# OmniRe: Omni Urban Scene Reconstruction

**Ziyu Chen**[*,¶]    **Jiawei Yang**[¶]    **Jiahui Huang**[§]    **Riccardo de Lutio**[§]
**Janick Martinez Esturo**[§]    **Boris Ivanovic**[§]    **Or Litany**[†,§]    **Zan Gojcic**[§]
**Sanja Fidler**[‡,§]    **Marco Pavone**[§,§]    **Li Song**[*]    **Yue Wang**[§,¶]

[*]Shanghai Jiao Tong University    [†]Technion    [‡]University of Toronto
[§]Stanford University    [§]NVIDIA Research    [¶]University of Southern California

## Abstract

We introduce `OmniRe`, a holistic approach for efficiently reconstructing high-fidelity dynamic urban scenes from on-device logs. Recent methods for modeling driving sequences using neural radiance fields or Gaussian Splatting have demonstrated the potential of reconstructing challenging dynamic scenes, but often overlook pedestrians and other non-vehicle dynamic actors, hindering a complete pipeline for dynamic urban scene reconstruction. To that end, we propose a comprehensive 3DGS framework for driving scenes, named `OmniRe`, that allows for accurate, full-length reconstruction of diverse dynamic objects in a driving log. `OmniRe` builds dynamic neural scene graphs based on Gaussian representations and constructs multiple local canonical spaces that model various dynamic actors, including vehicles, pedestrians, and cyclists, among many others. This capability is unmatched by existing methods. `OmniRe` allows us to holistically reconstruct different objects present in the scene, subsequently enabling the simulation of reconstructed scenarios with all actors participating in real-time (~60 Hz). Extensive evaluations on the Waymo dataset show that our approach outperforms prior state-of-the-art methods quantitatively and qualitatively by a large margin. We believe our work fills a critical gap in driving reconstruction. See the project page for code, video results and demos: ziyc.github.io/omnire.

## 1 Introduction

As autonomous driving increasingly adopts end-to-end models, the need for scalable and domain-gap-free simulation environments, where these systems can be evaluated in closed-loop, is becoming more evident. While the traditional way of using artist-generated assets is reaching its limits in terms of scale, diversity, and realism, the progress in data-driven methods for generating digital twins offers a strong alternative through reconstruction of simulation environments from on-device logs. Indeed, neural radiance fields (NeRFs) [26, 2, 48, 10, 47, 43] and Gaussian Splatting (GS) [17, 46] have emerged as powerful tools for reconstructing 3D scenes with high levels of visual and geometric fidelity. However, accurately and holistically reconstructing dynamic driving scenes remains a significant challenge, especially due to the complexity of real-world environments with diverse actors and types of motion.

Several works have already tried to tackle this challenge. Early methods typically ignore dynamic actors and focus only on reconstructing static parts of the scene [35, 25, 32, 10]. Subsequent works aim to reconstruct the dynamic scenes by either **(i)** modeling the scenes as a combination of a static and time-dependent neural field, where the decomposition of different scene parts is an emergent property [47, 38], or **(ii)** building a scene graph, in which dynamic actors and the static background are represented as nodes and reconstructed in their canonical frame and represented. The nodes of the scene graph are connected with edges that encode relative transformation parameters representing

(a) Scene Decomposition

(b) Corner Case Handling
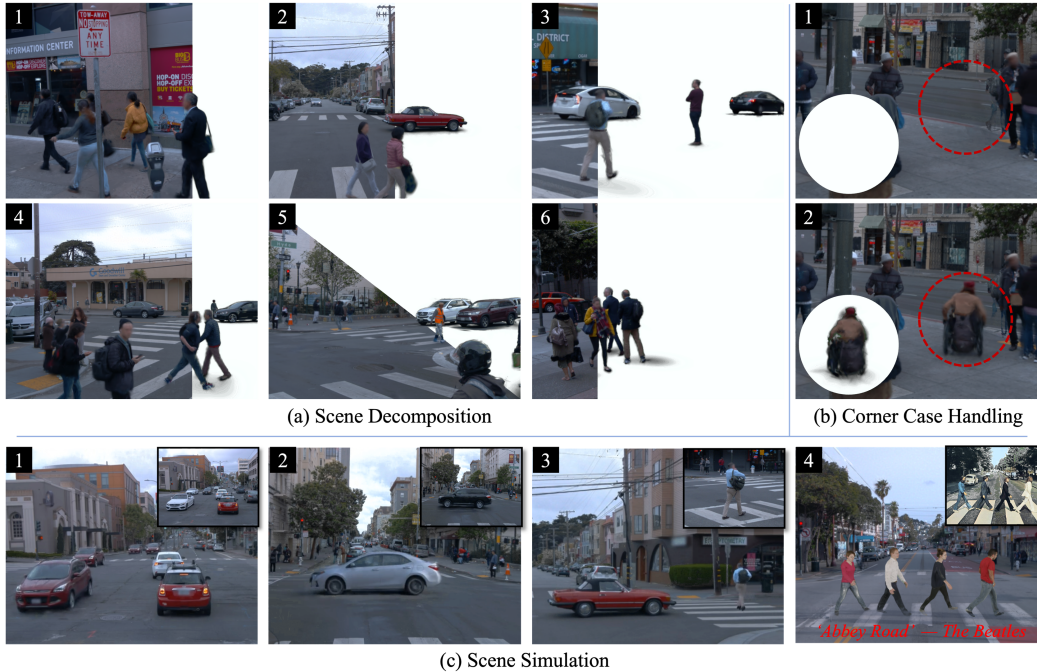
(c) Scene Simulation

Figure 1: (a) Decomposition of different parts of a scene. (b) Out-of-distribution categories that are overlooked by previous methods can be accurately handled by `OmniRe`. (c) `OmniRe` enables diverse applications including vehicle manipulation(c1, c2), traffic simulation (c3), human behavior simulation (c4), etc.

the motion of each actor through time [28, 19, 48, 43, 36, 7]. While the former is a more general formulation, the latter provides a higher level of editability and can be directly controlled with classical behavior models. However, the scene graph methods, still focus only on vehicles that can be represented as rigid bodies and thereby largely neglect other vulnerable road users (VRUs) such as pedestrians and cyclists that are critical in driving simulation.

To fill this critical gap, our work aims to model all dynamic actors, including vehicles, pedestrians, and cyclists, among many others. Unlike modeling objects from multi-view systems in a studio [11, 12, 33], reconstructing dynamic actors from in-the-wild scenarios is extremely challenging [5, 52, 53, 9, 16, 34]. Consider humans for example. Reconstructing humans from partial observations is a challenging problem on its own [16], and is made even more complex in driving scenarios with unfavorable distribution of the sensor observations and highly cluttered environments with frequent occlusions [39, 49, 40]. Indeed, even state-of-the-art human pose prediction models [9] often struggle to predict accurate poses, especially for pedestrians who are distant or occluded by others (*e.g.* Fig. 3). Additionally, there are other dynamic actors such as individuals on a wheelchair or pushing strollers, which cannot be modeled as simply with parametric models.

To address these mutually reinforced challenges, we propose an "omni" system capable of handling diverse actors. Our method `OmniRe` efficiently reconstructs high-fidelity dynamic driving scenes that encompass static backgrounds, driving vehicles, and non-rigidly moving dynamic actors (see Fig. 1). Specifically, we build a dynamic neural scene graph [28] based on Gaussian Splatting representation [17] and construct dedicated local canonical spaces for different dynamic actors. Driven by the "horses for courses" principle, `OmniRe` leverages the collective strengths of different representations: **(i)** vehicles are modeled as static Gaussians, transformed using rigid body transformations to simulate their motion over time, **(ii)** close-range walking pedestrians are fitted with a template-based SMPL model [23], enabling joint-level control using linear blend skinning weights, and **(iii)** far-range and other template-less dynamic actors are reconstructed using self-supervised deformation fields. This combination allows for accurate representation and controllable reconstruction of most objects of interest in the scene. More importantly, our representation is directly amenable to behavior and animation models that are commonly used in AV simulation (*e.g.*, Fig. 1-(c)).

To summarize, we make the following contributions:

- We introduce `OmniRe`, a holistic framework for dynamic driving scene reconstruction that embodies the "omni" principle in terms of actor coverage and representation flexibility. `OmniRe` leverages dynamic neural scene graphs based on Gaussian representations to unify the reconstruction of static backgrounds, driving vehicles, and non-rigidly moving dynamic actors (§ 4). It enables high-fidelity scene reconstruction, sensor simulation from novel viewpoints, and controllable scenario editing in real-time (§ 5).

- We address the challenges of modeling humans and other dynamic actors from driving logs such as occlusion, cluttered environments, and the limitations of existing human pose prediction models (§ 4.2). Our findings are based on AV scenes, but can be generalized to other domains.

- We perform extensive experiments and ablations to demonstrate the benefits of our holistic representation. `OmniRe` achieves state-of-the-art performance in scene reconstruction and novel view synthesis (NVS), significantly outperforming previous methods in terms of full image metrics (+1.88 PSNR for reconstruction and +2.38 PNSR for NVS). The differences are pronounced for dynamic actors, such as vehicles (+1.18 PSNR), and humans (+4.09 PSNR for reconstruction and +3.06 PSNR for NVS) (Tab. 1).

## 2 Related Work

**Dynamic Scene Reconstruction.** Neural representations are dominating novel view synthesis [26, 3, 2, 27, 8, 17]. These have been extended in different ways to enable dynamic scene reconstruction. *Deformation-based* approaches [31, 29, 37, 30, 4] and recently DeformableGS [50] and [41] propose to model dynamic scenes using a 3D neural representation for the canonical space, coupled with a deformation network mapping time-dependent observations to canonical deformations. These are generally limited to small scenes with limited movement, making them inadequate for challenging urban dynamic scenes. *Modulation-based* techniques operate by directly feeding the image time-stamps (or latent codes) as an additional input to a neural representation [44, 22, 21, 24]. However, this generally results in an underconstrained formulation, therefore requiring additional supervision, such as depth and optical flow (Video-NeRF [44] and NSFF [22]), or multi-view inputs captured from synchronized cameras (DyNeRF [21] and Dynamic3DGS [24]). $D^2$NeRF [42] proposed to expand on this formulation by partitioning the scene into static and dynamic fields. Following this, SUDS [38] and EmerNeRF [47] have shown impressive reconstruction ability for dynamic autonomous driving scenes. However, they model all dynamic elements using a single dynamic field, rather than modeling each separately, thus they lack controllability, limiting their practicality as sensor simulators. *Explicit decomposition* of the scene into separate agents enables controlling them individually. These agents can be represented as bounding boxes in a scene graph as in Neural Scene Graphs [28] (NSG) that is widely adopted in UniSim [48], MARS [43], NeuRAD [36], ML-NSG [7] and recent Gaussian-based works StreetGaussians [46], DrivingGaussians [55], and HUGS [54]. However, these approaches handle only rigid objects due to limitations of time-independent representations [28, 43, 48, 55, 54, 46, 36, 7] or limitations of deformation-based techniques [50, 14]. To address them, `OmniRe` proposes a Gaussian scene graph that incorporates various Gaussian representations for both rigid and non-rigid objects, providing extra flexibility and controllability for diverse actors.

**Human Reconstruction.** Human bodies have variable appearance and complex motions, calling for dedicated modeling techniques. NeuMan [15] proposes to employ the SMPL body model [23] to warp ray points to a canonical space. This approach enables the reconstruction of non-rigid human bodies and warrants fine control. Similarly, recent works such as GART [20], GauHuman [13] and HumanGaussians [18] have combined the Gaussian representation and the SMPL model. However, these methods are not directly applicable in-the-wild. As for recovering human dynamics in driving scenes, [49] focuses on shape and pose reconstruction for LiDAR simulation, while [40] and [39] aim to recreate natural and accurate human motion from partial observations. However, these methods focus solely on shape and pose estimation and are limited in appearance modeling. In contrast, our method not only models human appearance but also integrates this modeling within a holistic scene framework, to achieve a comprehensive solution. Urban scenes typically involve numerous pedestrians, with sparse observation, often accompanied by severe occlusion. We analyze these challenges in detail and address them in § 4.2.

# 3 Preliminaries

**3D Gaussian Splatting.** First introduced in [17], 3D Gaussian Splatting (3DGS) represents scenes via a set of colored blobs $\mathcal{G} = \{g\}$ whose intensity distribution is a Gaussian. Each Gaussian (blob) $g = (o, \boldsymbol{\mu}, \mathbf{q}, \boldsymbol{s}, \boldsymbol{c})$ contains the following attributes: opacity $o \in (0, 1)$, mean position $\boldsymbol{\mu} \in \mathbb{R}^3$, rotation $\mathbf{q} \in \mathbb{R}^4$ represented as a quaternion, anisotropic scaling factors $\boldsymbol{s} \in \mathbb{R}^3_+$, and view-dependent colors $\boldsymbol{c} \in \mathbb{R}^F$ represented as spherical harmonics (SH) coefficients. To compute the color $C$ of a pixel, Gaussians overlapping with this pixel are sorted by their distance to the camera center (sorted by $i \in \mathcal{N}$) and $\alpha$-blended:

$$C = \sum_{i \in \mathcal{N}} \boldsymbol{c}_i \alpha_i \prod_{j=1}^{i-1}(1 - \alpha_j). \tag{1}$$

Where $\alpha_i$ is computed as $\alpha_i = o_i \exp(-\frac{1}{2}(\mathbf{p} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{p} - \boldsymbol{\mu}_i))$, $\boldsymbol{\Sigma}_i$ is the 2D projection covariance. We further define the application of a rigid (affine) transformation $\mathbf{T} = (\mathbf{R}, \mathbf{t}) \in \mathbb{SE}(3)$ to all Gaussians in the set as:

$$\mathbf{T} \otimes \mathcal{G} = (o, \mathbf{R}\boldsymbol{\mu} + \mathbf{t}, \mathrm{Rot}(\mathbf{R}, \mathbf{q}), \boldsymbol{s}, \boldsymbol{c}), \tag{2}$$

where $\mathrm{Rot}(\cdot)$ denotes rotating the quaternion by the rotation matrix.

**Skinned Multi-Person Linear (SMPL) Model.** SMPL [23] is a parametric human body model that combines the advantages of a parametric mesh with linear blending skinning (LBS) to manipulate body shape and pose. At its core, SMPL uses a template mesh $\mathcal{M}_h = (\mathcal{V}, \mathcal{F})$ defined in a canonical rest pose, parameterized by $n_v$ vertices $\mathcal{V} \in \mathbb{R}^{n_v \times 3}$. The template mesh can be shaped and transformed using shape parameters $\boldsymbol{\beta}$ and pose parameters $\boldsymbol{\theta}$:

$$\mathcal{V}_S = \mathcal{V} + B_S(\boldsymbol{\beta}) + B_P(\boldsymbol{\theta}), \tag{3}$$

where $B_S(\boldsymbol{\beta}) \in \mathbb{R}^{n_v \times 3}$ and $B_P(\boldsymbol{\theta}) \in \mathbb{R}^{n_v \times 3}$ are the $xyz$ offsets to individual vertices [18] and $\mathcal{V}_S$ are the vertex locations in the shaped space.

To further deform the vertices $\mathcal{V}_S$ to achieve the desired pose $\boldsymbol{\theta}'$, SMPL utilizes pre-defined LBS weights $\boldsymbol{W} \in \mathbb{R}^{n_k \times n_v}$ and the joint transformations $\boldsymbol{G}$ to define the deformation of each vertex $\boldsymbol{v}_i$: $\boldsymbol{v}_i' = \left(\sum_k \boldsymbol{W}_{k,i} \boldsymbol{G}_k\right) \boldsymbol{v}_i$, where $n_k$ is the number of joints, and the joint transformations $\boldsymbol{G}$ are derived from the source pose $\boldsymbol{\theta}$, the target pose $\boldsymbol{\theta}'$ and shape $\boldsymbol{\beta}$. The pose parameters include the body pose component $\boldsymbol{\theta}_b \in \mathbb{R}^{23 \times 3 \times 3}$ and the global orientation component $\boldsymbol{\theta}_g \in \mathbb{R}^{3 \times 3}$. For more details of SMPL, we refer readers to [23]. Our method obtains pose parameters $\boldsymbol{\theta}$ for each pedestrian across all frames, as well as their individual shape parameters $\boldsymbol{\beta} \in \mathbb{R}^{10}$, these pose sequences initialize the non-rigid dynamics of pedestrians. The detailed process is described in § 4.2.

# 4 Method

As overviewed in Fig. 2, we build a *Gaussian Scene Graph* representation that holistically covers both the static background and diverse *movable* entities. We discuss our systematic strategy for representing different semantic classes in § 4.1, highlighting one of our primary contributions. Modeling humans in unconstrained environments is particularly challenging due to the non-rigid nature of the human body, the difficulty of accurate initialization, and the severe occlusions often present in the wild. We present our approach to this problem in § 4.2, which significantly boosts the performance. Lastly, we show how the scene representation is end-to-end optimized to obtain faithful and controllable reconstructions in § 4.3.

## 4.1 Dynamic Gaussian Scene Graph Modeling

**Gaussian Scene Graph.** To allow for flexible control of diverse *movable* objects in the scene without sacrificing reconstruction quality, we opt for a *Gaussian Scene Graph* representation. Our scene graph is composed of the following nodes: (1) a *Sky Node* representing the sky that is far away from the ego-car, (2) a *Background Node* representing the static scene background such as buildings, roads, and vegetation, (3) a set of *Rigid Nodes*, each representing a rigidly movable object such as a vehicle, (4) a set of *Non-rigid Nodes* that model pedestrians or cyclists. Nodes of type (2,3,4) can be converted directly into world-space Gaussians which we will introduce next. These Gaussians are concatenated and rendered using the rasterizer proposed in [17]. The Sky Node is represented by an optimizable
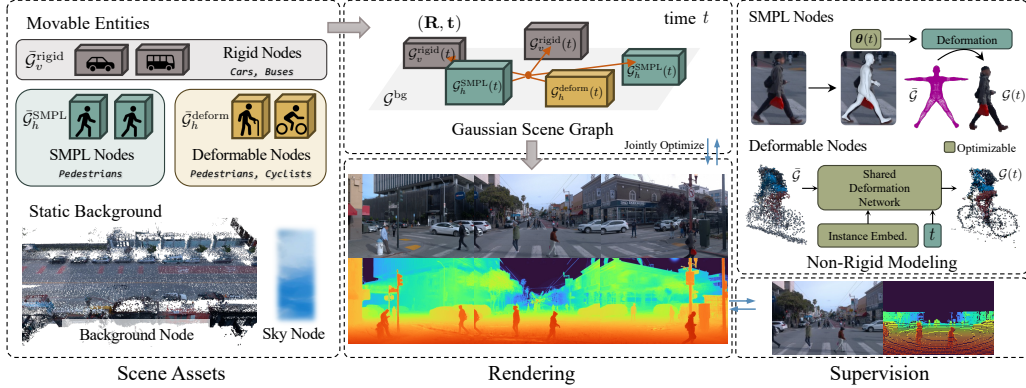
Figure 2: **Method Overview.** Gaussians of all foreground models are defined in their local or canonical spaces. At a given time $t$, the Gaussians are deformed and transformed into the world space, forming a Gaussian scene graph together with background Gaussians to model the entire scene. The Gaussians in the scene graph are rasterized to render images and depth, and are jointly optimized using reconstruction losses. We utilize SMPL Gaussians to model non-rigid human bodies and deformable Gaussians to handle out-of-distribution non-rigid categories.

environment texture map, similar to [6], rendered separately, and composited with the rasterized Gaussian image with simple alpha blending.

**Background Node.** The background node is represented by a set of static Gaussians $\mathcal{G}^{\text{bg}}$. These Gaussians are initialized by accumulating the LiDAR points and additional points generated randomly in accordance with the strategy described in [6].

**Rigid Nodes.** Gaussians representing the vehicles (*e.g.* cars or trucks) are defined as $\bar{\mathcal{G}}_v^{\text{rigid}}$ in the object's local space (denoted by the upper bar), where $v$ is the index of the vehicle/node. While the Gaussians within a vehicle will not change over time in the local space, the positions of Gaussians in world space will change according to the vehicle's pose $\mathbf{T}_v \in \mathbb{SE}(3)$. At a given time $t \in \mathbb{R}$, the Gaussians are transformed into world space by simply applying the pose transformation:

$$\mathcal{G}_v^{\text{rigid}}(t) = \mathbf{T}_v(t) \otimes \bar{\mathcal{G}}_v^{\text{rigid}}. \tag{4}$$

**Non-Rigid Nodes.** Unlike rigid vehicles, non-rigid dynamic classes such as pedestrians and cyclists, which are all human-related, require extra consideration of both their global movements in world space and their continuous deformations in local space to accurately reconstruct their dynamics. To enable a reconstruction that fully explains the underlying geometry, we further subdivide the non-rigid nodes into two categories: *SMPL Nodes* for walking or running pedestrians with SMPL templates that enable joint-level control and *Deformable Nodes* for out-of-distribution non-rigid instances (such as cyclists and other template-less dynamic entities).

**Non-Rigid SMPL Nodes.** As introduced in § 3, SMPL provides a parametric way of representing human poses and deformations, and we hence use the model parameters $(\boldsymbol{\theta}(t), \boldsymbol{\beta})$ to drive the 3D Gaussians within the nodes. Here $\boldsymbol{\theta}(t) \in \mathbb{R}^{24 \times 3 \times 3}$ represents the human posture that changes over time $t$. For each node indexed by $h$, We tessellate the SMPL template mesh $\mathcal{M}_h$ instantiated from the resting pose (the *'Da'* pose) with 3D Gaussians $\bar{\mathcal{G}}_h^{\text{SMPL}}$ using a strategy similar to GART [20], where each Gaussian is binded to its corresponding vertex of $\mathcal{M}_h$. The world-space Gaussians for each node can be then computed as:

$$\mathcal{G}_h^{\text{SMPL}}(t) = \mathbf{T}_h(t) \otimes \text{LBS}(\boldsymbol{\theta}(t), \bar{\mathcal{G}}_h^{\text{SMPL}}). \tag{5}$$

Here $\mathbf{T}_h(t) \in \mathbb{SE}(3)$ is the global pose of the node at time $t$, and $\text{LBS}(\cdot)$ is the linear blend skinning operation that deforms the Gaussians according to the SMPL pose parameters. In order to compute the LBS operator, one first precomputes the skinning weights of each Gaussian in $\bar{\mathcal{G}}_h^{\text{SMPL}}$ w.r.t. the SMPL key joints. Once $\boldsymbol{\theta}$ changes over time, the key joints' transformations are updated and linearly interpolated onto the Gaussians to obtain the final deformed positions and rotations, while other attributes in the Gaussian remain unchanged. Crucially, it is highly challenging to accurately optimize the SMPL poses $\boldsymbol{\theta}(t)$ from scratch just with sensor observations (even for single-person or indoor

5

scenarios [15, 20, 18]). Hence a rough initialization of $\boldsymbol{\theta}(t)$ is typically needed, whose details are deferred to a dedicated section § 4.2.

**Non-Rigid Deformable Nodes.** These nodes act as a fallback option for out-of-distribution non-rigid instances, e.g. extremely faraway pedestrians which even state-of-the-art SMPL predictors cannot provide an accurate estimation; or long-tail template-less non-rigid instances. Hence, we propose to use a general deformation network $\mathcal{F}_\varphi$ with parameter $\varphi$ to fit the non-rigid motions within the nodes. Specifically, for node $h$, the world-space Gaussians are defined as:

$$\mathcal{G}_h^{\text{deform}}(t) = \mathbf{T}_h(t) \otimes \left( \bar{\mathcal{G}}_h^{\text{deform}} \oplus \mathcal{F}_\varphi(\bar{\mathcal{G}}_h^{\text{deform}}, \boldsymbol{e}_h, t) \right), \tag{6}$$

where the deformation network generates the changes of the Gaussian attributes from time $t$ to the canonical space Gaussians $\bar{\mathcal{G}}_h^{\text{deform}}$, outputting the changes in position $\delta\boldsymbol{\mu}_h(t)$, rotation $\delta\mathbf{q}_h(t)$, and the scaling factors $\delta\boldsymbol{s}_h(t)$. The changes are applied back to $\bar{\mathcal{G}}_h^{\text{deform}}$ with the $\oplus$ operator that internally performs a simple arithmetic addition that results in $(o, \boldsymbol{\mu} + \delta\boldsymbol{\mu}(t), \mathbf{q} + \delta\mathbf{q}(t), \boldsymbol{s} + \delta\boldsymbol{s}(t), \boldsymbol{c})$. Notably, previous approaches such as [50] utilizes a single deformation network for the entire scene, and usually fail in highly complex outdoor dynamic scenes with rapid movements. On the contrary, in our work, we define a per-node deformation field which has much more representation power. To maintain computational efficiency, the network weights $\varphi$ are shared and the identities of the nodes are disambiguated via an instance embedding parameter $\boldsymbol{e}_h$. Experimental results in § 5.2 show that deformable Gaussians are essential for achieving good reconstruction quality.

**Sky Node.** We follow [6, 47] to use a separate environmental map to fit the sky color from viewing directions. Compositing the sky image $C_{\text{sky}}$ with the rendered Gaussians $C_{\mathcal{G}}$ consisting of $(\mathcal{G}^{\text{bg}}, \{\mathcal{G}_v^{\text{rigid}}\}, \{\mathcal{G}_h^{\text{SMPL}}\}, \{\mathcal{G}_h^{\text{deform}}\})$, we obtain the final rendering as:

$$C = C_{\mathcal{G}} + (1 - O_{\mathcal{G}})C_{\text{sky}}, \tag{7}$$

where $O_{\mathcal{G}} = \sum_{i=1}^{N} \alpha_i \prod_{j=1}^{i-1}(1 - \alpha_j)$ is the rendered opacity mask of Gaussians.

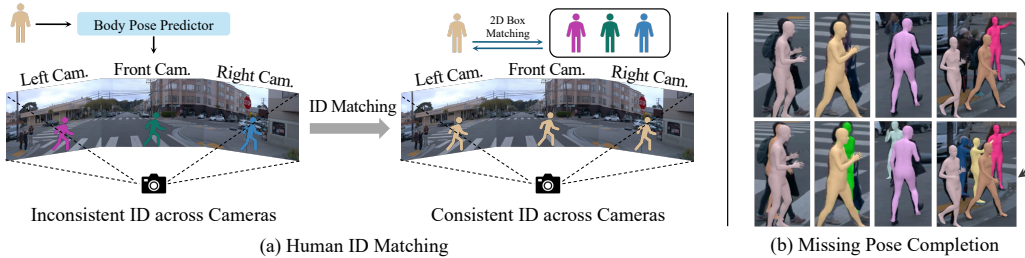## 4.2 Reconstructing In-the-Wild Humans



Figure 3: **Human Pose Processing.** (a) Human ID matching ensures consistent identification across cameras. (b) Missing pose completion to recover poses of occluded individuals.

To initialize the parameters $\boldsymbol{\theta}(t)$ of the non-rigid SMPL nodes, we extend an off-the-shelf predictor 4D-Humans [9] which estimates human body poses from raw video inputs. However, it presents several practical limitations that hinder its usability in our context. We discuss and address these challenges with following modules to predict accurate and temporally consistent human body poses from multi-view images captured in-the-wild, in the presence of frequent occlusions.

**Human ID Matching.** 4D-Humans [9] is designed to process single-camera videos only. In our multi-camera setup, this limitation leads to missed connections between the same person across views (Fig. 3(a)). To address this, we match the estimated poses of detected humans to the ground truth IDs in the dataset using the mean Intersection-over-Union (mIoU) between detections and GT boxes, ensuring that each pedestrian is consistently identified across multiple cameras.

**Missing Pose Completion.** 4D-Humans [9] struggles to predict SMPL poses for occluded individuals, which are common in AV scenes, leading to missing predictions. We recover missing poses by interpolating poses from neighboring frames. As visualized in Fig. 3(b), this process enables recovering accurate poses for occluded individuals, thus achieving temporally complete pose sequences.

6

**Scene-Pose Alignment.** As a general model designed to be camera-agnostic, 4D-Humans assumes a virtual camera with fixed parameters for all video inputs. In contrast, real cameras have different parameters. This leads to misalignment between the scale and position of predicted poses and the real-world coordinate systems. We utilize the box size and position data available for each individual to correct the scale and position of predicted poses.

**Pose Refinement.** Errors from the pose predictor, interpolation, and alignment estimation result in noisy human body poses. We utilize these noisy poses to initialize the dynamics of SMPL nodes and jointly refine the per-frame poses $\boldsymbol{\theta}(t)$ of each individual during training by optimizing with reconstruction losses. Our ablation studies (§ 5.2) show that human body pose refinement is crucial for improving reconstruction quality and pose accuracy.

### 4.3 Optimization

We simultaneously optimize all the parameters as mentioned in § 4.1 in *a single stage* to reconstruct the entire scene. These parameters include: **(1)** all the Gaussian attributes (opacity, mean positions, scaling, rotation, and appearance) in their local spaces, namely $\mathcal{G}^{\text{bg}}, \{\bar{\mathcal{G}}_v^{\text{rigid}}\}, \{\bar{\mathcal{G}}_h^{\text{SMPL}}\}, \{\bar{\mathcal{G}}_h^{\text{deform}}\}$, **(2)** the poses of both rigid and non-rigid nodes for each frame $t$, i.e., $\{\mathbf{T}_v(t)\}, \{\mathbf{T}_h(t)\}$, **(3)** the human poses of all the SMPL nodes for each frame $t$, i.e., $\{\boldsymbol{\theta}(t)\}$, and the corresponding skinning weights, **(4)** the weight $\varphi$ of the deformation network $\mathcal{F}$, **(5)** the weight of the sky model.

We use the following objective function for optimization:

$$\mathcal{L} = (1 - \lambda_r)\,\mathcal{L}_1 + \lambda_r \mathcal{L}_{\text{SSIM}} + \lambda_{\text{depth}}\mathcal{L}_{\text{depth}} + \lambda_{\text{opacity}}\mathcal{L}_{\text{opacity}} + \mathcal{L}_{\text{reg}}, \tag{8}$$

where $\mathcal{L}_1$ and $\mathcal{L}_{\text{SSIM}}$ are the L1 and SSIM losses on rendered images, $\mathcal{L}_{\text{depth}}$ compares the rendered depth of Gaussians with sparse depth signals from LiDAR, $\mathcal{L}_{\text{opacity}}$ encourages the opacity of the Gaussians to align with the non-sky mask, and $\mathcal{L}_{\text{reg}}$ represents various regularization terms applied to different Gaussian representations. Detailed descriptions of loss terms are provided in the Appendix.

## 5 Experiments

**Dataset.** We conduct experiments on the Waymo Open Dataset [33], which comprises real-world driving logs. However, most of these logs depict scenes with relatively simple dynamics and rarely focus on non-rigid classes. Therefore, we select eight highly complex dynamic scenes that, in addition to typical vehicles, include diverse dynamic classes such as pedestrians and cyclists. Each selected segment contains approximately 150 frames. The segment IDs are listed in Tab. 8.

**Baselines.** We compare our method against several Gaussian Splatting approaches: 3DGS [17], DeformableGS [50], StreetGS [46], HUGS [54], and PVG [6]. Additionally, we compare our method with NeRF-based approach EmerNeRF [47]. Among methods compared, for StreetGS [46], we use our own reimplementation. For 3DGS [17] and DeformableGS [50], we use the implementation with LiDAR supervision to ensure the comparison fairness. For other methods, we use their official code. For training, we utilize data from the three front-facing cameras, resized to a resolution of $640{\times}960$ for all methods, along with LiDAR data for supervision. We utilize the instance bounding boxes provided by the dataset to transform objects and refine them via pose optimization during training. For further implementation details, please refer to Appendix.

### 5.1 Main Results

**Appearance.** We evaluate our method on scene reconstruction and novel view synthesis (NVS) tasks, using every 10th frame as the held-out test set for NVS. We report PSNR and SSIM scores for full images, as well as human-related and vehicle-related regions, to assess dynamic reconstruction capabilities. The quantitative results in Tab. 1 show that `OmniRe` outperforms all other methods, with a significant margin in human-related regions, validating our holistic modeling of dynamic actors. Additionally, while StreetGS [46] and our method model vehicles in a similar way, we observe that `OmniRe` is slightly better than StreetGS even in vehicle regions. This is due to the absence of human modeling in StreetGS, which allows supervision signals from human regions (e.g., colors, LiDAR depth) to incorrectly influence vehicle modeling. The issues StreetGS faces are one of our motivations for modeling almost everything in a scene holistically, aiming to eliminate erroneous supervision and unintended gradient propagation.

Table 1: **Comparison on Waymo Open Dataset.** We compute PSNR and SSIM for both the full image and dynamic regions. *Vehicle* indicates regions corresponding to vehicle-related classes, while *Human* indicates regions corresponding to human-related classes. *Box* indicates methods that utilize bounding boxes for dynamic modeling. *LiDAR* means method using LiDAR information.

| | | | Scene Reconstruction | | | | | | Novel View Synthesis | | | | | |
| | | | Full Image | | Human | | Vehicle | | Full Image | | Human | | Vehicle | |
| Methods | Box | LiDAR | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| EmerNeRF[47] | | ✓ | 31.93 | 0.902 | 22.88 | 0.578 | 24.65 | 0.723 | 29.67 | 0.883 | 20.32 | 0.454 | 22.07 | 0.609 |
| 3DGS[17] | | ✓ | 26.00 | 0.912 | 16.88 | 0.414 | 16.18 | 0.425 | 25.57 | 0.906 | 16.62 | 0.387 | 16.00 | 0.407 |
| DeformGS[50] | | ✓ | 28.40 | 0.929 | 17.80 | 0.460 | 19.53 | 0.570 | 27.72 | 0.922 | 17.30 | 0.426 | 18.91 | 0.530 |
| PVG[6] | | ✓ | 32.37 | 0.937 | 24.06 | 0.703 | 25.02 | 0.787 | 30.19 | 0.919 | 21.30 | 0.567 | 22.28 | 0.679 |
| HUGS[54] | ✓ | ✓ | 28.26 | 0.923 | 16.23 | 0.404 | 24.31 | 0.794 | 27.65 | 0.914 | 15.99 | 0.378 | 23.27 | 0.748 |
| StreetGS[46] | ✓ | ✓ | 29.08 | 0.936 | 16.83 | 0.420 | 27.73 | 0.880 | 28.54 | 0.928 | 16.55 | 0.393 | 26.71 | 0.846 |
| Ours | ✓ | ✓ | **34.25** | **0.954** | **28.15** | **0.845** | **28.91** | **0.892** | **32.57** | **0.942** | **24.36** | **0.727** | **27.57** | **0.858** |



(a) Ground Truth      (b) Ours

(c) EmerNeRF      (d) StreetGS
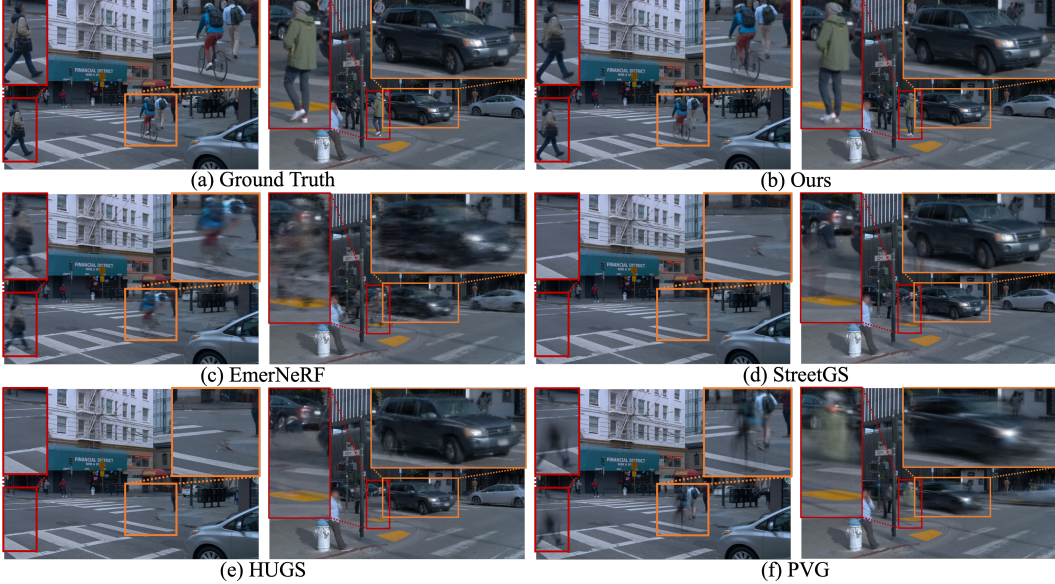
(e) HUGS      (f) PVG

Figure 4: **Qualitative Comparison of Novel View Synthesis.** The insets highlight the details of the reconstructed dynamic objects. `OmniRe` manages to recover very fine details, achieving high-quality reconstruction of various common dynamic objects, including vehicles, pedestrians, and cyclists.

In addition, we show visualizations in Fig. 4 to assess model performance qualitatively. Although PVG [6] performs well on the scene reconstruction task, it struggles with the novel view synthesis task in highly dynamic scenes, resulting in blurry dynamic objects in novel views (Fig. 4-(f)). HUGS [54] (Fig. 4-(e)), StreetGS [46](Fig. 4-(d)) and 3DGS [17] (Fig. 10-(h)) fail to recover the pedestrians because they are not capable of modeling non-rigid objects. DeformableGS [50] (Fig. 10-(g)) suffers from extreme motion blur for outdoor dynamic scenes with rapid movements, despite achieving reasonable performance for indoor scenes and cases with small motion. EmerNeRF [47] reconstructs coarse structures of moving humans and vehicles to a certain level, but struggles with fine-grained details(Fig. 4-(c)). In contrast to all these methods in comparison, our method faithfully reconstructs fine details for any object in the scene, handling occlusion, deformation, and extreme motion. We recommend readers to check our project page for video comparisons of these methods.

**Geometry.** In addition to appearance, we also investigate whether our method can reconstruct fine geometry of urban scenes. We evaluate the Root Mean Squared Error (RMSE) and two-way Chamfer Distances (CD) for LiDAR depth reconstruction on both training frames and novel frames. Details about evaluation procedures are provided in the Appendix. Tab. 4 reports the results. Our method outperforms others by a large margin. Fig. 5 illustrates the accurate reconstruction of dynamic actors achieved by our method in comparison to other approaches.
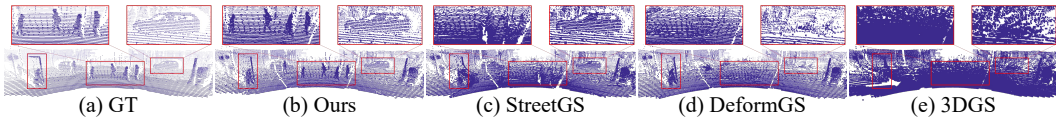
## 5.2 Ablation Studies & Applications

Figure 5: **Visualizations of Rendered LiDAR.** Our method accurately reconstructs LiDAR data for humans and vehicles compared to other approaches.
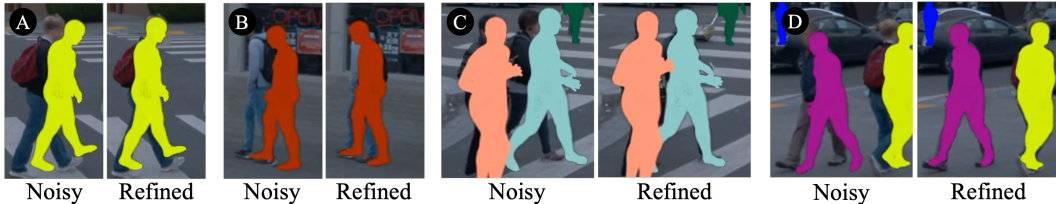


| Noisy | Refined | Noisy | Refined | Noisy | Refined | Noisy | Refined |

Figure 6: **Ablation of Human Body Pose Refinement.**



| w/o Deformed Actors | w/ Deformed Actors | w/o SMPL | w/ SMPL |

Figure 7: **Ablation of Human Modeling.**

**SMPL Modeling.** SMPL modeling is important to model the local, continuous movements of humans. We study its impact by disabling the human pose transformation enabled by SMPL and report the results in Tab. 3 ((a) v.s.

Table 2: **Ablation on GT Boxes Refinement.**

|  | Full PSNR | | Human PSNR | | Vehicle PSNR | |
| --- | --- | --- | --- | --- | --- | --- |
|  | Recon. | NVS | Recon. | NVS | Recon. | NVS |
| Complete Model | **34.25** | **32.57** | **28.15** | **24.36** | **28.91** | **27.57** |
| w/o Box Refine. | 33.04 | 31.72 | 26.53 | 23.67 | 25.57 | 24.78 |

(b)) and illustrate these effects in Fig. 7-(B). Without template-based modeling, the reconstructed human renderings appear highly blurred, particularly around the legs, thus failing to accurately reconstruct human body movements. This contrasts sharply with the precise leg reconstruction observed in our default setting. Moreover, SMPL modeling provides joint-level control, improving the controllability (Fig. 1-(c,3), (c,4)).

**Human Body Pose Refinement.** As discussed in (§ 4.2), the human body pose predictions suffer from prediction errors and scale ambiguity, which subsequently leads to significant misalignment errors as shown in Fig. 6 (Noisy). We improve this by jointly optimizing the human poses and Gaussians via the same reconstruction losses. Tab. 3-(a) v.s. (c) ablates this design choice, and Fig. 6 showcases the refined poses. These results verify the effectiveness of our refinement strategy.

**Deformable Nodes.** Deformable nodes are important for accurately reconstructing out-of-distribution or template-less actors. Our approach addresses this challenge by learning a self-supervised deformation field that transforms Gaussians from their canonical space to the shape space. Tab. 3 ((a) v.s. (d)) proves the importance of this component. In Fig. 7-(A) shows that without deformable nodes, some dynamic actors are either ignored or incorrectly blended into the background.

Table 3: **Ablation on Non-Rigid Modeling.**

|  | Full PSNR | | Human PSNR | |
| --- | --- | --- | --- | --- |
|  | Recon. | NVS | Recon. | NVS |
| (a) Ours default | **34.25** | **32.57** | **28.15** | **24.36** |
| (b) w/o SMPL actors | 32.80 | 31.76 | 24.71 | 23.18 |
| (c) w/o Body pose refine | 33.84 | 32.44 | 26.97 | 24.04 |
| (d) w/o Deformed actors | 33.64 | 32.17 | 25.26 | 22.41 |

Table 4: **Evaluation of LiDAR Depth Accuracy.**

| Methods | Training Frames | | Novel Frames | |
| --- | --- | --- | --- | --- |
|  | CD↓ | RMSE↓ | CD↓ | RMSE↓ |
| 3DGS[17] | 0.415 | 2.804 | 0.467 | 2.896 |
| DeformableGS[50] | 0.384 | 2.965 | 0.383 | 2.990 |
| StreetGS[46] | 0.274 | 2.199 | 0.286 | 2.228 |
| Ours | **0.242** | **1.894** | **0.244** | **1.909** |

9

**Boxes Refinement.** In practice, we observe that the instance bounding boxes provided by the dataset are imprecise. These noisy ground truth boxes can be harmful to rendering quality. To address this, we jointly refine the bounding box parameters during training. Tab. 2 and Fig. 12 show the practical benefits of this simple yet effective step, which results in improved numeric metrics and reduced blurriness of foreground objects.

**Applications to Simulation.** Thanks to the decomposition nature of `OmniRe`, each instance is modeled separately. After joint training, we obtain reconstructed assets that can be flexibly edited in terms of position and rotation. Beyond editing within a single scene, we can also transfer assets from one scene to another, adding variety and complexity to the reconstructed environments. Fig. 1(c,left) demonstrates a swap of the black vehicle originally in the scene (inset) with a reconstructed vehicle from another scene; and (c,right) an insertion of a pedestrian from the scene in the inset to a new scene. Additional car swap edits are shown in Fig. 11. Through explicit modeling of pedestrians and other non-rigid individuals, we achieve the simulation of reenacted scenarios involving detailed pedestrian-vehicle interaction. As demonstrated in Fig. 8, we simulate a moving vehicle stopping at a crossing, waiting for a pedestrian who slowly crosses. The pedestrian is reconstructed from another scene. This simulation of humans is extremely challenging for previous methods.



Figure 8: A sample of human-vehicle interaction simulation in driving scenarios. For video demos, we refer readers to visit our project page.

# 6 Conclusion

Our method, `OmniRe`, tackles comprehensive urban scene modeling using Gaussian Scene Graphs. It achieves fast, high-quality reconstruction and rendering, suggesting promise for autonomous driving and robotics simulation. We also present solutions for human modeling in complex environments. Future work includes self-supervised learning, improved scene representations, and safety/privacy considerations.

**Broader impact.** Our method aims to address a significant problem in autonomous driving—simulation. This approach has the potential to enhance the development and testing of autonomous vehicles, potentially leading to safer and more efficient AV systems. Simulation, in a safe and controllable manner, remains an open and challenging research question.

**Ethics&Privacy.** Our work does not involve the collection or annotation of new data. We utilize well-established public datasets that adhere to strict ethical guidelines. These datasets ensure that sensitive information, including identifiable human features, is blurred or anonymized to protect individual privacy. We are committed to ensuring that our method, as well as future applications, are employed responsibly and ethically to maintain safety and preserve privacy.

**Limitations.** `OmniRe` still has some limitations. First, our method does not explicitly model lighting effects, which may lead to visual harmony issues during simulations, particularly when combining elements reconstructed under varying lighting conditions. Addressing this non-trivial challenge requires dedicated efforts beyond the scope of our current work. Further research into modeling light effects and enhancing simulation realism remains crucial for achieving more convincing and harmonious results. Second, similar to other per-scene optimization methods, `OmniRe` produces less satisfactory novel views when the camera deviates significantly from the training trajectories. We believe that incorporating data-driven priors, such as image or video generative models, represents a promising direction for future exploration.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.

[3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.

[4] Hongrui Cai, Wanquan Feng, Xuetao Feng, Yan Wang, and Juyong Zhang. Neural surface reconstruction of dynamic scenes with monocular rgb-d camera. In *Thirty-sixth Conference on Neural Information Processing Systems (NeurIPS)*, 2022.

[5] Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11594–11604, 2021.

[6] 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*, 2023.

[7] Tobias Fischer, Lorenzo Porzi, Samuel Rota Bulo, Marc Pollefeys, and Peter Kontschieder. 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.

[8] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.

[9] Shubham Goel, Georgios Pavlakos, Jathushan Rajasegaran, Angjoo Kanazawa, and Jitendra Malik. Humans in 4D: Reconstructing and tracking humans with transformers. In *ICCV*, 2023.

[10] 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.

[11] Marc Habermann, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Livecap: Real-time human performance capture from monocular video. *ACM Transactions On Graphics (TOG)*, 38(2):1–17, 2019.

[12] Marc Habermann, Weipeng Xu, Michael Zollhofer, Gerard Pons-Moll, and Christian Theobalt. Deepcap: Monocular human performance capture using weak supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5052–5063, 2020.

[13] Shoukang Hu and Ziwei Liu. Gauhuman: Articulated gaussian splatting from monocular human videos. *arXiv preprint arXiv:*, 2023.

[14] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. Sc-gs: Sparse-controlled gaussian splatting for editable dynamic scenes. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.

[15] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. Neuman: Neural human radiance field from a single video. In *European Conference on Computer Vision*, pages 402–418. Springer, 2022.

[16] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5614–5623, 2019.

[17] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023.

[18] Muhammed Kocabas, Jen-Hao Rick Chang, James Gabriel, Oncel Tuzel, and Anurag Ranjan. HUGS: Human gaussian splatting. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

[19] 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.

[20] Jiahui Lei, Yufu Wang, Georgios Pavlakos, Lingjie Liu, and Kostas Daniilidis. Gart: Gaussian articulated template models. *arXiv preprint arXiv:2311.16099*, 2023.

[21] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, and Zhaoyang Lv. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.

[22] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. 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)*, 2021.

[23] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. volume 34, pages 248:1–248:16. ACM, October 2015.

[24] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In *3DV*, 2024.

[25] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021.

[26] 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. In *ECCV*, 2020.

[27] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–15, 2022.

[28] 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.

[29] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021.

[30] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *ACM Trans. Graph.*, 40(6), dec 2021.

[31] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.

[32] 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.

[33] 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.

[34] Qingping Sun, Yanjun Wang, Ailing Zeng, Wanqi Yin, Chen Wei, Wenjia Wang, Haiyi Mei, Chi-Sing Leung, Ziwei Liu, Lei Yang, et al. Aios: All-in-one-stage expressive human pose and shape estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1834–1843, 2024.

[35] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022.

[36] 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.

[37] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2021.

[38] Haithem Turki, Jason Y Zhang, Francesco Ferroni, and Deva Ramanan. Suds: Scalable urban dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12375–12385, 2023.

[39] Jingbo Wang, Zhengyi Luo, Ye Yuan, Yixuan Li, and Bo Dai. Pacer+: On-demand pedestrian animation controller in driving scenarios. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

[40] Jingbo Wang, Ye Yuan, Zhengyi Luo, Kevin Xie, Dahua Lin, Umar Iqbal, Sanja Fidler, and Sameh Khamis. Learning human dynamics in autonomous driving scenarios. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20796–20806, 2023.

[41] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Wang Xinggang. 4d gaussian splatting for real-time dynamic scene rendering. *arXiv preprint arXiv:2310.08528*, 2023.

[42] Tianhao Wu, Fangcheng Zhong, Andrea Tagliasacchi, Forrester Cole, and Cengiz Oztireli. $D^2$nerf: Self-supervised decoupling of dynamic and static objects from a monocular video. *Advances in Neural Information Processing Systems*, 2022.

[43] Zirui Wu, Tianyu Liu, Liyi Luo, Zhide Zhong, Jianteng Chen, Hongmin Xiao, Chao Hou, Haozhe Lou, Yuantao Chen, Runyi Yang, Yuxin Huang, Xiaoyu Ye, Zike Yan, Yongliang Shi, Yiyi Liao, and Hao Zhao. Mars: An instance-aware, modular and realistic simulator for autonomous driving. *CICAI*, 2023.

[44] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9421–9431, 2021.

[45] 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.

[46] 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.

[47] Jiawei Yang, Boris Ivanovic, Or Litany, Xinshuo Weng, Seung Wook Kim, Boyi Li, Tong Che, Danfei Xu, Sanja Fidler, Marco Pavone, and Yue Wang. Emernerf: Emergent spatial-temporal scene decomposition via self-supervision. In *International Conference on Learning Representations*, 2023.

[48] Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. Unisim: A neural closed-loop sensor simulator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1389–1399, 2023.

[49] Ze Yang, Sivabalan Manivasagam, Ming Liang, Bin Yang, Wei-Chiu Ma, and Raquel Urtasun. Recovering and simulating pedestrians in the wild. In *Conference on Robot Learning*, pages 419–431. PMLR, 2021.

[50] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. *arXiv preprint arXiv:2309.13101*, 2023.

[51] Zongxin Ye, Wenyu Li, Sidun Liu, Peng Qiao, and Yong Dou. Absgs: Recovering fine details for 3d gaussian splatting. *arXiv preprint arXiv:2404.10484*, 2024.

[52] Zhengming Yu, Wei Cheng, Xian Liu, Wayne Wu, and Kwan-Yee Lin. Monohuman: Animatable human neural field from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16943–16953, 2023.

[53] Ye Yuan, Umar Iqbal, Pavlo Molchanov, Kris Kitani, and Jan Kautz. Glamr: Global occlusion-aware human mesh recovery with dynamic cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[54] 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. *arXiv preprint arXiv:2403.12722*, 2024.

[55] Xiaoyu Zhou, Zhiwei Lin, Xiaojun Shan, Yongtao Wang, Deqing Sun, and Ming-Hsuan Yang. Drivinggaussian: Composite gaussian splatting for surrounding dynamic autonomous driving scenes. *arXiv preprint arXiv:2312.07920*, 2023.

# Supplemental Material

## A  Implementation Details

**Initialization:** For the background model, we refer to PVG [6], combining $6 \times 10^5$ LiDAR points with $4 \times 10^5$ random samples, which are divided into $2 \times 10^5$ near samples uniformly distributed by distance to the scene's origin and $2 \times 10^5$ far samples uniformly distributed by inverse distance. To initialize the background, we filter out the LiDAR samples of dynamic objects. For rigid nodes and non-rigid deformable nodes, we utilize their bounding boxes to accumulate the LiDAR points, while for non-rigid SMPL nodes, we initialize the Gaussians on the template mesh in their canonical space. To determine the initial color of Gaussians, we project LiDAR points onto the image plane, whereas random samples are initialized with random colors. The initial human body pose sequences of non-rigid SMPL Nodes are obtained through the process described in § 4.2.

**Training:** Our method trains for 30,000 iterations with all scene nodes optimized jointly. The learning rate for Gaussian properties aligns with the default settings of 3DGS [17], but varies slightly across different node types. Specifically, we set the learning rate for the rotation of Gaussians to $5 \times 10^{-5}$ for non-rigid SMPL nodes and $1 \times 10^{-5}$ for other nodes. The degrees of spherical harmonics are set to 3 for background nodes, rigid nodes, and non-rigid deformable nodes, while it is set to 1 for non-rigid SMPL nodes. The learning rate for the rotation of instance boxes is $1 \times 10^{-5}$, decreasing exponentially to $5 \times 10^{-6}$. The learning rate for the translation of instance boxes is $5 \times 10^{-4}$, decreasing exponentially to $1 \times 10^{-4}$. The learning rate for human body poses of non-rigid SMPL nodes is $5 \times 10^{-5}$, decreasing exponentially to $1 \times 10^{-5}$. For the Gaussian densification strategy, we utilize the absolute gradient of Gaussians introduced in [51] to control memory usage. We set the densification threshold of position gradient to $3 \times 10^{-4}$. This use of absolute gradient has a minimal impact on performance, as discussed in detail in Appendix D.3. The densification threshold for scaling is $3 \times 10^{-3}$. Our method runs on a single NVIDIA RTX 4090 GPU, with training for each scene taking about 1 hour. Training time varies with different training settings. For more implementation details, please visit our project page for code.

**Optimization:** We utilize the loss function introduced in Eq (8) to jointly optimize all learnable parameters. The image loss is computed as:

$$\mathcal{L}_{\text{image}} = (1 - \lambda_r)\,\mathcal{L}_1 + \lambda_r \mathcal{L}_{\text{SSIM}} \tag{9}$$

due to sparse temporal-spatial observation of the dynamic part, its supervision signal is insufficient. To address this, we apply a higher image loss weight to the dynamic regions identified by the rendered dynamic mask. This weight is set to 5. The depth map loss is computed as:

$$\mathcal{L}_{\text{depth}} = \frac{1}{hw} \sum \left\| \mathcal{D}^s - \hat{\mathcal{D}} \right\|_1 \tag{10}$$

where $\mathcal{D}^s$ is the inverse of the sparse depth map. We project LiDAR points onto the image plane to generate the sparse LiDAR map, and $\hat{\mathcal{D}}$ is the inverse of the predicted depth map.

The mask loss $\mathcal{L}_{\text{opacity}}$ is computed as:

$$\mathcal{L}_{\text{opacity}} = -\frac{1}{hw} \sum O_{\mathcal{G}} \cdot \log O_{\mathcal{G}} - \frac{1}{hw} \sum M_{\text{sky}} \cdot \log(1 - O_{\mathcal{G}}) \tag{11}$$

where $M_{\text{sky}}$ is the sky mask, and $O_{\mathcal{G}}$ is the rendered opacity map.

In addition to the reconstruction losses, we introduce various regularization terms for different Gaussian representations to improve quality. Among these, an important regularization term is $\mathcal{L}_{\text{pose}}$, designed to ensure smooth human body poses $\boldsymbol{\theta}(t)$. This term is defined as:

$$\mathcal{L}_{\text{pose}} = \frac{1}{2} \left\| \boldsymbol{\theta}(t - \delta) + \boldsymbol{\theta}(t + \delta) - 2\boldsymbol{\theta}(t) \right\|_1 \tag{12}$$

where $\delta$ is a randomly chosen integer from $\{1, 2, 3, 4, 5\}$. We set the weight of the SSIM loss, $\lambda_r$, to 0.2, the depth loss, $\lambda_{\text{depth}}$, to 0.1, the opacity loss, $\lambda_{\text{opacity}}$, to 0.05, and the pose smoothness loss, $\lambda_{\text{pose}}$, to 0.01. More details can be found in our open-sourced codebase.

## B    Baselines

• **EmerNeRF** [47] is a state-of-the-art NeRF-Based method for dynamic driving scene reconstruction. EmerNeRF uses a static field represented by a 3D Hash-Grid to model the static parts of the scene and a dynamic field with a 4D Hash-Grid to model the dynamic parts. Additionally, it employs a flow field to aggregate the dynamic features. This self-supervised decomposition approach yields good results on dynamic scene modeling and static-dynamic decomposition, with the scene flow emerging in the flow field.

• **DeformableGS** [50] defines a canonical space to represent scenes with Gaussians. To model dynamics, it uses a deformation network to predict offsets of Gaussian properties. These offsets then deform the Gaussians to fit the scene dynamics. DeformableGS works well in synthetic and indoor datasets. We compare it to our method to evaluate its ability to model challenging out-door dynamic scenes.

• **StreetGS** [46] is a dynamic scene modeling method based on Gaussian Splatting for driving scenes. StreetGS models the components of dynamic scenes separately: the static background and the foreground vehicles. It utilizes boxes predicted by an off-the-shelf model to warp the Gaussians of foreground vehicles and refine them during training. StreetGS yields good results on driving scenes but ignores other non-rigid dynamic objects in the scene.

• **HUGS** [54] is a GS-based method for driving scene modeling and understanding. It not only models the appearance of a scene but also distills 2D flow maps and semantic maps into the 3D scene to enable holistic urban scene understanding. Similar to StreetGaussian [46], HUGS uses object boxes for compositing dynamic elements. HUGS achieves good performance in both scene modeling and semantic modeling. However, it primarily focuses on rigid backgrounds and objects, without addressing non-rigid dynamics.

• **PVG** [6] introduces Periodic Vibration Gaussians that vibrate over time with optimizable vibration directions, life span, and life peak (the moment of highest opacity) to represent dynamic scenes. These Gaussians are optimized using a self-supervised approach. The method achieves static-dynamic decomposition by categorizing Gaussians based on their life spans. We compare PVG with our method to evaluate our capability in modeling highly complex dynamic scenes.

Among all the compared methods, HUGS [54] and StreetGaussians [46] require bounding boxes of foreground objects. PVG [6], StreetGaussians [46], and EmerNeRF [47] utilize LiDAR data for depth supervision. While the original implementations of 3DGS [17] and DeformableGS [50] do not include depth supervision, we added LiDAR depth supervision in experiments to ensure a fair comparison.

## C    Evaluation

**Appearance.** For the Novel View Synthesis task, we select every 10th frame from the original sequence as the test set. We use PSNR and SSIM to evaluate the quality of the rendered images. Since we focus on dynamic scenes, we also compute the PSNR and SSIM for regions with vehicles and humans. To identify regions of vehicles and humans, we use Segformer [45] to obtain semantic masks. We further identify the movable dynamic parts using projections of moving object bounding boxes, utilizing their velocity information. One example of dynamic masks cam bee seen in Fig. 9.



GT Image                    Human Region                    Vehicle Region

Figure 9: An example of the dynamic masks for computing dynamic region metrics.

**Geometry.** Our method uses LiDAR data to initialize Gaussians and supervise scene depth by comparing the rendered depth map with the sparse LiDAR depth map. Post-training, Gaussians typically deviate from their initial state through densification or optimization, Therefore, comparing the LiDAR depth reconstruction is still a valid comparison. We follow the depth evaluation method of StreetSurf [10]: render a depth map and match depth pixels to LiDAR rays. For Chamfer Distance, re-project the predicted depth to 3D using the LiDAR ray direction and origin. For RMSE, compare the GT and predicted ranges for LiDAR rays.

# D Additional Results

## D.1 Qualitative Comparison

We refer readers to our our project page for video comparisons of the methods.



(g) DeformableGS                                          (h) 3DGS
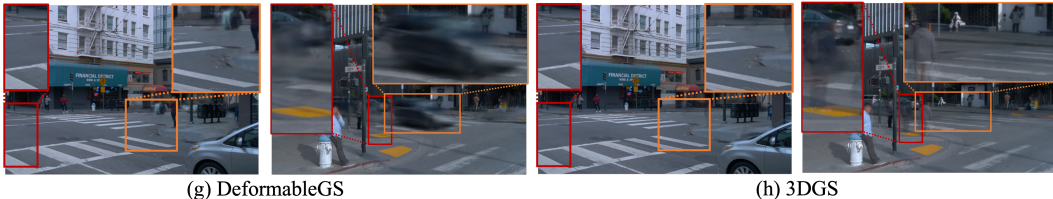
Figure 10: Additional Qualitative Comparison of Novel View Synthesis.

## D.2 Quantitative Comparison

To further validate our method's effectiveness, we tested our method against StreetGS [46] and EmerNeRF [47] on 32 dynamic scenes from the Waymo dataset, with results reported in Tab. 5.

Table 5: We expanded our evaluation to 32 dynamic scenes from the Waymo dataset, comparing our method with StreetGS [46] and EmerNeRF [47]. The segment IDs are listed in Tab. 6.

| | Scene Reconstruction | | | | | | Novel View Synthesis | | | | | |
| | Full Image | | Human | | Vehicle | | Full Image | | Human | | Vehicle | |
| Methods | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| EmerNeRF | 31.29 | 0.877 | 23.14 | 0.581 | 24.47 | 0.709 | 29.04 | 0.851 | 20.76 | 0.467 | 21.80 | 0.582 |
| StreetGS | 29.93 | 0.931 | 19.63 | 0.524 | 27.48 | 0.871 | 28.73 | 0.910 | 18.77 | 0.470 | 26.18 | 0.825 |
| Ours | **33.73** | **0.946** | **28.28** | **0.855** | **28.02** | **0.880** | **31.71** | **0.924** | **24.57** | **0.730** | **26.55** | **0.833** |

Table 6: Segment IDs of 32 dynamic scenes of Waymo Dataset used in the test for Tab. 5.

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| seg104554... | seg125050... | seg169514... | seg584622... | seg776165... | seg138251... | seg448767... | seg965324... |
| seg119252... | seg122514... | seg132544... | seg134024... | seg166004... | seg173881... | seg215148... | seg391164... |
| seg454855... | seg560223... | seg571325... | seg587066... | seg842457... | seg952165... | seg952995... | seg112365... |
| seg152664... | seg411445... | seg123218... | seg102252... | seg148106... | seg265611... | seg179934... | seg104859... |

## D.3 Ablation Studies

**Absolute Gradient.** In our implementation, we applied AbsGrad for 3DGS densification across all reproduced methods (StreetGS, DeformableGS, and 3DGS) as a standard practice. To quantify the impact of AbsGrad, we conducted a comparative study with and without its application. The results of this analysis are presented in Tab. 7. We see that disabling AbsGrad leads to a marginal performance decrease ( 0.1 PSNR) for all methods, proving that AbsGrad is not the key factor contributing to our performance advantage over others. Note that DeformableGS fails to run due to out-of-memory issues when AbsGrad was disabled. Based on these findings, we recommend the incorporation of AbsGrad as a standard practice in 3DGS densification and related methods.

Table 7: **Ablation on AbsGrad.** By default, we apply AbsGrad to all GS-based approaches reproduced by us. We now disable it to analyze its impact. We mark methods with AbsGrad enabled with grey background. We observe that 1) DeformableGS fails under w/o. AbsGrad setting because of out of memory issue; 2) enabling AbsGrad is a good practice ( +0.1 PNSR for all methods) but not an enabling factor for our performance lead.

| | Scene Reconstruction | | | | | | Novel View Synthesis | | | | | |
| | Full Image | | Human | | Vehicle | | Full Image | | Human | | Vehicle | |
| Methods | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| EmerNeRF | 31.93 | 0.902 | 22.88 | 0.578 | 24.65 | 0.723 | 29.67 | 0.883 | 20.32 | 0.454 | 22.07 | 0.609 |
| 3DGS* | 26.00 | 0.912 | 16.88 | 0.414 | 16.18 | 0.425 | 25.57 | 0.906 | 16.62 | 0.387 | 16.00 | 0.407 |
| 3DGS | 25.84 | 0.910 | 16.69 | 0.405 | 16.02 | 0.415 | 25.61 | 0.905 | 16.52 | 0.383 | 15.97 | 0.405 |
| DeformGS* | 28.40 | 0.929 | 17.80 | 0.460 | 19.53 | 0.570 | 27.72 | 0.922 | 17.30 | 0.426 | 18.91 | 0.530 |
| DeformGS | – | – | – | – | – | – | – | – | – | – | – | – |
| PVG | 32.37 | 0.937 | 24.06 | 0.703 | 25.02 | 0.787 | 30.19 | 0.919 | 21.30 | 0.567 | 22.28 | 0.679 |
| HUGS | 28.26 | 0.923 | 16.23 | 0.404 | 24.31 | 0.794 | 27.65 | 0.914 | 15.99 | 0.378 | 23.27 | 0.748 |
| StreetGS* | 29.08 | 0.936 | 16.83 | 0.420 | 27.73 | 0.880 | 28.54 | 0.928 | 16.55 | 0.393 | 26.71 | 0.846 |
| StreetGS | 28.89 | 0.932 | 16.70 | 0.409 | 28.07 | 0.878 | 28.46 | 0.926 | 16.41 | 0.387 | 26.86 | 0.845 |
| Ours* | 34.25 | 0.954 | 28.15 | 0.845 | 28.91 | 0.892 | 32.57 | 0.942 | 24.36 | 0.727 | 27.57 | 0.858 |
| Ours | 34.11 | 0.953 | 28.00 | 0.842 | 28.83 | 0.890 | 32.46 | 0.941 | 24.28 | 0.726 | 27.55 | 0.857 |

Table 8: Segment IDs of 8 dynamic scenes of Waymo Dataset used in the test for Tab. 1. and Tab. 7

| seg104554... | seg125050... | seg169514... | seg584622... seg776165... | seg138251... | seg448767... | seg965324... |

**Additional Results** The Tab. 10 is the full table of Tab. 3 that includes evaluation on SSIM. The Tab. 9 is the full table of Tab. 2 that includes evaluation on SSIM.

Table 9: **Ablation on GT Boxes Refinement.**

| | Scene Reconstruction | | | | | | Novel View Synthesis | | | | | |
| | Full Image | | Human | | Vehicle | | Full Image | | Human | | Vehicle | |
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Complete Model | 34.25 | 0.954 | 28.15 | 0.845 | 28.91 | 0.892 | 32.57 | 0.942 | 24.36 | 0.727 | 27.57 | 0.858 |
| w/o Box Refine. | 33.04 | 0.947 | 26.53 | 0.790 | 25.57 | 0.813 | 31.72 | 0.936 | 23.67 | 0.686 | 24.78 | 0.785 |

# E `OmniRe` In Practice

**Bounding Boxes.** Similar to other scene-graph-based approaches [28, 48, 36, 7, 55, 54, 46], we utilize bounding boxes for driving scene reconstruction, as they are widely used for producing superior reconstruction results compared to methods that do not employ them. Additionally, bounding boxes offer significant controllability. It allows precise manipulation of both rigid objects like vehicles and non-rigid objects such as individual human body movements—an ability lacking in self-supervised methods like EmerNeRF [47] and PVG [6] that do not use instance information. This level of controllability is crucial for tasks like scene simulation, which require the ability to manage the movement of all participating agents. Lastly, bounding box annotation is a standard and generally straightforward process in the autonomous driving field, with most popular datasets already providing these annotations via established auto-labeling tools, thereby minimizing manual effort and making the resource both efficient and accessible. For real-world driving logs, these auto-labeling tools can generate precise bounding boxes at little cost.

**How to Determine Gaussian Representations for Humans?** We categorize pedestrians into two groups for modeling. Near-range pedestrians, detected by our human pose processing module introduced in § 4.2, are modeled using SMPL nodes. Far-range pedestrians, typically undetected due to distance, are modeled using deformable nodes. This approach naturally distinguishes between near and far-range pedestrians based on human detection capabilities.

Other individuals, such as those using wheelchairs, skateboards, or bicycles, are often labeled as "cyclists" in the datasets we study. However, these labels may be specific to the dataset used, and in some cases, the annotations might not be accurate. For instance, in the Waymo Dataset, a person on a motorcycle may be labeled as a "vehicle". This reliance on dataset-specific labels could potentially limit the generalization of our method to other scenarios with imperfect labels.

To address this issue, we conducted preliminary experiments using GPT-4o [1] to classify individuals (cropped by bounding boxes) into two categories: pedestrians and humans using personal transportation devices (e.g., wheelchairs, bicycles, motorcycles). Testing on 60 individuals (30 from each

Table 10: **Ablation on Non-Rigid Modeling.**

| | Scene Reconstruction | | | | Novel View Synthesis | | | |
| | Full Image | | Human | | Full Image | | Human | |
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
|---|---|---|---|---|---|---|---|---|
| (a) Ours default | **34.25** | **0.954** | **28.15** | **0.845** | **32.57** | **0.942** | **24.36** | **0.727** |
| (b) w/o SMPL actors | 32.80 | 0.949 | 24.71 | 0.770 | 31.76 | 0.939 | 23.18 | 0.694 |
| (c) w/o Body Pose Refine. | 33.84 | 0.952 | 26.97 | 0.815 | 32.44 | 0.941 | 24.04 | 0.712 |
| (d) w/o Deformed actors | 33.64 | 0.953 | 25.26 | 0.766 | 32.17 | 0.941 | 22.41 | 0.653 |



Reconstructed Background     Original Car     Insert Car A     Insert Car B
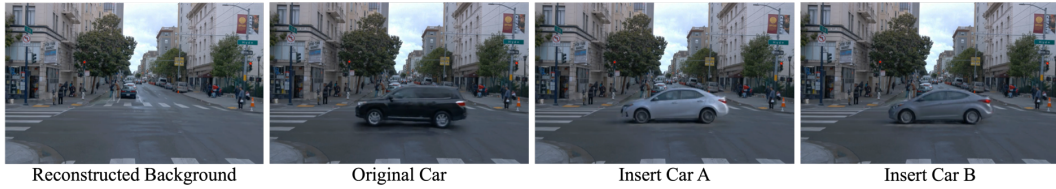
Figure 11: Our method allows for flexible editing of scene assets.

category), GPT-4o [1] achieved 100% accuracy. This suggests that accurate labels can be obtained relatively easily, thanks to the development of vision-language models.
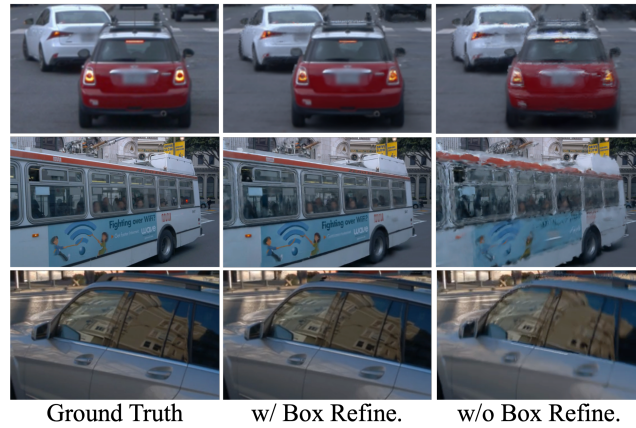


Ground Truth     w/ Box Refine.     w/o Box Refine.

Figure 12: **Ablation of Boxes Refinement.**