (Barron et al. 2021, 2022; Müller et al. 2022; Liu et al. 2020). However, these methods face challenges in dynamic scene reconstruction due to the lack of temporal modeling capabilities. A common solution involves introducing time-dependent modeling through additional time conditions and deformation networks, as adopted by (Park et al. 2021a; Pumarola et al. 2020; Park et al. 2021b). Additionally, several works also improve NeRF’s reconstruction quality in both static and dynamic scenes by introducing additional supervision signals, such as semantic segmentation (Kundu et al. 2022; Zhi et al. 2021; Jain, Tancik, and Abbeel 2021), depth maps (Xu et al. 2022; Guo et al. 2022; Deng et al. 2022a; Niemeyer et al. 2022; Deng et al. 2022b), and optical flow (Chen and Tsukada 2022; Li et al. 2023). Despite these advances, NeRF-based methods struggle with long training times and limited ability to handle complex motions. Recently, 3D Gaussian Splatting (3DGS) (Kerbl et al. 2023) has demonstrated remarkable potential by representing scenes with explicit 3D Gaussian primitives, achieving both high training efficiency and superior rendering quality. Inspired by NeRF-based approaches in handling dynamic scenes, several works have extended 3DGS to model temporal changes. Specifically, (Luiten et al. 2024) proposes dynamic modeling by tracking Gaussian properties (position, orientation) while maintaining persistent appearance, while (Yang et al. 2024b) extends the 3D representation to 4D by introducing temporal Gaussian distributions. Building upon these ideas, (Yang et al. 2023b) designs dedicated deformation networks for motion modeling, (Wu et al. 2024) advances the efficiency by maintaining only canonical 3D Gaussians with HexPlane and (Bae et al. 2024) introduces embeddings to model fast and slow motions. These methods lay the groundwork for dynamic scene reconstruction in autonomous driving scenarios.

**Autonomous Driving Simulation** Traditional autonomous driving simulators like AirSim (Shah et al. 2018) and CARLA (Dosovitskiy et al. 2017) require extensive manual effort for environment creation while struggling to achieve photorealistic rendering. To address this, Neural based approaches have emerged as promising solutions for simulating street scenes. Early NeRF-based methods (Ost et al. 2021; Turki, Ramanan, and Satyanarayanan 2021; Rematas et al. 2022; Huang et al. 2023) introduced neural scene representation for large-scale urban environments, followed by improvements in efficiency and scalability (Tancik et al. 2022; Turki et al. 2023; Lu et al. 2023; Liu et al. 2023). For dynamic urban scene reconstruction, NSG (Ost et al. 2021) introduced supervised scene decomposition using learned scene graphs and latent object representations. This supervised approach has since become a new paradigm, adopted by MARS (Wu et al. 2023), NeuRAD (Tonderski et al. 2023) and Multi-Level Neural

Scene Graphs (Fischer et al. 2024). With the advent of 3D Gaussian Splatting, DrivingGaussian (Zhou et al. 2024b) also adopted this supervised paradigm and developed hierarchical scene representations combining dynamic object graphs with incrementally updated static elements, while StreetGaussian (Yan et al. 2024) enhanced this framework with trajectory optimization and appearance modeling using 4D spherical harmonics. Recent work HUGS (Zhou et al. 2024a) further advances this line of research by incorporating physical constraints into the joint optimization of geometry, appearance, and semantics. To eliminate the need for expensive supervision, VDG (Li et al. 2024) introduced pose-free reconstruction by integrating self-supervised visual odometry for pose estimation and depth initialization, while employing motion mask supervision for static-dynamic scene decomposition. Inspired by D2NeRF’s (Wu et al. 2022) self-supervised decomposition in general dynamic scenes, SUDS (Turki et al. 2023) introduced this paradigm to autonomous driving by leveraging optical flow guidance, followed by EmerNeRF (Yang et al. 2023a) proposing a flow-free approach with hash-grid based scene organization. Despite the dominance of NeRF-based methods, this self-supervised paradigm was later extended to 3D Gaussian Splatting by several works. PVG (Chen et al. 2023) pioneered this extension by introducing periodic vibration-based temporal dynamics for unified representation of both static and dynamic elements, followed by S3Gaussian (Huang et al. 2024).

However, current methods typically employ binary static-dynamic classification, which oversimplifies the complex motion patterns in street scenes. This limitation becomes particularly evident when reconstructing objects with varying velocities, as shown in our evaluations on the Waymo-NOTR dataset. To address this challenge, we propose Explicit Motion Decomposition (EMD), a plug-and-play approach that better captures the continuous spectrum of motion in street scenes. Before presenting our method, we first review the fundamentals of 3D and 4D Gaussian Splatting.

### 3 Preliminaries

#### 3.1 3D Gaussian Splatting

3D Gaussian Splatting (3D-GS) (Kerbl et al. 2023) proposes an explicit rendering approach to represent 3D scenes through a collection of 3D Gaussian primitives  $\mathbb{G} = \{(\mu_k, \Sigma_k, \alpha_k, \mathbf{c}_k)\}_{k=1}^K$ , where  $K$  is the total number of Gaussians. Each Gaussian primitive represents a probabilistic distribution of density in 3D space, defined by its probability density function:

$$G_k(x) = e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1}(x-\mu_k)}, \quad (1)$$

where  $x \in \mathbb{R}^3$  represents any point in the 3D world space,  $\mu_k \in \mathbb{R}^3$  and  $\Sigma_k \in \mathbb{R}^{3 \times 3}$  denote the mean position and covariance matrix in world space, respectively, where the covariance matrix determines the shape and orientation of the Gaussian,  $\alpha_k \in [0, 1]$  is the opacity, and  $\mathbf{c}_k$  encodes the view-dependent color information.

**Rendering Process** For rendering, each 3D Gaussian is projected onto the image plane, where the 3D mean  $\mu_k$  is

transformed to 2D mean  $\mu_k^{2D}$ , and the world space covariance matrix  $\Sigma_k$  is transformed to screen space as  $\Sigma'_k = JW\Sigma_k W^T J^T$ , with  $W$  and  $J$  being the viewing transformation and projective transformation Jacobian matrices, respectively. The final pixel color at screen space position  $x$  is computed through front-to-back alpha compositing:

$$C(x) = \sum_{k \in \mathcal{N}(\mathbf{x})} \mathbf{c}_k \alpha_k(x) \prod_{j=1}^{k-1} (1 - \alpha_j(x)), \quad (2)$$

where  $\alpha_k(x) = \alpha_k \exp\left(-\frac{1}{2}(x - \mu_k^{2D})^T \Sigma'_k^{-1}(x - \mu_k^{2D})\right)$  represents the opacity contribution of the  $k$ -th Gaussian at pixel  $x$ , computed as the product of the Gaussian’s base opacity  $\alpha_k$  and its projected 2D Gaussian evaluation.  $\mathcal{N}(\mathbf{x})$  represents the set of indices of Gaussians intersecting pixel  $x \in \mathbb{R}^2$ .

In practice, for each Gaussian  $k$ , we parameterize its covariance matrix  $\Sigma_k$  using rotation quaternion  $\mathbf{q}_k$  and scaling vector  $\mathbf{s}_k$  as:

$$\Sigma_k = R(\mathbf{q}_k) S(\mathbf{s}_k) S(\mathbf{s}_k)^T R(\mathbf{q}_k)^T, \quad (3)$$

where  $R(\mathbf{q}_k)$  is the rotation matrix defined by quaternion  $\mathbf{q}_k$ , and  $S(\mathbf{s}_k)$  is the diagonal scaling matrix defined by scaling vector  $\mathbf{s}_k$ . The view-dependent color  $\mathbf{c}_k$  is encoded using spherical harmonics (SH) coefficients:

$$\mathbf{c}_k(\mathbf{d}) = \sum_{l=0}^L \sum_{m=-l}^l k_{l,m}^{(k)} Y_{l,m}(\mathbf{d}), \quad (4)$$

where  $\mathbf{d}$  is the viewing direction,  $Y_{l,m}$  are the spherical harmonic basis functions, and  $k_{l,m}^{(k)}$  are the corresponding coefficients for the  $k$ -th Gaussian.

#### 3.2 4D Gaussian Splatting

4D Gaussian Splatting extends the static 3D-GS framework to handle dynamic scenes by incorporating temporal information. For a dynamic scene captured at different timestamps  $t \in [0, T]$ , each Gaussian primitive is now characterized by time-varying parameters:  $\mathbb{G}(t) = \{(\mu_k(t), \Sigma_k(t), \alpha_k(t), \mathbf{c}_k(t))\}_{k=1}^K$ . The probability density function of each Gaussian at time  $t$  becomes:

$$G_k(x, t) = e^{-\frac{1}{2}(x - \mu_k(t))^T \Sigma_k(t)^{-1}(x - \mu_k(t))}, \quad (5)$$

The temporal evolution of these parameters is typically modeled through one of two main approaches:

**Deformation-based Approach** This approach models temporal changes by applying a deformation field to the base Gaussians. For each timestamp  $t$ , the position of each Gaussian is updated as:

$$\mu_k(t) = \mu_k(0) + \Delta\mu_k(t), \quad (6)$$

where  $\mu_k(0)$  is the initial position and  $\Delta\mu_k(t)$  is the displacement predicted by a deformation network. Similarly, other parameters such as rotation, scaling, and color can be modeled as temporal offsets from their initial states.

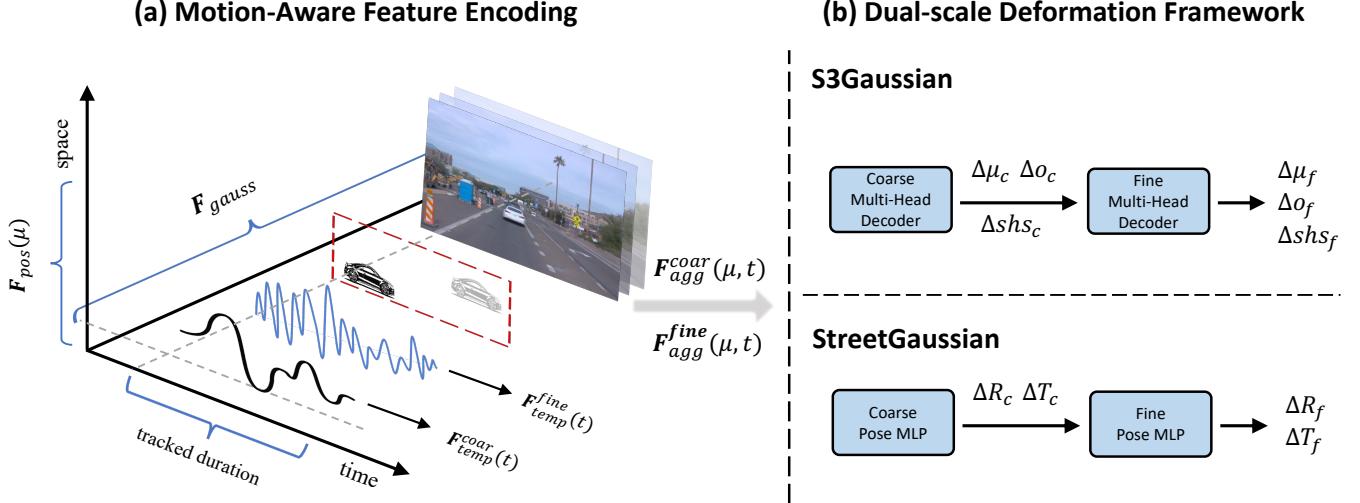


Figure 2: Overview of our Explicit Motion Decomposition (EMD) framework. Given input Gaussian primitives, our method processes them through two main components: (a) Motion-aware Feature Encoding, which combines spatial, temporal, and Gaussian-specific information to capture motion characteristics; and (b) Dual-scale Deformation Framework, which hierarchically models fast global motions and slow local deformations. The framework can be seamlessly integrated into both supervised (StreetGaussian) and self-supervised (S3Gaussian) approaches through our proposed integration strategies.

**Trajectory-based Approach** Alternatively, the temporal evolution can be represented by explicitly modeling the continuous trajectory of each parameter. For instance, the position trajectory can be parameterized using a set of  $N$  control points  $\{\mathbf{p}_i\}_{i=1}^N$  and basis functions  $\{\phi_i(t)\}_{i=1}^N$ :

$$\mu_k(t) = \sum_{i=1}^N \mathbf{p}_i \phi_i(t), \quad (7)$$

where the basis functions  $\phi_i(t)$  can be B-splines or other temporal interpolation functions. The rendering process remains similar to static 3D-GS, but now Gaussian parameters are evaluated at the specific timestamp  $t$  before projection and compositing.

## 4 Methodology

4D Gaussian Splatting has shown promising results in dynamic scene reconstruction. However, modeling dynamic street scenes remains challenging due to diverse motion patterns. To better handle the varying motion patterns in street scenes, we propose Explicit Motion Decomposition (EMD), a plug-and-play module that can be seamlessly integrated into existing 4D Gaussian-based frameworks to enhance their capability in handling dynamic scenarios, as illustrated in Fig. 2.

### 4.1 Problem Formulation

Given a set of static 3D Gaussian primitives  $\mathbb{G} = \{(\mu_k, \mathbf{s}_k, \mathbf{q}_k, \alpha_k, \mathbf{c}_k)\}_{k=1}^K$  and a timestamp  $t$ , our goal is to learn a deformation field  $\mathcal{D}$  that maps each Gaussian's parameters from their canonical states to their corresponding deformed states at time  $t$ . For notational simplicity, we omit the Gaussian index  $k$  in the following formulation:

$$\{\mu_t, \mathbf{s}_t, \mathbf{q}_t, \alpha_t, \mathbf{c}_t\} = \mathcal{D}(\{\mu, \mathbf{s}, \mathbf{q}, \alpha, \mathbf{c}\}, t), \quad (8)$$

To effectively handle the diverse motion patterns in street scenes, especially the distinct movements between vehicles and pedestrians, we propose a motion-aware deformation module that processes input Gaussian parameters through two key components: motion-aware feature encoding and dual-scale deformation prediction.

**Motion-aware Feature Encoding** Different types of objects in street scenes exhibit distinct motion characteristics. To capture these varied patterns, we first encode the input Gaussian parameters into a comprehensive feature space that combines spatial, temporal, and Gaussian-specific information:

$$\mathbf{F}_{agg}(μ, t) = [\mathbf{F}_{pos}(μ), \mathbf{F}_{temp}(t), \mathbf{F}_{gauss}], \quad (9)$$

The spatial component employs multi-frequency positional encoding:

$$\mathbf{F}_{pos}(μ) = [μ, \{\sin(2^i π μ), \cos(2^i π μ)\}_{i=0}^{P-1}], \quad (10)$$

where  $P$  is the number of frequency bands, enabling the network to capture both fine geometric details and global structures. For temporal information, we design an adaptive temporal embedding function:

$$\mathbf{F}_{temp}(t) = \mathcal{T}(t, N(i)) = \text{Interp}(\mathbf{W}, t, N(i)), \quad (11)$$

where  $\mathbf{W} \in \mathbb{R}^{N_{max} \times D}$  is a learnable embedding matrix that captures motion patterns, and  $N(i)$  progressively increases from  $N_{min}$  to  $N_{max}$  temporal samples during training iteration  $i$ , allowing the model to gradually capture finer temporal dynamics. For Gaussian-specific features, we assign a learnable latent embedding  $\mathbf{z}_k \in \mathbb{R}^D$  to each Gaussian  $k$ , where  $\mathbf{F}_{gauss} = \mathbf{z}_k$ , enabling the model to learn and represent individual motion characteristics.

Table 1: Comparative performance of our framework and baseline approaches on the Waymo-NOTR dataset. The best performances are highlighted in **bold**, and the second-best are indicated with underlining.  $\uparrow$  indicates higher is better, while  $\downarrow$  indicates lower is better.

Dataset	Methods	Scene Reconstruction						Novel View Synthesis					
		Full Image			Vehicle			Full Image			Vehicle		
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$
D32	EmerNeRF	28.16	0.806	0.228	24.32	0.682	25.14	0.747	0.313	<b>23.49</b>	0.660		
	3DGS	28.47	0.876	0.136	23.26	0.716	25.14	0.813	0.165	20.48	<b>0.753</b>		
	MARS	28.24	0.866	0.252	23.37	0.701	26.61	0.796	0.305	22.21	0.697		
	S3Gaussian	30.69	0.900	0.121	26.23	0.804	<b>26.62</b>	0.824	0.159	22.61	0.681		
S3Gaussian+Ours		<b>32.50</b>	<u>0.933</u>	<b>0.082</b>	<b>29.04</b>	<b>0.879</b>	<u>26.55</u>	<b>0.833</b>	<b>0.126</b>	<u>23.39</u>	<u>0.703</u>		

**Dual-scale Deformation Framework** Given the diverse motion patterns in street scenes, ranging from large vehicular movements to subtle pedestrian motions, we design a hierarchical deformation framework that can effectively handle both scales of motion:

$$\begin{aligned} \mathcal{D}(\mu, t) = & \mathcal{D}_{coarse}(\mathbf{F}_{agg}(\mu, t)) \\ & + \mathcal{D}_{fine}(\mathbf{F}_{agg}(\mu + \Delta\mu_{coarse}, t)), \end{aligned} \quad (12)$$

The final deformed parameters combine both coarse and fine scale predictions:

$$\begin{aligned} \mu_t &= \mu + \Delta\mu_{coarse} + \Delta\mu_{fine} \\ \mathbf{s}_t &= \mathbf{s} + \Delta\mathbf{s}_{coarse} + \Delta\mathbf{s}_{fine} \\ \mathbf{q}_t &= \mathbf{q} \otimes \Delta\mathbf{q}_{coarse} \otimes \Delta\mathbf{q}_{fine}, \end{aligned} \quad (13)$$

where  $\mathcal{D}_{coarse}$  focuses on modeling large-scale motions such as vehicle translations, while  $\mathcal{D}_{fine}$  captures local deformations like articulated movements. Similar to position updates, other Gaussian parameters including opacity  $\alpha_t$  and spherical harmonics coefficients  $\mathbf{c}_t$  are also updated through this dual-scale framework.

## 4.2 Integration with Existing Frameworks

As mentioned in the introduction, current approaches for street scene reconstruction generally fall into supervised and self-supervised paradigms. Having formalized our EMD framework, we now demonstrate its integration into representative methods from both paradigms: StreetGaussian for supervised learning and S3Gaussian for self-supervised learning.

**Self-supervised Integration: S3Gaussian** For self-supervised scenarios, we enhance S3Gaussian’s architecture with our motion-aware features. While S3Gaussian originally employs a Multi-head Gaussian Decoder for deformation prediction, we augment each Gaussian with our learnable embedding  $\mathbf{z}_k$  and restructure its decoder into our dual-scale framework. Specifically, both coarse and fine stages predict deformations in position ( $\Delta\mu$ ), opacity ( $\Delta\alpha$ ), and spherical harmonics coefficients ( $\Delta\mathbf{c}$ ), enabling more precise motion modeling through hierarchical refinement.

**Supervised Integration: StreetGaussian** For supervised settings, StreetGaussian provides a framework that represents dynamic objects through tracked vehicle poses and object-specific Gaussians. Each object is characterized by tracked

poses  $\{R_t, T_t\}_{t=1}^{N_t}$  that transform object Gaussians from local coordinates  $(\mu_o, R_o)$  to world coordinates  $(\mu_w, R_w)$ . To enhance its motion modeling capability, we also augment each Gaussian with our learnable embedding  $\mathbf{z}_k$  and apply temporal embedding only within each object’s tracked duration, effectively capturing object-specific temporal dynamics.

To address the challenge of noisy tracked poses while maintaining their geometric meaning, we incorporate our dual-scale framework into their pose optimization:

$$\begin{aligned} R'_t &= R_t(\Delta R_t^c + \Delta R_t^f) \\ T'_t &= T_t + (\Delta T_t^c + \Delta T_t^f), \end{aligned} \quad (14)$$

where  $\Delta R_t^c, \Delta T_t^c$  handle large pose corrections and  $\Delta R_t^f, \Delta T_t^f$  capture subtle adjustments. This decomposition allows the model to effectively correct tracking errors while preserving the physical meaning of the tracked poses. Similarly, we apply the dual-scale framework to appearance modeling, where spherical harmonics coefficients are refined through both coarse and fine stages, enabling more accurate dynamic appearance representation.

We adopt the same loss function as S3Gaussian and StreetGaussian to train the entire pipeline. For further training details, please refer to the supplementary materials.

## 5 Experiments

In this section, we first describe the datasets and evaluation metrics used in our experiments (Sec. 5.1). Then, we demonstrate the experimental results of our approach under both self-supervised and supervised settings (Sec. 5.2). Finally, we present comprehensive ablation studies to validate the effectiveness of individual components in our framework (Sec. 5.3).

### 5.1 Datasets and Metrics

**Dataset.** We conduct extensive experiments on the Waymo Open Dataset (Sun et al. 2020) to comprehensively evaluate our method. To validate the effectiveness of our plug-and-play module, we select representative state-of-the-art methods from both supervised and self-supervised paradigms: StreetGaussian (Yan et al. 2024) and S3Gaussian (Huang et al. 2024). For self-supervised evaluation, we use the dynamic32 (D32) split introduced by EmerNeRF (Yang et al. 2023a), containing 32 sequences with vehicle motion, where each

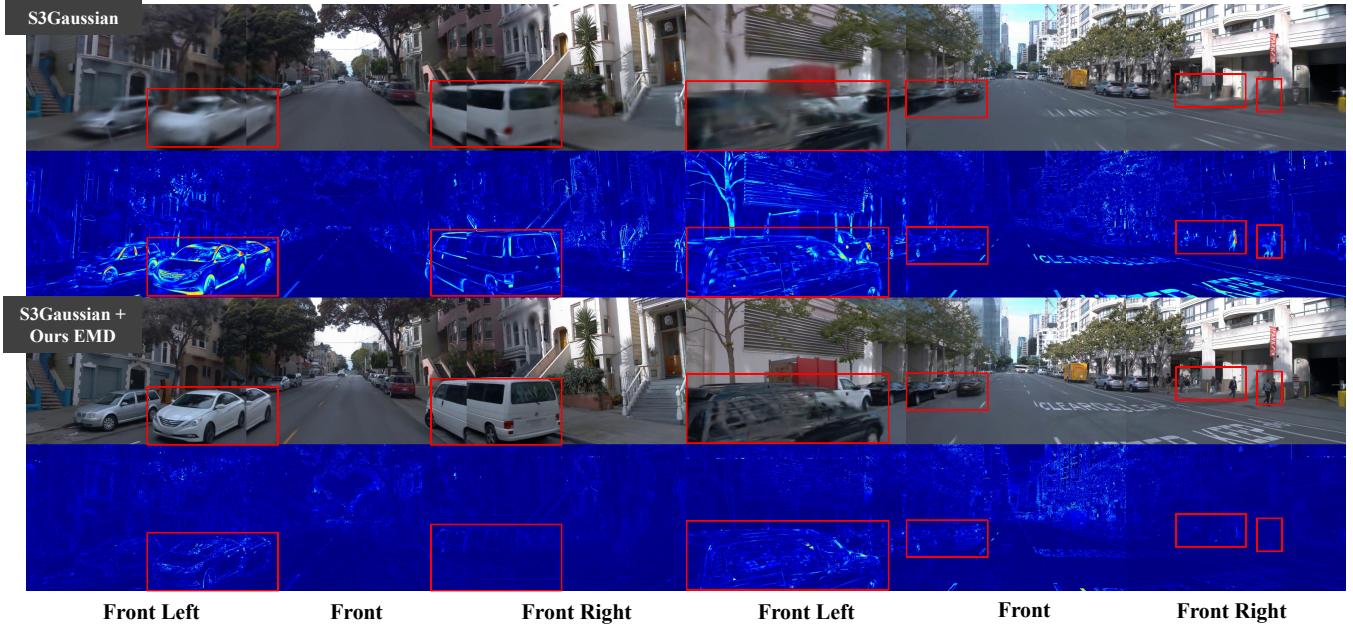


Figure 3: Qualitative comparison on the self-supervised setting between S3Gaussian and S3Gaussian+ours EMD. We also visualize the error maps between the rendered images and ground truth to provide further insights.



Figure 4: Motion deformation comparison between S3Gaussian + Ours and S3Gaussian. Please zoom in for more details.

sequence consists of approximately 50-100 frames captured by three cameras under various conditions. For supervised evaluation, we follow the data split protocol established in StreetGaussian (Yan et al. 2024) to enable direct comparison.

**Evaluation Metrics.** We employ comprehensive metrics to evaluate both reconstruction quality and novel view synthesis capability. For scene reconstruction, we use PSNR and SSIM to evaluate both full-scene and vehicle-specific reconstruction quality. Additionally, we compute LPIPS for perceptual quality assessment. Following previous protocols, we evaluate novel view synthesis on every 10th frame for self-supervised settings and every 4th frame for supervised settings.

## 5.2 Main Results

**Self-supervised Performance** Tab. 1 presents the comparative results on the D32 split, where no 3D bounding box annotations were used. Our method significantly outperforms previous self-supervised approaches, achieving notable improvements in both full-scene and object-specific metrics. Specifically, we observe a substantial increase in PSNR for the full scene (32.50 vs. 30.69) and for vehicle-specific metrics (PSNR: 29.04 vs. 26.23) when compared to S3Gaussian. These results demonstrate the enhanced ability of our method

to accurately model complex street scenes without explicit 3D box annotations. For novel view synthesis, our method continues to perform competitively, achieving a PSNR of 26.55, which reflects its robust generalization to previously unseen viewpoints.

In addition, we present a side-by-side visualization comparison between S3Gaussian and S3Gaussian+EMD in Fig. 3. The error maps, which compare the ground truth to the rendered results, clearly demonstrate that S3Gaussian+EMD outperforms S3Gaussian in modeling dynamic objects with varying motion speeds. S3Gaussian implicitly models dynamic vehicles, but it fails to capture changes in speed during motion and the differences in movement between vehicles, leading to blurred reconstructions. On the contrary, our method captures the distinct motion characteristics of different dynamic objects, leading to more accurate and consistent scene reconstructions. This emphasizes the effectiveness of Explicit Motion Decomposition (EMD) in modeling the motion of dynamic objects, improving the overall decomposition and photorealistic rendering of street scenes.

We also present a visual comparison of motion between S3Gaussian and S3Gaussian+EMD. As shown in Fig. 4, the proposed dual-scale deformation network generates a coarse deformation to model slower motion and larger-scale geome-



Figure 5: Qualitative ablation study results across three camera views from the Waymo dataset. (a) Our complete model achieves sharp and consistent reconstruction. (b) Removing Gaussian embedding  $\mathcal{F}_{gauss}$  leads to blurry object boundaries. (c) Without temporal embedding  $\mathcal{F}_{temp}$ , results show motion artifacts. (d) Without coarse deformation  $\mathcal{D}_{coarse}$ , geometric consistency is lost. (e) Absence of fine deformation  $\mathcal{D}_{fine}$  causes detail degradation.

try and a fine deformation to capture faster motion and finer geometric details in the scene. With the incorporation of EMD, S3Gaussian can capture detailed features of dynamic vehicles, including the car brand logo, as demonstrated in the figure. In contrast, S3Gaussian treats the entire moving car as a single dynamic object, failing to produce clear synthesis results for the dynamic vehicles.

**Supervised Performance** To demonstrate the versatility of our framework, we also evaluate it in the supervised setting using the same scenes as StreetGaussian (Yan et al. 2024) (Tab. 2). When incorporating 3D bounding box supervision through our adaptive training scheme, our method achieves superior performance in full scene reconstruction, improving PSNR by 1.42dB (36.03 vs. 34.61), SSIM by 1.1% (0.949 vs. 0.938), and LPIPS by 15.2% (0.067 vs. 0.079). These results demonstrate that our approach effectively enhances the overall scene reconstruction quality while maintaining competitive performance in dynamic object modeling. Please refer to the supplementary for visualization on the supervised setting.

Table 2: Comparative performance of our framework and baseline approaches on the StreetGaussian dataset. The best performances are highlighted in **bold**.  $\uparrow$  indicates higher is better, while  $\downarrow$  indicates lower is better.

Methods	Full Image			Vehicle
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$
3DGS	29.64	0.918	0.117	21.25
NSG	28.31	0.862	0.346	24.32
MARS	29.75	0.886	0.264	26.54
EmerNeRF	30.87	0.905	0.133	21.67
StreetGaussian	34.61	0.938	0.079	<b>30.23</b>
StreetGaussian + Ours	<b>36.03</b>	<b>0.949</b>	<b>0.067</b>	29.81

### 5.3 Ablation Studies

To assess the contribution of each component in our framework, we perform comprehensive ablation studies, as shown in Tab. 3 and Fig. 5. The results reveal the critical role of the Gaussian embedding, as its removal leads to the performance drop, with a reduction of 0.29 PSNR (32.21 vs. 32.50).

Table 3: Ablation study on the D32 dataset showing the impact of different components in our framework. All variants are evaluated using the self-supervised setting. The best performances are highlighted in **bold**. ↑ indicates higher is better, while ↓ indicates lower is better.

Variant	Full Image			Vehicle
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑
Full Model	<b>32.50</b>	<b>0.933</b>	<b>0.082</b>	<b>29.04</b>
w/o Gaussian Embedding	32.21	0.928	0.089	28.80
w/o Temporal Embedding	32.23	0.922	0.091	28.08
w/o Coarse Deformation	29.40	0.890	0.146	24.54
w/o Fine Deformation	32.45	0.931	0.118	28.80

This indicates that the Gaussian embedding is essential for effectively capturing the motion characteristics for each dynamic gaussian. The temporal embedding also plays a crucial role, with its absence leading to a 0.27 PSNR drop (32.23 vs. 32.50), underscoring its importance in modeling the temporal variation of object motion over time.

Both the coarse and fine deformation components are integral to the final performance, with their removal leading to considerable performance degradation. Specifically, excluding the coarse deformation component causes a significant 3.10 PSNR drop (29.40 vs. 32.50), suggesting that the coarse adjustments are vital for maintaining the overall scene structure. On the other hand, removing the fine deformation component results in a 0.036 LPIPS increase (0.118 vs. 0.082), implying that fine deformation refines the details and its absence worsens the perceptual quality. The ablation study demonstrates that our proposed motion-aware feature encoding and dual-scale deformation effectively model dynamic objects with varying motion speeds, which is crucial for enhancing the reconstruction quality of existing street Gaussians.

## 6 Limitation

Although EMD effectively addresses the challenge of modeling dynamic objects with varying speeds by incorporating learnable embeddings, some limitations remain. Existing street Gaussian methods do not account for environmental lighting, yet lighting effects play a crucial role in the quality of reconstructions under different lighting conditions. In future work, we plan to explore the possibility of developing a plug-and-play technique to enhance lighting effects in existing methods.

## 7 Conclusion

In this paper, we present EMD, a plug-and-play module that effectively handles varying motion speeds in street scene reconstruction. By introducing motion-aware feature encoding and dual-scale deformation modeling, our approach successfully captures the continuous spectrum of motion patterns inherent in real-world scenarios. Comprehensive experiments on the Waymo Open Dataset demonstrate that EMD significantly improves reconstruction quality when integrated with both supervised and self-supervised frameworks. Our method substantially improves full scene reconstruction (+1.81 PSNR) and vehicle-specific regions (+2.81 PSNR),

validating the effectiveness of explicit motion modeling. We believe our work opens up new possibilities for high-quality dynamic scene reconstruction in autonomous driving applications, and the plug-and-play nature of our approach makes it readily applicable to future developments in neural scene representations.

## References

- Bae, J.; Kim, S.; Yun, Y.; Lee, H.; Bang, G.; and Uh, Y. 2024. Per-Gaussian Embedding-Based Deformation for Deformable 3D Gaussian Splatting. *arXiv preprint arXiv:2404.03613*.
- Barron, J. T.; Mildenhall, B.; Tancik, M.; Hedman, P.; Martin-Brualla, R.; and Srinivasan, P. P. 2021. Mip-NeRF: A Multiscale Representation for Anti-Aliasing Neural Radiance Fields. *arXiv:2103.13415*.
- Barron, J. T.; Mildenhall, B.; Verbin, D.; Srinivasan, P. P.; and Hedman, P. 2022. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5470–5479.
- Cao, A.; and Johnson, J. 2023. HexPlane: A Fast Representation for Dynamic Scenes. *CVPR*.
- Chen, Q.-A.; and Tsukada, A. 2022. Flow Supervised Neural Radiance Fields for Static-Dynamic Decomposition. In *2022 International Conference on Robotics and Automation (ICRA)*, 10641–10647.
- Chen, Y.; Gu, C.; Jiang, J.; Zhu, X.; and Zhang, L. 2023. Periodic Vibration Gaussian: Dynamic Urban Scene Reconstruction and Real-time Rendering. *arXiv:2311.18561*.
- Cleac'h, S. L.; Yu, H.; Guo, M.; Howell, T. A.; Gao, R.; Wu, J.; Manchester, Z.; and Schwager, M. 2022. Differentiable Physics Simulation of Dynamics-Augmented Neural Objects. *IEEE Robotics and Automation Letters*, 8: 2780–2787.
- Deng, K.; Liu, A.; Zhu, J.-Y.; and Ramanan, D. 2022a. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12882–12891.
- Deng, K.; Liu, A.; Zhu, J.-Y.; and Ramanan, D. 2022b. Depth-supervised NeRF: Fewer Views and Faster Training for Free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Dosovitskiy, A.; Ros, G.; Codella, F.; Lopez, A.; and Koltun, V. 2017. CARLA: An open urban driving simulator. In *Conference on robot learning*, 1–16. PMLR.
- Fischer, T.; Porzi, L.; Bulo, S. R.; Pollefeyns, M.; and Kortschieder, P. 2024. Multi-level neural scene graphs for dynamic urban environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 21125–21135.
- Guo, Y.-C.; Kang, D.; Bao, L.; He, Y.; and Zhang, S.-H. 2022. Nerfren: Neural radiance fields with reflections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 18409–18418.
- Huang, N.; Wei, X.; Zheng, W.; An, P.; Lu, M.; Zhan, W.; Tomizuka, M.; Keutzer, K.; and Zhang, S. 2024. S3Gaussian:

- Self-Supervised Street Gaussians for Autonomous Driving. *arXiv preprint arXiv:2405.20323*.
- Huang, S. Y.; Gojcic, Z.; Wang, Z.; Williams, F.; Kasten, Y.; Fidler, S.; Schindler, K.; and Litany, O. 2023. Neural LiDAR Fields for Novel View Synthesis. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, 18190–18200.
- Jain, A.; Tancik, M.; and Abbeel, P. 2021. Putting NeRF on a Diet: Semantically Consistent Few-Shot View Synthesis. *arXiv:2104.00677*.
- Kerbl, B.; Kopanas, G.; Leimkühler, T.; and Drettakis, G. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *arXiv:2308.04079*.
- Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A. C.; Lo, W.-Y.; Dollár, P.; and Girshick, R. 2023. Segment Anything. *arXiv:2304.02643*.
- Kundu, A.; Genova, K.; Yin, X.; Fathi, A.; Pantofaru, C.; Guibas, L.; Tagliasacchi, A.; Dellaert, F.; and Funkhouser, T. 2022. Panoptic Neural Fields: A Semantic Object-Aware Neural Scene Representation. In *CVPR*.
- Li, H.; Li, J.; Zhang, D.; Wu, C.; Shi, J.; Zhao, C.; Feng, H.; Ding, E.; Wang, J.; and Han, J. 2024. VDG: Vision-Only Dynamic Gaussian for Driving Simulation. *arXiv preprint arXiv:2406.18198*.
- Li, Z.; Niklaus, S.; Snavely, N.; and Wang, O. 2021. Neural Scene Flow Fields for Space-Time View Synthesis of Dynamic Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Li, Z.; Wang, Q.; Cole, F.; Tucker, R.; and Snavely, N. 2023. DynIBaR: Neural Dynamic Image-Based Rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Liu, J. Y.; Chen, Y.; Yang, Z.; Wang, J.; Manivasagam, S.; and Urtasun, R. 2023. Real-Time Neural Rasterization for Large Scenes. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, 8382–8393.
- Liu, L.; Gu, J.; Lin, K. Z.; Chua, T.-S.; and Theobalt, C. 2020. Neural Sparse Voxel Fields. *NeurIPS*.
- Lu, F.; Xu, Y.; Chen, G.-S.; Li, H.; Lin, K.-Y.; and Jiang, C. 2023. Urban Radiance Field Representation with Deformable Neural Mesh Primitives. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, 465–476.
- Luiten, J.; Kopanas, G.; Leibe, B.; and Ramanan, D. 2024. Dynamic 3D Gaussians: Tracking by Persistent Dynamic View Synthesis. In *3DV*.
- Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. *arXiv:2003.08934*.
- Müller, T.; Evans, A.; Schied, C.; and Keller, A. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.*, 41(4): 102:1–102:15.
- Niemeyer, M.; Barron, J. T.; Mildenhall, B.; Sajjadi, M. S.; Geiger, A.; and Radwan, N. 2022. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5480–5490.
- Ost, J.; Mannan, F.; Thurey, N.; Knodt, J.; and Heide, F. 2021. Neural Scene Graphs for Dynamic Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2856–2865.
- Park, K.; Sinha, U.; Barron, J. T.; Bouaziz, S.; Goldman, D. B.; Seitz, S. M.; and Martin-Brualla, R. 2021a. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5865–5874.
- Park, K.; Sinha, U.; Hedman, P.; Barron, J. T.; Bouaziz, S.; Goldman, D. B.; Martin-Brualla, R.; and Seitz, S. M. 2021b. HyperNeRF: A Higher-Dimensional Representation for Topologically Varying Neural Radiance Fields. *ACM Trans. Graph.*, 40(6).
- Pumarola, A.; Corona, E.; Pons-Moll, G.; and Moreno-Noguer, F. 2020. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Rabich, S.; Stotko, P.; and Klein, R. 2023. FPO++: Efficient Encoding and Rendering of Dynamic Neural Radiance Fields by Analyzing and Enhancing Fourier PlenOctrees. *ArXiv*, abs/2310.20710.
- Rematas, K.; Liu, A.; Srinivasan, P. P.; Barron, J. T.; Tagliasacchi, A.; Funkhouser, T.; and Ferrari, V. 2022. Urban Radiance Fields. *CVPR*.
- Shah, S.; Dey, D.; Lovett, C.; and Kapoor, A. 2018. Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In *Field and Service Robotics: Results of the 11th International Conference*, 621–635. Springer.
- Sun, P.; Kretzschmar, H.; Dotiwalla, X.; Chouard, A.; Patnaik, V.; Tsui, P.; Guo, J.; Zhou, Y.; Chai, Y.; Caine, B.; et al. 2020. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2446–2454.
- Sun, Y.; and Hariharan, B. 2023. Dynamo-Depth: Fixing Unsupervised Depth Estimation for Dynamical Scenes. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Tancik, M.; Casser, V.; Yan, X.; Pradhan, S.; Mildenhall, B.; Srinivasan, P. P.; Barron, J. T.; and Kretzschmar, H. 2022. Block-NeRF: Scalable Large Scene Neural View Synthesis. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 8238–8248.
- Taud, H.; and Mas, J.-F. 2018. Multilayer perceptron (MLP). *Geomatic approaches for modeling land change scenarios*, 451–455.
- Tonderski, A.; Lindström, C.; Hess, G.; Ljungbergh, W.; Svensson, L.; and Petersson, C. 2023. NeuRAD: Neural Rendering for Autonomous Driving. *arXiv preprint arXiv:2311.15260*.
- Tretschk, E.; Tewari, A.; Golyanik, V.; Zollhöfer, M.; Lassner, C.; and Theobalt, C. 2020. Non-Rigid Neural Radiance Fields: Reconstruction and Novel View Synthesis of a Dynamic Scene From Monocular Video. *arXiv:2012.12247*.

- Turki, H.; Ramanan, D.; and Satyanarayanan, M. 2021. Mega-NeRF: Scalable Construction of Large-Scale NeRFs for Virtual Fly-Throughs. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 12912–12921.
- Turki, H.; Zhang, J. Y.; Ferroni, F.; and Ramanan, D. 2023. SUDS: Scalable Urban Dynamic Scenes. In *Computer Vision and Pattern Recognition (CVPR)*.
- Wang, L.; Zhang, J.; Liu, X.; Zhao, F.; Zhang, Y.; Zhang, Y.; Wu, M.; Yu, J.; and Xu, L. 2022. Fourier PlenOctrees for Dynamic Radiance Field Rendering in Real-Time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 13524–13534.
- Wu, G.; Yi, T.; Fang, J.; Xie, L.; Zhang, X.; Wei, W.; Liu, W.; Tian, Q.; and Wang, X. 2024. 4D Gaussian Splatting for Real-Time Dynamic Scene Rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 20310–20320.
- Wu, T.; Zhong, F.; Tagliasacchi, A.; Cole, F.; and Öztireli, C. 2022. D2NeRF: Self-Supervised Decoupling of Dynamic and Static Objects from a Monocular Video. *ArXiv*, abs/2205.15838.
- Wu, Z.; Liu, T.; Luo, L.; Zhong, Z.; Chen, J.; Xiao, H.; Hou, C.; Lou, H.; Chen, Y.; Yang, R.; Huang, Y.; Ye, X.; Yan, Z.; Shi, Y.; Liao, Y.; and Zhao, H. 2023. MARS: An Instance-aware, Modular and Realistic Simulator for Autonomous Driving. *CICAI*.
- Xian, W.; Huang, J.-B.; Kopf, J.; and Kim, C. 2020. Space-time Neural Irradiance Fields for Free-Viewpoint Video. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 9416–9426.
- Xu, D.; Jiang, Y.; Wang, P.; Fan, Z.; Shi, H.; and Wang, Z. 2022. Sinnerf: Training neural radiance fields on complex scenes from a single image. In *European Conference on Computer Vision*, 736–753. Springer.
- Xu, H.; Alldieck, T.; and Sminchisescu, C. 2021. H-NeRF: Neural Radiance Fields for Rendering and Temporal Reconstruction of Humans in Motion. In *Neural Information Processing Systems*.
- Yan, Y.; Lin, H.; Zhou, C.; Wang, W.; Sun, H.; Zhan, K.; Lang, X.; Zhou, X.; and Peng, S. 2024. Street Gaussians: Modeling Dynamic Urban Scenes with Gaussian Splatting. In *ECCV*.
- Yang, J.; Ivanovic, B.; Litany, O.; Weng, X.; Kim, S. W.; Li, B.; Che, T.; Xu, D.; Fidler, S.; Pavone, M.; and Wang, Y. 2023a. EmerNeRF: Emergent Spatial-Temporal Scene Decomposition via Self-Supervision. *arXiv preprint arXiv:2311.02077*.
- Yang, L.; Kang, B.; Huang, Z.; Xu, X.; Feng, J.; and Zhao, H. 2024a. Depth Anything: Unleashing the Power of Large-Scale Unlabeled Data. In *CVPR*.
- Yang, Z.; Gao, X.; Zhou, W.; Jiao, S.; Zhang, Y.; and Jin, X. 2023b. Deformable 3D Gaussians for High-Fidelity Monocular Dynamic Scene Reconstruction. *arXiv preprint arXiv:2309.13101*.
- Yang, Z.; Yang, H.; Pan, Z.; and Zhang, L. 2024b. Real-time Photorealistic Dynamic Scene Representation and Rendering with 4D Gaussian Splatting. In *International Conference on Learning Representations (ICLR)*.
- Zhi, S.; Laidlow, T.; Leutenegger, S.; and Davison, A. J. 2021. In-Place Scene Labelling and Understanding with Implicit Scene Representation. In *ICCV*.
- Zhou, H.; Shao, J.; Xu, L.; Bai, D.; Qiu, W.; Liu, B.; Wang, Y.; Geiger, A.; and Liao, Y. 2024a. HUGS: Holistic Urban 3D Scene Understanding via Gaussian Splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 21336–21345.
- Zhou, X.; Lin, Z.; Shan, X.; Wang, Y.; Sun, D.; and Yang, M.-H. 2024b. Drivinggaussian: Composite gaussian splatting for surrounding dynamic autonomous driving scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 21634–21643.