

3DREALCAR: AN IN-THE-WILD RGB-D CAR DATASET WITH 360-DEGREE VIEWS

Anonymous authors

Paper under double-blind review

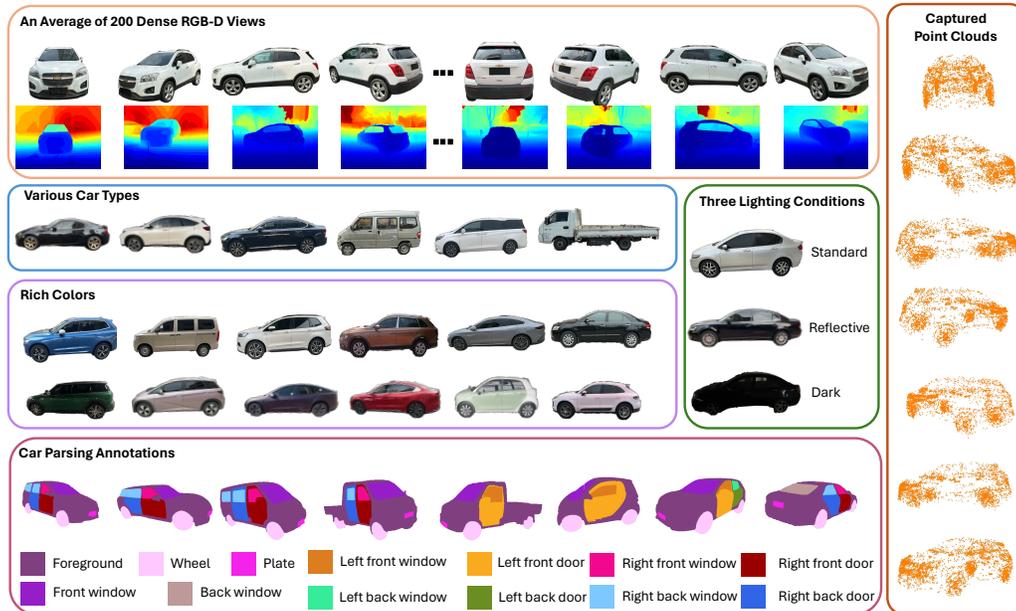


Figure 1: **Characteristics of our curated high-quality 3DRealCar dataset.** 3DRealCar contains detailed annotations for various colors, car types, brands, and even car parsing maps. In particular, our dataset contains three lighting conditions on car surfaces, bringing challenges to existing methods.

ABSTRACT

3D cars are commonly used in self-driving systems, virtual/augmented reality, and games. However, existing 3D car datasets are either synthetic or low-quality, presenting a significant gap toward the high-quality real-world 3D car datasets and limiting their applications in practical scenarios. In this paper, we propose the first large-scale 3D real car dataset, termed 3DRealCar, offering three distinctive features. (1) **High-Volume**: 2,500 cars are meticulously scanned by smartphones, obtaining car images and point clouds with real-world dimensions; (2) **High-Quality**: Each car is captured in an average of 200 dense, high-resolution 360-degree RGB-D views, enabling high-fidelity 3D reconstruction; (3) **High-Diversity**: The dataset contains various cars from over 100 brands, collected under three distinct lighting conditions, including reflective, standard, and dark. Additionally, we offer detailed car parsing maps for each instance to promote research in car parsing tasks. Moreover, we remove background point clouds and standardize the car orientation to a unified axis for the reconstruction only on cars and controllable rendering without background. We benchmark 3D reconstruction results with state-of-the-art methods across each lighting condition in 3DRealCar. Extensive experiments demonstrate that the standard lighting condition part of 3DRealCar can be used to produce a large number of high-quality 3D cars, improving various 2D and 3D tasks related to cars. Notably, our dataset brings insight into the fact that recent 3D reconstruction methods face challenges in reconstructing high-quality 3D cars under reflective and dark lighting conditions. **Our dataset is available here.**

Table 1: **The comparison of existing 3D car datasets.** Our dataset contains unique characteristics compared with existing 3D car datasets. Lighting means the lighting conditions of the surfaces of cars. Point Cloud represents the point clouds with actual sizes in real-world scenes.

Dataset	Instances	Type	Views	Resolution	Brand	Lighting	Car Parsing	Depth	Point Cloud
SRN-Car	2151	Synthetic	250	128×128	×	×	×	×	×
Objaverse-car	511	Synthetic	-	-	×	×	×	×	×
MVMC	576	Real	~10	600×450	~40	×	×	×	×
3DRealCar (Ours)	2500	Real	~200	1920×1440	100+	3	13	✓	✓



Figure 2: **Visual comparisons of 3D car datasets and the results of a 3D generative method.** Our 3DRealCar is captured in real-world scenes and contains more densely captured views. In addition, our dataset has annotations for three different lighting conditions on the car surface. We also compare a recent state-of-the-art text-to-3D model, MVDream (Shi et al., 2023b) with a prompt “*a modern sedan*”, demonstrating its failure to generate high-quality 3D car models.

1 INTRODUCTION

Cars, as both daily objects and vehicles, are of significant interest to researchers, especially in the field of autonomous driving. Autonomous perception systems are typically trained on daily scene datasets that are collected frequently. However, these datasets often exhibit long-tailed distributions, with far fewer instances of corner-case scenarios, like car accidents. Consequently, this imbalance leads to the autonomous perception system generalizing well in the most frequently occurring scenes. This means that the system is likely to perform poorly in rare situations, posing significant safety risks to drivers. To build a reliable system, it is essential to have a simulator that can simulate photorealistic hazardous scenes. Moreover, high-quality 3D cars are necessary for a realistic simulator.

Recent 3D car reconstruction methods (Wang et al., 2023a; Zhou et al., 2023; Xie et al., 2023) mainly reconstruct cars from self-driving datasets (Sun et al., 2020; Caesar et al., 2020; Geiger et al., 2013). To apply reconstructed cars to real-world scenes, the reconstructed 3D car should be high-quality. However, it is very challenging to obtain such high-quality 3D cars for the following reasons: (1) Previous 3D car reconstruction methods produce low-quality 3D cars, primarily because they train on self-driving datasets with low-resolution car images and a limited number of trainable views. (2) Manually crafting a high-quality 3D car model requires specialized artists, which is time-consuming. (3) There is no large-scale 3D real car dataset that can be utilized to produce a bulk of 3D cars.

Moreover, existing 3D car datasets are either synthetic or only contain a few posed images, as shown in Figure 2. SRN-Car (Chang et al., 2015) and Objaverse-Car (Deitke et al., 2023) collect 3D car

108 computer-aided design (CAD) models from the Internet, but these models are synthetic and contain
109 non-photorealistic texture. Although MVMC (Zhang et al., 2021) is a real car dataset, it collects only
110 ten views on average for each car. On the contrary, our collected 3DRealCar dataset provides an
111 average of 200 dense RGB-D views per car for high-quality 3D car reconstruction.

112 We also show that the recent state-of-the-art 3D generative method, MVDream (Shi et al., 2023b), as
113 depicted in Figure 2, fails to generate high-quality cars due to the multi-view inconsistency introduced
114 by generative models (Rombach et al., 2022; Stability.AI, 2023; Liu et al., 2023c; Sun et al., 2023).
115 Thus, the existing 3D generation methods cannot be employed to generate high-quality 3D real car
116 assets.

117 In this work, we collect a large-scale 3D real car dataset in the wild, termed 3DRealCar, which
118 contains dense high-quality views and rich diversity. **During data collection, we employ smartphones
119 with ARKit (Apple, 2021) to scan cars parked on roadsides or parking lots, obtaining posed RGB-D
120 images and point clouds of cars.** In particular, we scan around the cars in three loops to obtain dense
121 views. Note that we collect car data with the consent of owners. In Table 1 and Figure 1, we show
122 our dataset possesses striking characteristics compared with previous 3D car datasets. We capture
123 dense RGB-D images in high resolution, which promotes the reconstruction of high-quality 3D cars.
124 Furthermore, we scan cars under three different lighting conditions, resulting in the surfaces of cars
125 having different lighting effects, such as reflective, standard, and dark, where we denote the standard
126 as the smooth lighting condition without obvious specular highlight. Figure 2 shows some examples
127 of three lighting conditions in our dataset. Note that the number of instances in our dataset is the
128 largest in existing datasets. Therefore, our collected 3DRealCar dataset has a rich diversity in terms
129 of car types, colors, brands, and lighting conditions. We also provide car parsing map annotations
130 with thirteen classes for each instance, which enable our dataset to be applied in car component
131 understanding tasks.

132 To construct a high-quality dataset, we filter out the images that are out of focus, occluded, or
133 blurred. To facilitate the 3D reconstruction solely on cars, we remove the point clouds of the
134 background. We also adjust the orientation of the car facing along the x-axis before the reconstruction
135 for controllable rendering. Based on the high-quality posed RGB-D images, point clouds, and multi-
136 grained annotations, we can apply the dataset to various tasks related to cars. Figure 3 shows our
137 dataset supports over 10 tasks, including several popular 2D and 3D tasks to promote the advancement
138 of car-related research.

139 We leverage existing state-of-the-art methods to benchmark 3D reconstruction and car parsing tasks
140 of our 3DRealCar dataset. We also conduct extensive experiments to demonstrate that the reflective
141 and dark lighting conditions in our dataset are challenging to existing methods, which brings a new
142 challenge for 3D reconstruction in awful lighting conditions. Furthermore, we demonstrate that
143 our 3DRealCar dataset can bring real-car prior and enhance existing 3D generation and downstream
144 methods. Overall, the contributions of this work can be summarized below:

- 145 • We propose the first large-scale 3D real car dataset, named 3DRealCar, which contains 2,500 car
146 instances and their point clouds with actual sizes in real-world scenes.
- 147 • 3DRealCar contains RGB-D images and point clouds with detailed annotations, supporting re-
148 searchers to investigate various tasks in both 2D and 3D.
- 149 • We conduct 3D reconstruction and car parsing benchmarks to advance car-related tasks. Notably,
150 we observe that existing methods face challenges under the extreme lighting conditions of 3DRealCar.
- 151 • Extensive experiments demonstrate our 3DRealCar dataset can enhance real-car prior and improve
152 the performance of existing 3D generation and novel view synthesis methods.

153 2 RELATED WORK

154
155 **3D Car Datasets.** There are several well-known large-scale autonomous driving datasets so far,
156 such as Nuscenes (Caesar et al., 2020), KITTI (Geiger et al., 2013), Waymo (Sun et al., 2020),
157 Pandaset (Xiao et al., 2021), ApolloScape (Huang et al., 2018), and Cityscape (Cordts et al., 2016).
158 These datasets are captured by multi-view cameras and lidars mounted on ego cars. Various works
159 (Wang et al., 2023a; García Orellana et al., 2001; Liu et al., 2024; Xie et al., 2023) attempt to
160 reconstruct 3D cars in these datasets. However, these methods fall short of reconstructing high-
161

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

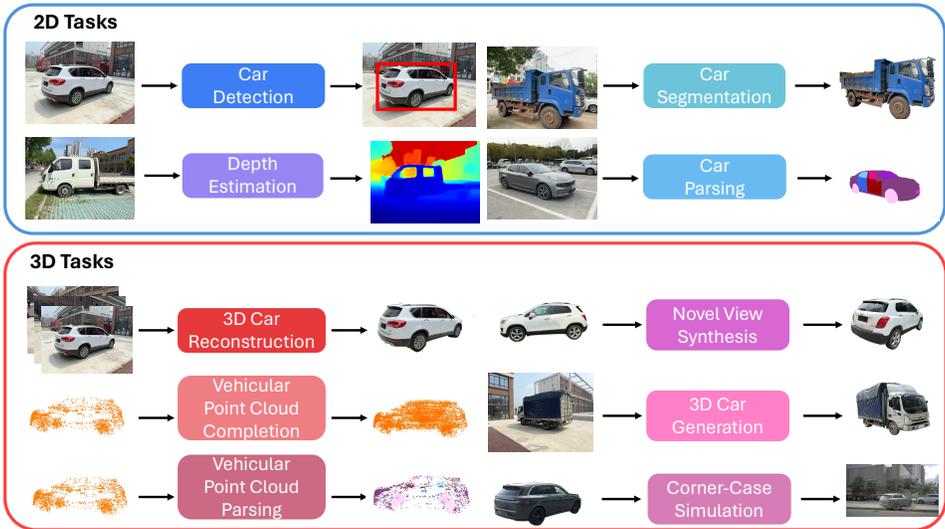


Figure 3: **The applicable tasks of our dataset.** Our proposed 3DRealCar dataset containing RGB-D images, point clouds, and rich annotations can be applied to various popular 2D and 3D tasks to support the construction of safe and reliable self-driving system.

quality 3D cars due to the lack of sufficient and dense training views. DeepMANTA (Chabot et al., 2017) provides car component segmentation maps, but this dataset is based on the synthetic CAD model that cannot be accurately used in real-world settings. SRN-Car (Chang et al., 2015) and Objaverse (Deitke et al., 2023) collect 3D car models from existing repositories and Internet sources. However, these datasets only contain synthetic cars, which cannot produce realistic textures and geometry. MVMC (Zhang et al., 2021) is collected from car advertising websites, which contain a series of car images, especially multi-view images of each car. However, the views of images per car in MVMC are unposed and sparse, which is adverse to reconstructing high-quality 3D car models. In this paper, we collect a high-quality 3D real car dataset to fill the above gaps.

3D Reconstruction with Neural Field. 3D reconstruction aims to create a 3D structure digital representation of an object or a scene from its multi-view images, which is a long-standing task in computer vision. One of the most representative works in 3D reconstruction is Neural Radiance Fields (NeRFs) (Mildenhall et al., 2021), which demonstrates promising performance for novel view synthesis. Afterward, this method inspires a new wave of 3D reconstruction methods using the volume rendering method, with subsequent works focusing on improving its quality (Verbin et al., 2021; Barron et al., 2021; 2022b; Guo et al., 2022; Suhail et al., 2022; Chen et al., 2022b; Wang et al., 2023b; Barron et al., 2023), efficiency (Fridovich-Keil et al., 2022; Müller et al., 2022; Reiser et al., 2021; Sun et al., 2022; Kerbl et al., 2023a; Chen et al., 2022a; Garbin et al., 2021), applying artistic effects (Fan et al., 2022; Wang et al., 2022; Jain et al., 2022; Zhang et al., 2022b), and generalizing to unseen scenes (Yu et al., 2021; Chen et al., 2021; Wang et al., 2021; Johari et al., 2022; T et al., 2023; Chibane et al., 2021). Particularly, Kilonerf (Reiser et al., 2021) accelerates the training process of NeRF by dividing a large MLP into thousands of tiny MLPs. Furthermore, Mip-NeRF (Barron et al., 2021) proposes a conical frustum rather than a single ray to ameliorate aliasing. Mip-NeRF 360 (Barron et al., 2022a) further improves the application scenes of NeRF to the unbounded scenes. Although these NeRF-based methods demonstrate powerful performance on various datasets, the training time always requires several hours even one day more. Instant-NGP (Müller et al., 2022) uses a multi-resolution hash encoding method, which reduces the training time by a large margin. 3DGS (Kerbl et al., 2023a) proposes a new representation based on 3D Gaussian Splatting, which reaches real-time rendering for objects or unbounded scenes. 2DGS (Huang et al., 2024) proposes a perspective-accurate 2D splatting process that leverages ray-splat intersection and rasterization to further enhance the quality of the reconstructions. Scaffold-GS (Lu et al., 2023) proposes an anchor growing and pruning strategy to accelerate the scene coverage, which effectively reduces redundant Gaussians and improves rendering quality. However, there is not yet a large-scale 3D real car dataset so far. Therefore, we present a 3D real car dataset, named 3DRealCar in this work.

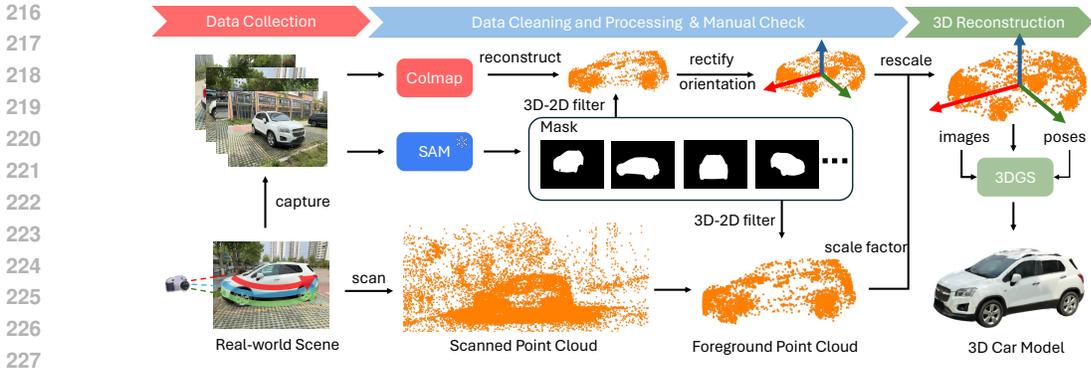


Figure 4: **Illustration of our data collection and preprocessing.** We first circle a car three times while scanning the car with a **smartphone** for the attainment of RGB-D images and its point clouds. Then we use Colmap (Schonberger & Frahm, 2016) and SAM (Kirillov et al., 2023) to obtain poses and remove the background point clouds. Finally, we use the 3DGS (Kerbl et al., 2023b) trained on the processed data to obtain 3D car model.

3D Generation with Diffusion Prior. Some current works (Jun & Nichol, 2023; Nichol et al., 2022) leverage a 3D diffusion model to learn the representation of 3D structure. However, these methods lack generalization ability due to the scarcity of 3D data. To facilitate 3D generation without direct supervision of 3D data, image or multi-view diffusion models are often used to guide the 3D creation process. Notable approaches like DreamFusion (Poole et al., 2022b) and subsequent works (Metzer et al., 2023; Lin et al., 2023) use an existing image diffusion model as a scoring function, applying Score Distillation Sampling SDS loss to generate 3D objects from textual descriptions. These methods, however, suffer from issues such as the Janus problem (Poole et al., 2022b; Metzer et al., 2023) and overly saturated textures. Inspired by Zero123 (Liu et al., 2023c), several recent works (Stability.AI, 2023; Shi et al., 2023a; Liu et al., 2023e; Kong et al., 2024; Zheng & Vedaldi, 2023; Melas-Kyriazi et al., 2024; Liu et al., 2023b;a; Qian et al., 2023) refine image or video diffusion models to better guide the 3D generation by producing more reliable multi-view images. However, these generative methods fail to generate high-quality cars due to the lack of the prior of real cars.

3 PROPOSED 3DREALCAR DATASET

3.1 DATA COLLECTION AND ANNOTATION

As shown in Figure 4, our dataset is collected using smartphones, specifically iPhone 14 models, adopting ARKit APIs (Apple, 2021) to scan cars for their point clouds and RGB-D images. The data collection process is conducted under three distinct lighting conditions, such as standard, reflective, and dark. These lighting conditions represent the lighting states of vehicle surfaces. It is important to note that all data collection is performed with the consent of owners. During the scanning process, the car should be stationary while we meticulously circle the car three times to capture as many views as possible. For each loop, we adjust the height of the smartphone to obtain images from different angles. Furthermore, we try our best to make sure captured images contain the entire car body without truncation. To preserve the privacy of owners, we make license plates and other private information obfuscated. To construct a high-quality dataset, we filter out some instances with blurred, out-of-focus, and occluded images. We also provide detailed annotations for car brands, types, and colors. Particularly, we provide the car parsing maps for each car with thirteen classes in our dataset as shown in Figure 1 for the advancement of car component understanding tasks.

3.2 DATA PREPROCESSING

Background Removal. Since we only reconstruct cars for the 3D car reconstruction task, the background should be removed. Recent Segment Anything Model (SAM) (Kirillov et al., 2023) demonstrates powerful context recognition and segmentation performance. However, SAM needs a bounding box, text, or point as a driving factor for accurate segmentation. Therefore, we employ Grounding DINO (Liu et al., 2023d) as a text-driven detector with a detection prompt with “car” for

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

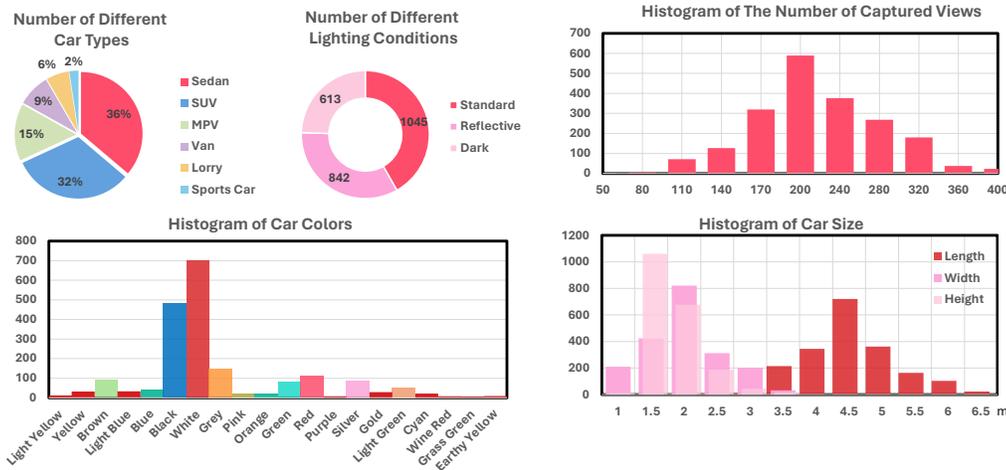


Figure 5: **The distributions of our 3DRealCar dataset.** We show distributions of car types, lighting conditions, captured views, car colors, and car size. We try our best to capture cars with various colors and types for the diversity of our dataset.

the attainment of car bounding boxes. With these bounding boxes, we use SAM to obtain the masks from captured images. The point cloud initialization is demonstrated useful for the convergence of 3D Gaussian Splatting (Kerbl et al., 2023b). Except for the removal of the background in 2D images, we still need to remove the background point clouds. Therefore, we first project the 3D point clouds into 2D space with camera parameters. Then, we can eliminate background point clouds with masks and save them for further processing.

Orientation Rectification. As shown in Figure 4, we utilize Colmap (Schonberger & Frahm, 2016) to reconstruct more dense point clouds and obtain accurate camera poses and intrinsics because we find that the estimated poses by the smartphone are not accurate. However, after the removal of the background point clouds, we find that the car orientation of the point cloud is random, which leads to the subsequent render task being uncontrollable. Given camera poses $P = \{p_i\}_1^{\mathcal{N}}$, where \mathcal{N} is the number of poses, we use Principal Component Analysis (PCA) (Abdi & Williams, 2010) to obtain a PCA component $\mathcal{T} \in \mathbb{R}^{3 \times 3}$. The PCA component is the principal axis of the data in 3D space, which represents rotation angles to each axis. Therefore, we leverage it to rectify the postures of cars parallel to the x-axis. However, this process cannot guarantee cars facing along the x-axis. Therefore, in some failure cases, we manually interfere and adjust the orientation along the x-axis. With the fixed car orientation, we can control rendered poses for the subsequent tasks.

Point Cloud Rescaling. The size of the point clouds reconstructed by Colmap (Schonberger & Frahm, 2016) does not match the real-world size, which inhibits the reconstruction of a practically sized 3D car. To address this, we calculate the bounding box of the scanned foreground point clouds to obtain its actual size in the real-world scene. Then, we rescale the rectified point clouds into the real size. In addition to the rescaling of the point clouds, we also need to adjust the camera poses. We rescale translations of camera poses using a scale factor calculated by the ratio of scanned point cloud size and Colmap point cloud size. After these rescaling processes, we use rescaled point clouds to reconstruct a 3D car model through recent state-of-the-art methods, like 3DGS (Kerbl et al., 2023b).

3.3 DATA STATISTICS

In our 3DRealCar, we provide detailed annotations for researchers to leverage our dataset for different tasks. During the data annotating, we discard the data with the number of views less than fifty. As we can observe in Figure 1 and 2, we collect our dataset under real-world scenes and meticulously scan dense views. Therefore, cars in our dataset possess dense views and realistic texture, which is necessary for the application in a real-world setting.

As shown in Figure 5, we conduct detailed statistical analysis to show the features of our dataset. Our dataset mainly contains six different car types, such as Sedan, SUV, MPV, Van, Lorry, and Sports Car.

Among them, sedans and SUVs are common to collect in real life, so their volume dominates in our dataset. We also count the number of different lighting conditions on cars. The standard condition means the car is well-lit and without strong specular highlights. The reflective condition means the car has strong specular highlights. Glossy materials bring huge challenges to recent 3D reconstruction methods. The dark condition means the car is captured in an underground parking so not well-lit. To promote high-quality reconstruction, we save the captured images in high resolution (1920×1440) and also capture as many views as possible. The number of captured images per car is an average of 200. The number of views ranges from 50 to 400. To enrich the diversity of our dataset, we try our best to collect as many different colors as possible. Therefore, our dataset contains more than twenty colors, but the white and black colors still take up most of our dataset. In addition, we also show the distribution of car size, in terms of their length, width, and height. We obtain their sizes by computing the bounding boxes of the scanned point clouds. Thanks to different car types, the sizes of cars are also diverse.

4 OVERVIEW OF 3DREALCAR TASKS

4.1 2D TASKS

Corner-case scene 2D Detection (Ultralytics, 2023; Zhang et al., 2022a; Zong et al., 2023): Given a serial of images $I = \{I_i\}_1^N$, this task aims to detect vehicles as accurately as possible. However, in some corner cases, like car accidents, detectors sometimes fail to detect target vehicles since this kind of scene is rare or not in the training set. Therefore, this task has crucial significance in building a reliable self-driving system, especially for accident scenarios.

2D Car Parsing (Chen et al., 2017; Hong et al., 2021; Xie et al., 2021; Kirillov et al., 2020): Given a serial of images $I = \{I_i\}_1^N$, this task aims to segment car parsing maps $S = \{S_i\}_1^N$. With annotated parsing maps, we can train a model to understand and segment each component of cars. This task can assist self-driving systems with more precise recognition.

4.2 3D TASKS

Neural Field-based Novel View Synthesis (Müller et al., 2022; Kerbl et al., 2023b; Huang et al., 2024): Given a serial of images $I = \{I_i\}_1^N$ and matched poses $P = \{p_i\}_1^N$, where N is the number of images and poses, the task of Neural Field-based Novel View Synthesis aims to reconstruct Neural Field model of a object or a scene. The reconstructed model is usually used to render 2D images with different views for the evaluation of the performance of novel view synthesis.

Diffusion-based Novel View Synthesis (Liu et al., 2023c;e; Stability.AI, 2023): Given a serial of reference images $I^{ref} = \{I_i^{ref}\}_1^N$, reference poses $P^{ref} = \{p_i^{ref}\}_1^N$, target images $I^{target} = \{I_i^{target}\}_1^N$, and target poses $P^{target} = \{p_i^{target}\}_1^N$, recent 3D generative models, such as Zero123 (Liu et al., 2023c), Syncdreamer (Liu et al., 2023e), and Stable-Zero123 (Stability.AI, 2023), take relative poses and reference images as inputs and generate target images. However, these models cannot generalize well to real car objects since they are trained on large-scale synthetic datasets (Deitke et al., 2023; 2024). In this work, we will demonstrate that our dataset can improve the robustness of these generative models to real cars.

Single Image to 3D Generation (Poole et al., 2022a; Sun et al., 2023; Tang et al., 2023): Given a text prompt or single image, recent 3D generation methods generate 3D objects with Score Distillation Sampling (SDS) (Poole et al., 2022a) and diffusion generative models (Rombach et al., 2022; Liu et al., 2023c; Stability.AI, 2023). However, these methods cannot generate high-quality 3D cars due to the lack of the prior of real cars in 3D-based diffusion models. Therefore, we would demonstrate the value of our dataset by improving recent 3D generation for real cars.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Corner-case 2D Detection. In this task, we leverage the reconstructed cars to simulate rare and corner-case scenes. To be specific, we use Nuscenets (Caesar et al., 2020) as background to

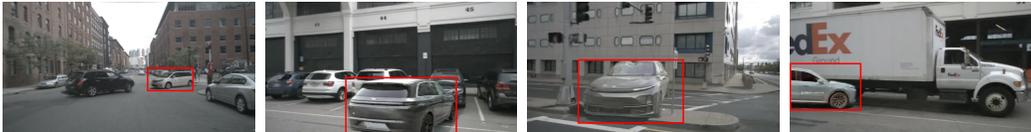


Figure 6: **The simulated corner-case scenes.** These scenes are rare but very important in real life. We use a red rectangle to highlight the simulated vehicles. These corner-case scenes show some vehicles have potential risks to traffic safety.

Table 2: **Detection improvements by simulated data for corner-case scenes.** We leverage lightweight YOLO serials models and recent state-of-the-art models for evaluation. We report the metric by calculating mAP@0.5 on the CODA dataset Li et al. (2022).

Simulated Data	YOLOv5n	YOLOv5s	YOLOv8n	YOLOv8s	DINO	CO-DETR
1000	0.285	0.341	0.299	0.371	0.437	0.465
2000	0.304	0.357	0.312	0.366	0.452	0.481
3000	0.345	0.389	0.357	0.403	0.495	0.517
4000	0.357	0.408	0.386	0.413	0.543	0.551
5000	0.361	0.426	0.386	0.435	0.571	0.582

simulate corner-case scenes with reconstructed cars and leverage recent popular detectors, like YOLOv8 (Ultralytics, 2023), as detectors for evaluation. To evaluate the robustness of detectors in corner-case scenes, we use the test part of the corner-case dataset, CODA (Li et al., 2022) as a testing set. Since we focus on the corner-case scenes of cars, so we only evaluate a car class.

2D Car Parsing. In this task, we utilize DeepLabV3(Chen et al., 2017), DDRNet (Hong et al., 2021), SegFormer (Xie et al., 2021), and PointRend (Kirillov et al., 2020) to benchmark our dataset. To be specific, we split 80% of our car parsing maps in 3DRealCar as the training set and the rest of 20% as the testing set.

Neural Field-based Novel View Synthesis. In this task, we randomly choose 100 instances from each lighting condition in our dataset and split 80% of the views per instance as the training set and the rest of 20% as the testing set. Specifically, we employ recent state-of-the-art neural field methods, including Instant-NGP (Müller et al., 2022), 3DGS (Kerbl et al., 2023b), GaussianShader (Jiang et al., 2023), and 2DGS (Huang et al., 2024) to benchmark our dataset.

Diffusion-based Novel View Synthesis. We finetune Zero123-XL (Liu et al., 2023c) on our 3DRealCar dataset to enhance its generalization to real cars. Note that since the training of diffusion-based models needs entire objects centered on images, we use the images rendered by our trained 3D models as training images.

Single Image to 3D Generation. In this task, we exploit Dreamcraft3D (Sun et al., 2023) as our baseline. Dreamcraft3D exploits Stable-Zero123 (Stability.AI, 2023) as a prior source for providing 3D generative prior. By fine-tuning Stable-Zero123 on our dataset, we enable it to obtain car-specific prior so it generalizes well to real cars.

5.2 2D TASKS

Corner-case 2D Detection. As shown in Table 2, we employ YOLOv5 and YOLOv8 serial models, DINO (Zhang et al., 2022a), and CO-DETR Zong et al. (2023) as our detectors for evaluation. To evaluate the performance of models in corner-case scenes, we leverage the test part of the CODA dataset (Li et al., 2022) as our testing set. In particular, when we increase the training simulated data from 500 to 5,000, the performance of detectors is also improved by a large margin. This phenomenon demonstrates that our simulated data is effective in improving a detector robust to corner-case scenes. We provide the visualizations of simulated corner-case scenes in Figure 6. The detailed simulation process and more visualizations can be seen in the supplementary.

2D Car Parsing. We conduct benchmarks for car parsing maps of our dataset using recent image segmentation methods, such as DeepLabV3(Chen et al., 2017), PointRend(Kirillov et al., 2020),

Table 3: **Benchmark results on 2D car parsing of our 3DRealCar dataset.** We use recent advanced image segmentation methods Chen et al. (2017); Hong et al. (2021); Xie et al. (2021); Kirillov et al. (2020) to benchmark our dataset.

Method	DeepLabV3	PointRend	DDRNet	SegFormer
mIOU \uparrow	0.556	0.562	0.603	0.613
mAcc \uparrow	0.616	0.619	0.659	0.663

Table 4: **Quantitative comparisons of SOTA 3D Generation method, Dreamcraft3D Sun et al. (2023) and its improved version by trained on our dataset.** CD denotes Chamfer Distance.

Method	CLIP-I \uparrow	Hausdorff \downarrow	CD \downarrow
Dreamcraft3D	0.812	1.572	0.587
+our dataset	0.847	1.364	0.371

Table 5: **Benchmark results on 3D reconstruction of our 3DRealCar dataset.** We present the 3D reconstruction performance of recent state-of-the-art methods in three lighting conditions, standard, reflective, and dark, respectively. The best results are highlighted.

Method	Standard			Reflective			Dark		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Instant-NGP Müller et al. (2022)	27.31	0.9315	0.1264	24.37	0.8613	0.1962	23.17	0.9152	0.1642
3DGS Kerbl et al. (2023b)	27.47	0.9367	0.1001	24.58	0.8647	0.1852	23.51	0.9181	0.1613
GaussianShader Jiang et al. (2023)	27.53	0.9311	0.1109	25.41	0.8684	0.1423	23.39	0.9172	0.1631
2DGS Huang et al. (2024)	27.34	0.9341	0.1095	23.19	0.8509	0.2041	22.63	0.9148	0.1681

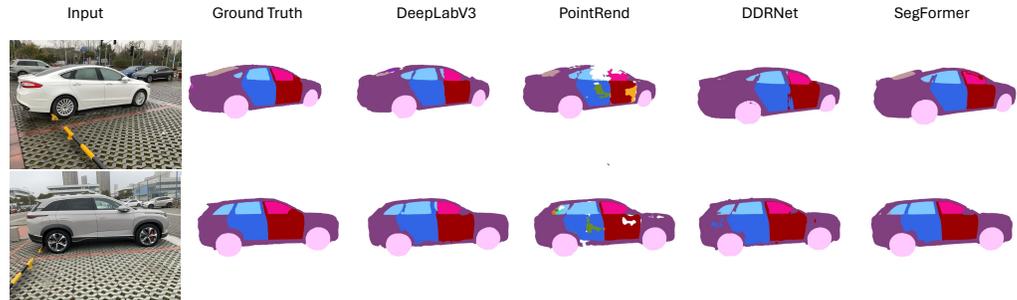


Figure 7: **Qualitative comparisons among recent advanced image segmentation methods.** We select the inputs from the testing set of our images and evaluate the capacity of car component understanding for each method.

DDRNet(Hong et al., 2021), and SegFormer(Xie et al., 2021). The quantitative performance for these methods on our dataset is summarized in Table 3. Visual comparisons are provided in Figure 7. Our high-quality dataset enables these methods to achieve promising performance, highlighting its potential for application in self-driving systems. In particular, our car parsing annotations encourage self-driving systems to recognize different components of cars in practical scenarios for safer automatic decisions.

5.3 3D TASKS

Neural Field-based Novel View Synthesis. As depicted in Table 5, we show benchmark results of recent state-of-the-art neural field methods, such as Instant-NGP (Müller et al., 2022), 3DGS (Kerbl et al., 2023b), GaussianShader (Jiang et al., 2023), and 2DGS (Huang et al., 2024) on our dataset. To the standard lighting condition, we can find that recent methods are capable of achieving PSNR more than 27 dB, which means these methods can reconstruct relatively high-quality 3D cars from our dataset. However, the reflective and dark condition results are lower than the standard. These two parts of our 3DRealCar bring two challenges to recent 3D methods. The first challenge is the reconstruction of specular highlights. Due to the particular property of cars, materials of car surfaces are generally glossy, which means it would produce plenty of specular highlights if cars are exposed to the sun or strong light. The second challenge is the reconstruction in a dark environment. The training images captured in the dark environment lose plenty of details for reconstruction. Therefore, how to achieve high-quality reconstruction results from these two extremely lighting conditions is a challenge to recent methods. 3D visualizations can be found on our project page.



497 **Figure 8: Visualizations of diffusion-based novel view synthesis.** we compare the results of the
 498 recent state-of-the-art diffusion-based method, Zero123-XL (Liu et al., 2023c) and its improvement
 499 by training on our dataset. Our dataset provides car-specific prior for the generative model to generate
 500 more photorealistic car images.



512 **Figure 9: Visualizations of single-image-to-3D generation.** we compare the results of the recent
 513 state-of-the-art single-image-to-3D method, Dreamcraft3D (Sun et al., 2023) and is enhanced version
 514 by training on our dataset.

517 **Diffusion-based Novel View Synthesis.** As illustrated in Figure 8, we show visual comparisons of
 518 Zero123-XL (Liu et al., 2023c) and our improved version by training on our dataset. As we can see,
 519 given input images, we use Zero123-XL and our improved version to synthesize novel views. In this
 520 figure, we can find that Zero123-XL prefers to generate synthetic results with unrealistic texture and
 521 geometry, due to the lack of prior for real objects. In contrast, our improved version of Zero123-XL
 522 can generate photorealistic geometry and texture, which demonstrates the effectiveness of our dataset.

523 **Single Image to 3D Generation.** As depicted in Figure 9, we visualize 3D generation results of the
 524 recent state-of-the-art single-image-to-3D method, Dreamcraft3D (Sun et al., 2023), along with its
 525 improved version by our dataset. This figure shows that Dreamcraft3D sometimes fails to generate
 526 complete geometry or realistic texture, due to the scarcity of the real car prior. As shown in Table 4,
 527 we also show quantitative comparisons of Dreamcraft3D and its improved version. CLIP-I means the
 528 similarity of rendered images with the original input. The quantitative and qualitative results indicate
 529 our dataset significantly improves 3D generation performance typically in terms of geometry and
 530 texture. These results underscore the effectiveness of our 3DRealCar dataset.

531 6 CONCLUSION

532

533 In this paper, we propose the first large-scale high-quality 3D real car dataset, named 3DRealCar.
 534 The collected dense and high-resolution 360-degree views for each car can be used to reconstruct a
 535 high-quality 3D car. Extensive experiments demonstrate the efficacy and challenges of our 3DRealCar
 536 in 3D reconstruction. Thanks to the reconstructed high-quality 3D cars from our dataset and car-part
 537 level annotations, our dataset can be utilized to support various tasks related to cars. In addition, the
 538 benchmarking results can serve as baselines for prospective research. Although 3DRealCar currently
 539 only has car exterior views, we intend to provide both exterior and interior views in the future to
 further promote the reconstruction of more intact 3D cars.

REFERENCES

- 540
541
542 Hervé Abdi and Lynne J Williams. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics*, 2(4):433–459, 2010.
- 543
544 Apple, 2021. URL <https://developer.apple.com/augmented-reality/arkit/>.
- 545
546 Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and
547 Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields.
548 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5855–5864,
549 2021.
- 550
551 Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-NeRF
552 360: Unbounded Anti-Aliased Neural Radiance Fields. *2022 IEEE/CVF Conference on Computer
553 Vision and Pattern Recognition (CVPR)*, pp. 5460–5469, 2022a.
- 554
555 Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf
556 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference
557 on Computer Vision and Pattern Recognition*, pp. 5470–5479, 2022b.
- 558
559 Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf:
560 Anti-aliased grid-based neural radiance fields. *ICCV*, 2023.
- 561
562 Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush
563 Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nusenes: A multimodal dataset for
564 autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
565 recognition*, pp. 11621–11631, 2020.
- 566
567 Florian Chabot, Mohamed Chaouch, Jaonary Rabarisoa, Céline Teuliere, and Thierry Chateau. Deep
568 manta: A coarse-to-fine many-task network for joint 2d and 3d vehicle analysis from monocular
569 image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
570 2040–2049, 2017.
- 571
572 Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li,
573 Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d
574 model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- 575
576 Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su.
577 Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings
578 of the IEEE/CVF International Conference on Computer Vision*, pp. 14124–14133, 2021.
- 579
580 Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields.
581 In *European Conference on Computer Vision*, pp. 333–350. Springer, 2022a.
- 582
583 Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous
584 convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- 585
586 Tianlong Chen, Peihao Wang, Zhiwen Fan, and Zhangyang Wang. Aug-nerf: Training stronger
587 neural radiance fields with triple-level physically-grounded augmentations. In *Proceedings of the
588 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15191–15202, 2022b.
- 589
590 Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision
591 transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022c.
- 592
593 Julian Chibane, Aayush Bansal, Verica Lazova, and Gerard Pons-Moll. Stereo radiance fields (srf):
Learning view synthesis from sparse views of novel scenes. In *IEEE Conference on Computer
Vision and Pattern Recognition (CVPR)*. IEEE, jun 2021.
- Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3213–3223, 2016.

- 594 Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig
595 Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of anno-
596 tated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
597 Recognition*, pp. 13142–13153, 2023.
- 598 Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan
599 Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of
600 10m+ 3d objects. *Advances in Neural Information Processing Systems*, 36, 2024.
- 601 Zhiwen Fan, Yifan Jiang, Peihao Wang, Xinyu Gong, Dejia Xu, and Zhangyang Wang. Unified
602 implicit neural stylization. In *European Conference on Computer Vision*, pp. 636–654. Springer,
603 2022.
- 604 Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo
605 Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF
606 Conference on Computer Vision and Pattern Recognition*, pp. 5501–5510, 2022.
- 607 Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf:
608 High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference
609 on Computer Vision*, pp. 14346–14355, 2021.
- 610 Carlos J García Orellana, Ramón Gallardo Caballero, Horacio M González Velasco, and Francisco J
611 López Aligué. Neusim: a modular neural networks simulator for beowulf clusters. In *Bio-Inspired
612 Applications of Connectionism: 6th International Work-Conference on Artificial and Natural
613 Neural Networks, IWANN 2001 Granada, Spain, June 13–15, 2001 Proceedings, Part II 6*, pp.
614 72–79. Springer, 2001.
- 615 Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti
616 dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- 617 Yuan-Chen Guo, Di Kang, Linchao Bao, Yu He, and Song-Hai Zhang. Nerfren: Neural radiance
618 fields with reflections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
619 Pattern Recognition*, pp. 18409–18418, 2022.
- 620 Yuanduo Hong, Huihui Pan, Weichao Sun, and Yisong Jia. Deep dual-resolution networks for
621 real-time and accurate semantic segmentation of road scenes. *arXiv preprint arXiv:2101.06085*,
622 2021.
- 623 Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for
624 geometrically accurate radiance fields. *arXiv preprint arXiv:2403.17888*, 2024.
- 625 Xinyu Huang, Xinjing Cheng, Qichuan Geng, Binbin Cao, Dingfu Zhou, Peng Wang, Yuanqing Lin,
626 and Ruigang Yang. The apollo-scape dataset for autonomous driving. In *Proceedings of the IEEE
627 conference on computer vision and pattern recognition workshops*, pp. 954–960, 2018.
- 628 Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided
629 object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer
630 Vision and Pattern Recognition*, pp. 867–876, 2022.
- 631 Yingwenqi Jiang, Jiadong Tu, Yuan Liu, Xifeng Gao, Xiaoxiao Long, Wenping Wang, and Yuexin
632 Ma. Gaussianshader: 3d gaussian splatting with shading functions for reflective surfaces. *arXiv
633 preprint arXiv:2311.17977*, 2023.
- 634 M. M. Johari, Y. Lepoittevin, and F. Fleuret. Geonerf: Generalizing nerf with geometry priors.
635 *Proceedings of the IEEE international conference on Computer Vision and Pattern Recognition
636 (CVPR)*, 2022.
- 637 Heewoo Jun and Alex Nichol. Shap-E: Generating conditional 3D implicit functions, 2023.
- 638 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting
639 for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023a.
- 640 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting
641 for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023b.

- 648 Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. Pointrend: Image segmentation as
649 rendering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
650 pp. 9799–9808, 2020.
- 651 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
652 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint*
653 *arXiv:2304.02643*, 2023.
- 654 Xin Kong, Shikun Liu, Xiaoyang Lyu, Marwan Taher, Xiaojuan Qi, and Andrew J Davison. Eschnet:
655 A generative model for scalable view synthesis. *arXiv preprint arXiv:2402.03908*, 2024.
- 656 Kaican Li, Kai Chen, Haoyu Wang, Lanqing Hong, Chaoqiang Ye, Jianhua Han, Yukuai Chen, Wei
657 Zhang, Chunjing Xu, Dit-Yan Yeung, et al. Coda: A real-world road corner case dataset for object
658 detection in autonomous driving. *arXiv preprint arXiv:2203.07724*, 2022.
- 659 Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten
660 Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content
661 creation. In *CVPR*, 2023.
- 662 Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen,
663 Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3D objects with
664 consistent multi-view generation and 3D diffusion. *arXiv preprint arXiv:2311.07885*, 2023a.
- 665 Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Zexiang Xu, Hao Su, et al. One-2-3-45: Any single
666 image to 3d mesh in 45 seconds without per-shape optimization. *arXiv preprint arXiv:2306.16928*,
667 2023b.
- 668 Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick.
669 Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International*
670 *Conference on Computer Vision*, pp. 9298–9309, 2023c.
- 671 Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei
672 Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for
673 open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023d.
- 674 Tianyu Liu, Hao Zhao, Yang Yu, Guyue Zhou, and Ming Liu. Car-studio: Learning car radiance
675 fields from single-view and unlimited in-the-wild images. *IEEE Robotics and Automation Letters*,
676 2024.
- 677 Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang.
678 Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint*
679 *arXiv:2309.03453*, 2023e.
- 680 Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. Scaffold-gs:
681 Structured 3d gaussians for view-adaptive rendering. *arXiv preprint arXiv:2312.00109*, 2023.
- 682 Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, Natalia Neverova, Andrea Vedaldi, Oran Gafni,
683 and Filippos Kokkinos. IM-3D: Iterative multiview diffusion and reconstruction for high-quality
684 3D generation. *arXiv preprint arXiv:2402.08682*, 2024.
- 685 Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for
686 shape-guided generation of 3d shapes and textures. In *Proceedings of the IEEE/CVF Conference*
687 *on Computer Vision and Pattern Recognition*, pp. 12663–12673, 2023.
- 688 Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren
689 Ng. Nerf: Representing Scenes As Neural Radiance Fields for View Synthesis. *Communications*
690 *of the ACM*, 65(1):99–106, 2021. doi: 10.1145/3503250.
- 691 Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics
692 primitives with a multiresolution hash encoding. *arXiv:2201.05989*, January 2022.
- 693 Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics
694 primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):
695 1–15, 2022.

- 702 Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-E: A System
703 for Generating 3D Point Clouds from Complex Prompts, 2022.
704
- 705 Li Niu, Wenyan Cong, Liu Liu, Yan Hong, Bo Zhang, Jing Liang, and Liqing Zhang. Making images
706 real again: A comprehensive survey on deep image composition. *arXiv preprint arXiv:2106.14490*,
707 2021.
- 708 Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d
709 diffusion. *arXiv preprint arXiv:2209.14988*, 2022a.
710
- 711 Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d
712 diffusion. *arXiv preprint arXiv:2209.14988*, 2022b.
- 713 Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee,
714 Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, et al. Magic123: One image to high-quality 3d
715 object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843*, 2023.
716
- 717 Erik Reinhard, Michael Adhikhmin, Bruce Gooch, and Peter Shirley. Color transfer between images.
718 *IEEE Computer graphics and applications*, 21(5):34–41, 2001.
- 719 Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural
720 radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International
721 Conference on Computer Vision*, pp. 14335–14345, 2021.
722
- 723 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
724 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-
725 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 726 Johannes L. Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *2016 IEEE
727 Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2016. doi: 10.1109/cvpr.
728 2016.445. URL <http://dx.doi.org/10.1109/cvpr.2016.445>.
729
- 730 Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen,
731 Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base
732 model. *arXiv preprint arXiv:2310.15110*, 2023a.
- 733 Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view
734 diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023b.
735
- 736 Stability.AI. Stable Zero123: Quality 3d object generation from single images. [https://
737 stability.ai/news/stable-zero123-3d-generation](https://stability.ai/news/stable-zero123-3d-generation), 2023.
- 738 Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Light field neural rendering.
739 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
740 8269–8279, 2022.
- 741 Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct Voxel Grid Optimization: Super-Fast
742 Convergence for Radiance Fields Reconstruction. In *Proceedings of the IEEE/CVF Conference on
743 Computer Vision and Pattern Recognition*, pp. 5459–5469, 2022.
744
- 745 Jingxiang Sun, Bo Zhang, Ruizhi Shao, Lizhen Wang, Wen Liu, Zhenda Xie, and Yebin Liu.
746 Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior. *arXiv preprint
747 arXiv:2310.16818*, 2023.
- 748 Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James
749 Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous
750 driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision
751 and pattern recognition*, pp. 2446–2454, 2020.
752
- 753 Mukund Varma T, Peihao Wang, Xuxi Chen, Tianlong Chen, Subhashini Venugopalan, and
754 Zhangyang Wang. Is attention all that nerf needs? In *The Eleventh International Confer-
755 ence on Learning Representations*, 2023. URL [https://openreview.net/forum?id=
xE-LtSE-xx](https://openreview.net/forum?id=xE-LtSE-xx).

- 756 Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. DreamGaussian: Generative
757 Gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.
758
- 759 Ultralytics. YOLOv8: A cutting-edge and state-of-the-art (sota) model that builds upon the success
760 of previous yolo versions. [https://github.com/ultralytics/ultralytics?tab=](https://github.com/ultralytics/ultralytics?tab=readme-ov-file)
761 [readme-ov-file](https://github.com/ultralytics/ultralytics?tab=readme-ov-file), 2023.
- 762 Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. *arXiv preprint arXiv:2112.03907*, 2021.
763
764
- 765 Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image
766 driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on*
767 *Computer Vision and Pattern Recognition*, pp. 3835–3844, 2022.
768
- 769 Jingkan Wang, Sivabalan Manivasagam, Yun Chen, Ze Yang, Ioan Andrei Bârsan, Anqi Joyce Yang,
770 Wei-Chiu Ma, and Raquel Urtasun. Cadsim: Robust and scalable in-the-wild 3d reconstruction for
771 controllable sensor simulation. *arXiv preprint arXiv:2311.01447*, 2023a.
- 772 Peng Wang, Yuan Liu, Zhaoxi Chen, Lingjie Liu, Ziwei Liu, Taku Komura, Christian Theobalt, and
773 Wenping Wang. F2-nerf: Fast neural radiance field training with free camera trajectories. *CVPR*,
774 2023b.
775
- 776 Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron,
777 Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view
778 image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
779 *Pattern Recognition*, pp. 4690–4699, 2021.
- 780 Pengchuan Xiao, Zhenlei Shao, Steven Hao, Zishuo Zhang, Xiaolin Chai, Judy Jiao, Zesong Li, Jian
781 Wu, Kai Sun, Kun Jiang, et al. Pandaset: Advanced sensor suite dataset for autonomous driving.
782 In *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pp. 3095–3101.
783 IEEE, 2021.
- 784 Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer:
785 Simple and efficient design for semantic segmentation with transformers. *Advances in neural*
786 *information processing systems*, 34:12077–12090, 2021.
787
- 788 Ziyang Xie, Junge Zhang, Wenye Li, Feihu Zhang, and Li Zhang. S-nerf: Neural radiance fields for
789 street views. *arXiv preprint arXiv:2303.00749*, 2023.
- 790 Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from
791 one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
792 *Recognition*, pp. 4578–4587, 2021.
793
- 794 Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung
795 Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv*
796 *preprint arXiv:2203.03605*, 2022a.
- 797 Jason Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. Ners: Neural reflectance
798 surfaces for sparse-view 3d reconstruction in the wild. *Advances in Neural Information Processing*
799 *Systems*, 34:29835–29847, 2021.
- 800 Kai Zhang, Nick Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snavely. Arf:
801 Artistic radiance fields, 2022b.
802
- 803 Chuanxia Zheng and Andrea Vedaldi. Free3D: Consistent novel view synthesis without 3D represen-
804 tation. *arXiv preprint arXiv:2312.04551*, 2023.
- 805 Xiaoyu Zhou, Zhiwei Lin, Xiaojun Shan, Yongtao Wang, Deqing Sun, and Ming-Hsuan Yang.
806 Drivinggaussian: Composite gaussian splatting for surrounding dynamic autonomous driving
807 scenes. *arXiv preprint arXiv:2312.07920*, 2023.
808
- 809 Zhuofan Zong, Guanglu Song, and Yu Liu. Detsr with collaborative hybrid assignments training. In
Proceedings of the IEEE/CVF international conference on computer vision, pp. 6748–6758, 2023.

A BROADER IMPACTS STATEMENT

The introduction of our 3DRealCar dataset has profound effects on self-driving research. We expect this dataset can encourage extensive research to promote the advancement of the community.

Research Impacts. By providing dense 360-degree views of cars with point clouds as initialization, our 3DRealCar can be used to reconstruct high-quality 3D real cars for 3D printing and simulation in corner-case scenes. By providing detailed car parsing map annotations, our dataset can be leveraged to segment 2D car components or point clouds. Note that our 3DRealCar is the first dataset providing 3D car parsing annotations. In our 3D reconstruction benchmarking experiments, the reflective and dark lighting conditions of our dataset bring challenges to existing methods to reconstruct 3D cars under awful lighting conditions. We expect our dataset to encourage widespread collaboration and accelerate the exploration of 3D real car reconstruction, parsing, and simulation.

Societal Impacts. We collect our 3DRealCar dataset with the consent of the owners. In addition, we blur license plates and other private information. We try our best to hide and preserve the privacy of owners. Therefore, our dataset would not have any privacy violation problems. Due to our dataset focusing on a car class, we believe our dataset has the potential to be employed in future self-driving research and improve self-driving systems further.

B LIMITATION AND DISCUSSION

Although our 3DRealCar is the largest dataset for the 3D real car dataset so far (2500 car instances with annotations), its scale is still limited compared to other datasets in the computer vision community. Therefore, we will further extend our dataset in the future. Moreover, our 3DRealCar dataset only provides the exterior views of cars without interior views. It is very crucial to reconstruct both exterior and interior views of cars for car marketing agencies. We will collect both exterior and interior views in the future to further extend our 3D real car dataset for intact 3D car models.

C EXPERIMENTAL SETTINGS

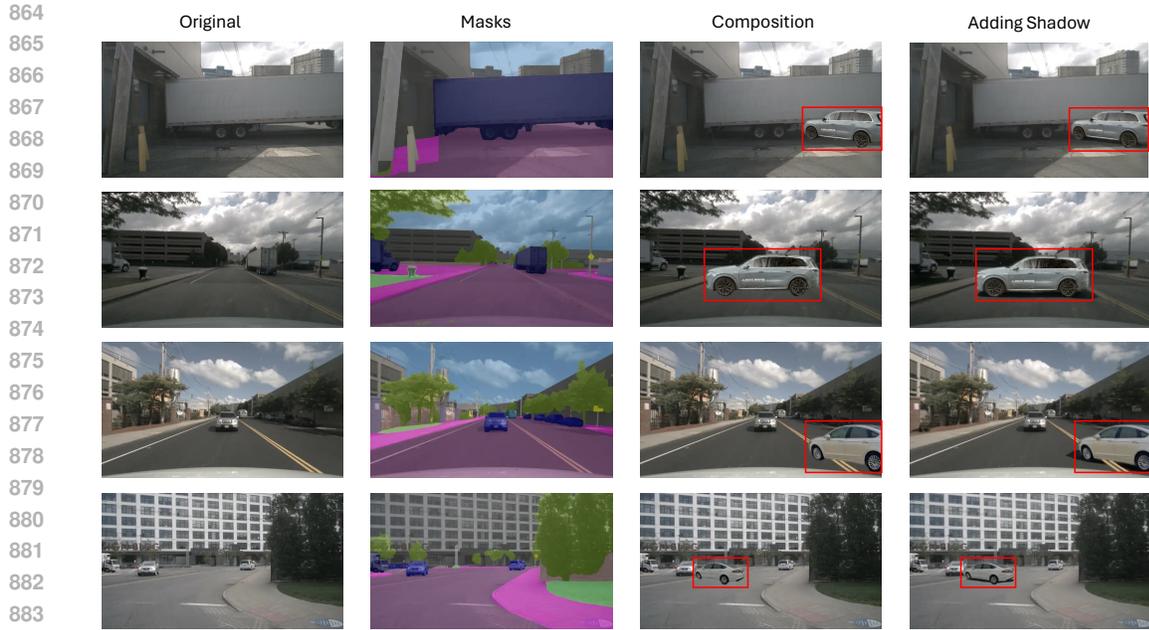
Note that all models used in this work are publicly available. Each model we use is linked below:

1. **3D Reconstruction:** Instant-NGP (Müller et al., 2022) [🔗](#), 3DGS (Kerbl et al., 2023b) [🔗](#), GaussianShader (Jiang et al., 2023) [🔗](#), and 2DGS (Huang et al., 2024) [🔗](#).
2. **2D Car Parsing:** MMsegmentation [🔗](#). This repository includes all 2D segmentation models (Chen et al., 2017; Hong et al., 2021; Xie et al., 2021; Kirillov et al., 2020) we used in this work.
3. **Novel View Synthesis:** Zero-123-XL (Liu et al., 2023c) [🔗](#).
4. **3D Generation:** DreamCraft3D (Sun et al., 2023) [🔗](#).
5. **Corner-case Simulation:** YOLOv5 and YOLOv8 (Ultralytics, 2023) [🔗](#), DINO (Zhang et al., 2022a) [🔗](#), CO-DETR (Zong et al., 2023) [🔗](#), and libcom (Niu et al., 2021) [🔗](#). Specifically, we use YOLOv5 and YOLOv8 serial models, DINO, and CO-DETR as detectors and libcom for the simulation of corner-case scenes.

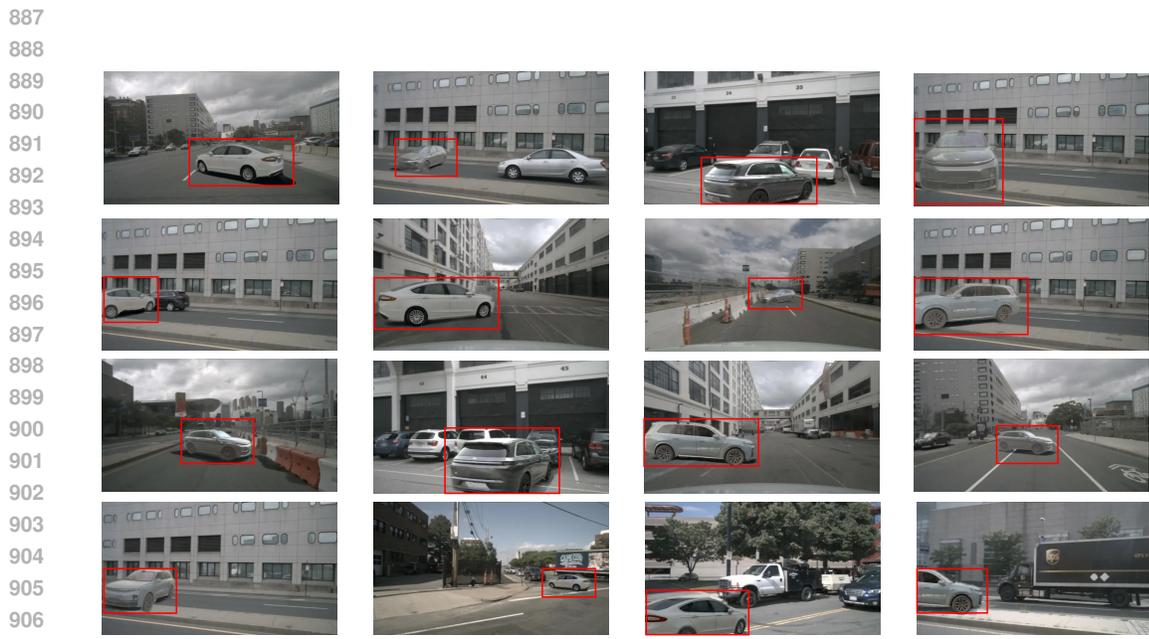
We express great appreciation to the authors of the aforementioned repositories for their invaluable contributions. For the GPU specification, we use 8 A100 GPUs for 3D reconstruction, 3D generation, and novel view synthesis. We utilize 2 3090 GPUs for other tasks. We use the default hyperparameters for training.

D DETAILED SIMULATION PROCESS AND ADDITIONAL VISUALIZATIONS

In this section, we show how we simulate corner-case scenes. As shown in Figure 10, we use images from Nuscenes (Caesar et al., 2020) as backgrounds and leverage ViT-Adapter (Chen et al., 2022c) to segment entire scenes for road masks. Then, we copy and paste the rendered images from the reconstructed high-quality 3D cars into the backgrounds with the guidance of road masks. In



885 **Figure 10: Visualizations of ablating simulation procedures.** We use a red rectangle to highlight
886 the simulated vehicles.



908 **Figure 11: More visualizations of simulated corner-case scenes.** We use a red rectangle to
909 highlight simulated vehicles. These corner-case scenes show some vehicles have potential risks to
910 traffic safety.

911
912
913 particular, we blur the edge between simulated foregrounds and backgrounds and then we use a
914 color transfer algorithm (Reinhard et al., 2001) to make the whole simulated scene look harmonious.
915 Finally, we use the shadow generation method in libcom (Niu et al., 2021) to add shadow for the
916 simulated cars such that the entire scene looks photorealistic. However, this simulation method
917 would generate some unreasonable scenes. Therefore, we manually intervene to select photorealistic
corner-case scenes. Additional simulation results are shown in Figure 11.

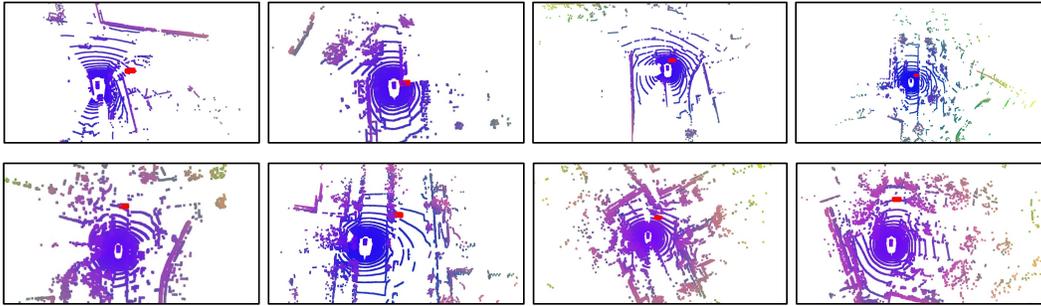


Figure 12: **Visualizations of point cloud inserting.** We use red color to annotate inserting vehicular point clouds with high density for better differentiation.



Figure 13: **Visualizations of 3D point cloud parsing.** With 2D car parsing map annotations, we lift the 2D car parsing maps into 3D point clouds and segment car components.

E SIMULATED LIDAR SCENES

As depicted in Figure 12, we can insert our car point clouds into lidar scenes to simulate corner-case scenes, like a car passing or parking horizontally in front of the ego car. To better differentiate the inserted cars, we set them with dense point clouds and red color. In a practical scene, the vehicular point clouds should be sparse and only have one side that could be scanned by the lidar. Therefore, when we apply the inserted vehicular point clouds into a scene, we should make the vehicular point clouds sparse and only contain one side. By training on a variety of simulated scenarios, including rare or dangerous situations that are difficult to collect in real life, the self-driving system can learn to handle unexpected events more effectively.

F 3D CAR PARSING

As shown in Figure 13, our dataset is the first to provide 3D car parsing annotations for parsing car components in 3D space. Thanks to that we provide 2D car parsing maps for every instance in our 3DRealCar dataset, we can lift 2D parsing maps to 3D and segment each component for point clouds and meshes. The primary purpose of these 3D car parsing maps is to enable precise and comprehensive analysis of vehicle structures, which is crucial for applications such as autonomous driving, vehicle design, vehicle editing, and virtual reality simulations. By using these detailed 3D parsing maps, developers and researchers can improve object recognition algorithms and enhance collision detection systems. Furthermore, this dataset facilitates the training of machine learning models to better understand the spatial relationships and physical attributes of car components, leading to more advanced and reliable automotive technologies.

G CONSENT FORM FOR 3DREALCAR

Since our dataset includes license plate information, we obtain consent from participants and require them to sign the consent form shown in Figure 14 before data collection. We ensure that no personally identifiable information, like the plate number, would appear in our published dataset. Additionally, it is crucial to note that our dataset is intended solely for academic use and is not permitted for commercial purposes.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Consent Form for Collection of 3D Real Car Dataset

Purpose:
The purpose of this research is to create a 3D real car dataset for the development of various automotive and machine learning applications. The data will be used to improve technologies such as autonomous driving, vehicle detection, and augmented reality.

Procedures:
Participants will provide cars for data collection where cars will be scanned using 3D imaging technology. This collection will take place at designated locations and last approximately several minutes. Participants' vehicles will be scanned from multiple angles to create a comprehensive dataset.

Risks:
There are minimal risks associated with participation in this study. Participants will not be exposed to any hazardous conditions, and the 3D scanning technology is non-invasive.

Benefits:
Participants will contribute to the advancement of automotive technology, potentially leading to safer and more efficient vehicle systems. Participants will be paid at a rate of \$ 200 per hour.

Confidentiality:
Identifying information will be removed from the dataset. Only authorized personnel will have access to the raw data. We will hide any private information, like the plate number.

Consent:
I, the undersigned, consent to participate in the recording of a 3D real car dataset. I understand the purpose of the study, the procedures involved, the risks, and the benefits. I understand that my participation is voluntary and that I can withdraw at any time without penalty.

Date: _____
Name of Participant: _____
Signature of Participant: _____
Name of Researcher: _____
Signature of Researcher: _____

Figure 14: Consent Form for Collection of 3DRealCar