# MV-DeepSDF: Implicit Modeling with Multi-Sweep Point Clouds for 3D Vehicle Reconstruction in Autonomous Driving

Yibo Liu<sup>1,2</sup>, Kelly Zhu<sup>1,3</sup>, Guile Wu<sup>1</sup>, Yuan Ren<sup>1</sup>, Bingbing Liu<sup>1</sup>, Yang Liu<sup>1</sup>, Jinjun Shan<sup>2</sup> <sup>1</sup>Huawei Noah's Ark Lab, <sup>2</sup>York University, <sup>3</sup>University of Toronto

{yorklyb,jjshan}@yorku.ca, kellyk.zhu@mail.utoronto.ca
{guile.wu, yuan.ren3, liu.bingbing, yang.liu9}@huawei.com

### Abstract

Reconstructing 3D vehicles from noisy and sparse partial point clouds is of great significance to autonomous driving. Most existing 3D reconstruction methods cannot be directly applied to this problem because they are elaborately designed to deal with dense inputs with trivial noise. In this work, we propose a novel framework, dubbed MV-DeepSDF, which estimates the optimal Signed Distance Function (SDF) shape representation from multisweep point clouds to reconstruct vehicles in the wild. Although there have been some SDF-based implicit modeling methods, they only focus on single-view-based reconstruction, resulting in low fidelity. In contrast, we first analyze multi-sweep consistency and complementarity in the latent feature space and propose to transform the implicit space shape estimation problem into an element-to-set feature extraction problem. Then, we devise a new architecture to extract individual element-level representations and aggregate them to generate a set-level predicted latent code. This set-level latent code is an expression of the optimal 3D shape in the implicit space, and can be subsequently decoded to a continuous SDF of the vehicle. In this way, our approach learns consistent and complementary information among multi-sweeps for 3D vehicle reconstruction. We conduct thorough experiments on two real-world autonomous driving datasets (Waymo and KITTI) to demonstrate the superiority of our approach over state-of-the-art alternative methods both qualitatively and quantitatively.

# 1. Introduction

3D vehicle reconstruction from sparse and partial point clouds is a fundamental need in the autonomous driving industry [12, 18]. It aims to infer the 3D structure of vehicles at arbitrary resolutions in the wild, which is of great significance to many downstream tasks in autonomous driving. Despite the many 3D reconstruction methods



Figure 1. An illustration of the motivation of the proposed approach. Our approach takes multi-sweep point clouds as input and simplifies 3D vehicle reconstruction from multi-sweeps into an element-to-set feature extraction problem. In this way, we infer an optimal estimation of the 3D shape described in the abstract implicit space for 3D vehicle reconstruction in autonomous driving.

[29, 36, 27, 40, 42, 1, 20, 4], most of them focus on imagebased inputs [29, 36, 27, 40, 42], dense point cloud inputs [20, 32, 11], or their combination [1, 4, 9] and thus, cannot be directly applied to this problem.

Recently, [12] has shown promising performance in applying implicit modeling to tackle this problem. Contrary to explicit modeling methods [33, 34, 48, 50, 44, 43, 10, 2, 39] that directly represent 3D object shape structure with points, voxels, or meshes, implicit modeling maps the 3D shape to a low-dimensional latent space and learns the projection from the latent space to a continuous function that describes the 3D shape. This presents two significant advantages: 1) the 3D shape can be stored as a low-dimensional memory-saving latent code [32, 11, 12, 6]; 2) the trained network outputs a continuous function in the 3D space which supports mesh extraction at any resolution [32]. To this end, we also choose to employ implicit modeling in our approach.

However, previous point cloud-based implicit modeling methods [32, 11, 12] mainly focus on recovering 3D shapes from a single-view partial point cloud and fail to leverage the multi-sweep information of vehicles. As a result, when noise or annotation errors exist in an individual sweep, single-view methods usually produce low fidelity results. In real-world datasets (*e.g.*, [37, 16]), multi-sweep



Figure 2. The framework of the proposed MV-DeepSDF. Farthest Point Sampling (FPS) [34] is applied to the raw point clouds for preprocessing and DeepSDF [32] is employed to generate a latent code for each observation. We then extract global features from the standardized point clouds (refer to yellow block) and concatenate the global features with the latent codes as *element-level representations*. Next, the element-level representations are transformed into a *set-level predicted latent code* (refer to red block). Finally, we employ a pre-trained DeepSDF decoder [32] to project the predicted latent code to the SDF of the vehicle in 3D space and recover the 3D mesh. Note that the decoder used in the final step (refer to white trapezoid) is identical to that of DeepSDF (refer to green block).

point clouds are usually available and contain richer shape information because of the various viewing angles offered by multiple observations [18]. Although there are some existing 3D reconstruction methods from multi-view point clouds [18, 20, 3, 35], implicit modeling with multi-sweep point clouds remains an unsolved problem.

In this work, we propose a novel framework, dubbed MV-DeepSDF, to exploit multi-sweep point clouds and implicitly generate high-fidelity 3D reconstructions of vehicles for autonomous driving. An illustration of the motivation for our proposed approach is depicted in Figure 1. Specifically, to shed light on the importance of exploring multi-sweep point clouds for 3D vehicle reconstruction, we first analyze multi-sweep consistency and complementarity in the latent feature space. We then propose to consider the problem of shape estimation in the implicit space as an element-to-set feature extraction problem, where a set represents a collection of multi-sweeps and an element represents a sparse and partial point cloud. Next, we devise a new architecture to simultaneously extract a global feature and latent code for each element in the multi-sweep and concatenate them as element-level representations. These element-level representations are then transformed into a set-level predicted latent code with the aid of average pooling and mapping. This set-level latent code is an optimal estimation of the 3D shape described in the abstract implicit space, which we subsequently decode to a continuous SDF of the vehicle with a pre-trained DeepSDF decoder [32] and recover the 3D mesh from the SDF. The contributions of this work are threefold:

· We analyze multi-sweep consistency and complemen-

tarity in the latent feature space and transform the problem of shape estimation in the implicit space into an element-to-set feature extraction problem. This simplifies 3D vehicle reconstruction from multi-sweep point clouds into an optimal estimation of the 3D shape described in the abstract implicit space.

- We propose a novel MV-DeepSDF framework (see Figure 2) for implicit modeling with multi-sweep point clouds. A new architecture is constructed to extract the element-level representations and generate the set-level predicted latent code.
- We qualitatively and quantitatively demonstrate the superior performance of our approach over the state-ofthe-art alternative methods through extensive experiments on real-world datasets [37, 16].

# 2. Related Work

# 2.1. Explicit Modeling-Based 3D Reconstruction

Conventional 3D explicit modeling can be classified into three categories, namely point-based [33, 34, 48, 50, 44], voxel-based [43, 10], and mesh-based [2, 39]. Point-based methods, such as Point Completion Network (PCN) [50], GRNet [44], and PoinTr [48], output a point cloud with limited resolution since the number of points is fixed and not suitable for generating watertight surfaces due to lack of topology description. Voxel-based methods [43, 10] describe volumes by subdividing the space into a 3D grid, but they are memory and compute intensive, leading to slow training and low resolution representations [32]. Meshbased methods only generate meshes with simple topologies [41] given by a fixed reference template from the same object class and cannot guarantee watertight surfaces [23, 17]. In contrast, we propose implicit modeling with multi-sweep point clouds in a novel MV-DeepSDF framework, which estimates the optimal SDF shape representation for reconstructing 3D structures of vehicles in autonomous driving.

#### 2.2. Implicit Modeling-Based 3D Reconstruction

Implicit modeling [28, 32, 11, 12] uses function-based decision boundaries to implicitly define surfaces for 3D representations. DeepSDF [32] is a classical approach to implicit modeling. It estimates the Signed Distance Function (SDF) of an object from a partial point cloud, where the SDF specifies whether a querying position is inside or outside the surface of the object and the distance of the querying position from the surface. There have been some improved versions [11, 12, 6, 46] of DeepSDF, but these methods only focus on 3D reconstruction from a singleview point cloud and thus, generally produce low fidelity results when noise or annotation errors exist in an individual sweep. In contrast, we propose a novel framework for implicit modeling with multi-sweep point clouds. We analyze multi-sweep consistency and complementarity in the latent feature space and propose to resolve this problem by an optimal estimation of the 3D shape described in the abstract implicit space. Moreover, other implicit modeling-based methods, such as Neural Radiance Fields (NeRF) [29], Point-NeRF [45], Direct Voxel Grid Optimization [36], NeuS [42], NEAT [27], and [1, 4, 9] exist, but they all require image input and cannot deal with only LiDAR data.

#### 2.3. Multi-View 3D Reconstruction

3D reconstruction from multi-view point clouds can be categorized into conventional approaches [3, 35, 30] and deep learning-based approaches [18, 20, 5]. Conventional geometric approaches, such as Iterative Closest Point (ICP), TSDF [30], and [3, 35], only leverage multi-view observations geometrically, so the resulting completion is merely an aligned stack of partial point clouds and cannot produce watertight shapes unless the sensor fully loops around the object. In deep learning-based approaches, [18] proposes a weakly-supervised framework to directly learn 3D shape completion from multi-view LiDAR sweeps, but its reconstruction result is not watertight. [20] presents a shape completion framework using a multi-view depth map-based shape representation approach, while [5] designs a network to use multiple partial point clouds encoded into the latent space to estimate the optimal shape of the object. However, the output of these networks [20, 5] are point-based representations and in contrast with continuous field functions such as SDFs [32], lack continuity in 3D space. Moreover, some methods, such as [13, 15, 47], exist for reconstructing vehicle shapes from multiple observations, but [13] relies on stereo images and motion priors to regularize the shape estimation of vehicles, [15] requires multi-view camera images to build semantic maps which contain 3D vehicle shapes, and [47] requires inputs to follow a time sequence to track and reconstruct 3D objects. In contrast, our approach simplifies 3D vehicle reconstruction from multisweep point clouds into an optimal estimation of the 3D shape described in the abstract implicit space, which supports mesh extraction at any resolution and can deal with noise and annotation errors that exist in an individual sweep.

# 3. Methodology

#### 3.1. Preliminaries

In this section, we introduce preliminaries regarding DeepSDF [32] since our approach utilizes its framework for extracting the latent code and decoding the continuous SDF for 3D mesh generation. Formally, when adopting the DeepSDF decoder into the shape completion task, points of partial point clouds are taken as given surface points and the SDF value is defined as the distance from the querying point to the surface of the object. Interior and exterior querying points are sampled along the normals of the surface points and the SDF value is computed for each querying point. Suppose the set of newly sampled interior and exterior points is denoted as  $\mathcal{P}_{\mathcal{M}}$ , where M is the number of points. Define the *m*th point in  $\mathcal{P}_{\mathcal{M}}$  as  $\boldsymbol{x}_m$  and its corresponding SDF value as  $s_m$ . In the domain of DeepSDF [32], a 3D shape can be represented by a latent code  $oldsymbol{z} \in \mathbb{R}^{256}$ in the latent space. The complete shape is obtained by optimizing z via Maximum a Posterior estimation:

$$\hat{\boldsymbol{z}} = \frac{1}{\sigma^2} \|\boldsymbol{z}\|_2^2 + \arg\min_{\boldsymbol{z}} \sum_{m=1}^M \mathcal{L}_1\left(f_\theta\left(\boldsymbol{z}, \boldsymbol{x}_m\right), s_m\right), \quad (1)$$

where  $\mathcal{L}_1$  represents the clamped  $\mathcal{L}_1$  distance and  $\sigma$  is a regularization hyperparameter. Once the optimal latent code  $\hat{z}$  is determined, the complete 3D shape can be recovered by finding the zero-SDF-value isosurface through the DeepSDF decoder  $f(\cdot)$ , which is defined as:

$$\hat{s} = f_{\theta} \left( \hat{\boldsymbol{z}}, \boldsymbol{x} \right), \tag{2}$$

where  $\theta$  represents the parameters of the decoder,  $x \in \mathbb{R}^3$  denotes the 3D coordinates of the querying point, and  $\hat{s}$  denotes the estimated SDF value given by  $f(\cdot)$ . Here, the sign of  $\hat{s}$ , either positive or negative, indicates whether the point lies on the exterior or interior of the surface, respectively. Thus, the surface of the object can be implicitly represented by the isosurface composed of points with zero SDF values.

Despite its simplicity and efficiency, DeepSDF only takes the given surface points into consideration when generating 3D shapes (*i.e.*, Eq. (1)). This can lead to suboptimal reconstruction results in areas not captured by the partial point clouds since DeepSDF will reconstruct arbitrarily based on prior knowledge learned during training. For instance, referencing the ground truth of Case 1 in Figure 3, the reconstruction results of DeepSDF for Cases 2, 3, and 4 only demonstrate high fidelity in areas with well-captured surface points. The missing areas are reconstructed arbitrarily. To resolve this problem, we analyze multi-sweep consistency and complementarity in the latent feature space and propose implicit modeling with multi-sweep point clouds in MV-DeepSDF.

#### **3.2.** Approach Overview

As shown in Figure 2, there are three main steps in MV-DeepSDF, namely preprocessing, optimal latent code prediction, and 3D mesh extraction.

1) **Preprocessing.** First, we carry out Farthest Point Sampling (FPS) [34] to sample a fixed number of points on each raw partial point cloud since our global feature extractor requires standardized point clouds as input. For FPS, the number of centroids is set as 256.

2) Optimal Latent Code Prediction. The post-FPS point clouds are used as input into the global feature extractor (yellow block in Figure 2) to extract global features for each individual point cloud. Meanwhile, the original point clouds are used to extract latent codes for each partial point cloud through the pre-trained DeepSDF (green block in Figure 2). Then, as shown in the red block of Figure 2, the global features and latent codes are concatenated as element-level representations, followed by an average pooling operation to aggregate the information into a single instance tensor. Finally, this instance tensor is transformed into the predicted latent code by a fully-connected layer as the set-level representation.

**3) 3D Mesh Extraction.** The predicted latent code is converted to an SDF using a pre-trained DeepSDF decoder [32]. The isosurface composed of querying points with zero SDF values represents the surface of the instance and the implicit surface can be rasterized to a 3D mesh using Marching Cubes [25]. Since the SDF is a continuous function, it supports 3D mesh extraction at any resolution.

#### 3.3. Consistency and Complementarity Analysis

**Background.** In literature [38, 51], multi-view consistency refers to the consistency of describing the same instance from different viewpoint observations and multiview complementarity refers to the complementary information of the same instance provided by different observations. Existing methods mainly explore multi-view consistency and complementarity for 3D reconstruction with par-



Figure 3. The 3D reconstruction results of DeepSDF [32] (row 2) using different point clouds (row 1). All point clouds are sampled from the same CAD model of ShapeNetV2 [7]. Case 1 shows the complete dense point cloud used to train the DeepSDF decoder, while Cases 2, 3, and 4 show the partial point clouds.

tial point clouds in explicit modeling [35, 3].

However, existing methods are not applicable for exploring multi-sweep/-view consistency and complementarity in implicit modeling because the abstract feature space operates differently than the explicit space. In implicit modeling, point clouds are expressed as low-dimensional feature vectors (a.k.a. latent vectors) in the latent feature space and the latent vectors are generated through a non-linear learning process that projects the 3D space to the feature space through a deep neural network [32, 33, 34, 50].

**Theoretical Analysis.** Despite these differences, we can still explore multi-sweep consistency and complementarity in the latent feature space for implicit modeling. Inspired by [26], consider a total of *B* observations. The latent code for the *i*th observation, denoted by  $z_i$ , can be represented as:

$$\boldsymbol{z}_i = \boldsymbol{z}_{i,c} + \boldsymbol{z}_{i,s} + \boldsymbol{e}_i, \tag{3}$$

where  $i = \{1, \dots, B\}$ ,  $z_{i,c}$  is the consistent component,  $z_{i,s}$  is the specific component, and  $e_i$  is the error component. In this work,  $z_{i,c}$  and  $z_{i,s}$  can be defined as follows:

$$\boldsymbol{z}_{i,c} = \boldsymbol{z}_i \cap \boldsymbol{z}_{gt} \quad ext{and} \quad (\boldsymbol{z}_{i,s} + \boldsymbol{e}_i) = \boldsymbol{z}_{gt} - \boldsymbol{z}_i, \quad (4)$$

where  $z_{gt}$  is the latent code corresponding to the ground truth shape and ' $\cap$ ' and '-' represent the similarity and difference in information captured by the two feature vectors, respectively. Now, consider any two consistent components  $z_{i,c}$  and  $z_{j,c}$ . In the first case where  $z_{i,c} \cap z_{j,c} = \emptyset$ , the two components contain completely different consistent information, which is complementary. In the second case where  $z_{i,c} \cap z_{j,c} \neq \emptyset$ , but  $z_{i,c} \not\subseteq z_{j,c}$  and  $z_{j,c} \not\subseteq z_{i,c}$ , the two components contain some different information, which can be aggregated as complementary information to approach the ground truth. In the final case where  $z_{i,c} \cap z_{j,c} \neq \emptyset$ , but  $z_{i,c} \subseteq z_{j,c}$  or  $z_{j,c} \subseteq z_{i,c}$ , either  $z_{i,c}$  or  $z_{j,c}$  contains redundant information from the other component. While redundant, this will not adversely impact complementary aggregation. Thus, to make the optimal  $\hat{z}$  approach the ground truth  $z_{gt}$ , it is desired to aggregate all consistent and complementary information among all individual features:

$$\hat{\boldsymbol{z}} = \boldsymbol{z}_{1,c} \cup \boldsymbol{z}_{2,c} \cdots \cup \boldsymbol{z}_{B,c} + \boldsymbol{p}, \tag{5}$$

where ' $\cup$ ' refers to the complementary aggregation of information among features and p denotes the predicted information for regions not captured in the multi-sweeps (*e.g.*, the side of the vehicle not captured by the multi-sweeps as shown in the Cases 3 and 4 of Figure 3). In summary, with implicit modeling, we can learn the optimal latent code  $\hat{z}$ by aggregating consistent (Eqs. (3)-(4)) and complementary (Eq. (5)) information among multi-sweeps.

# 3.4. Architecture Design of MV-DeepSDF

To take advantage of multi-sweep consistency and complementarity in the latent feature space, we design a new architecture for implicit modeling. Our proposed architecture uses multi-sweep point clouds as model input and generates an optimal estimation of the 3D shape in the implicit space (a.k.a. the optimal latent code) as the output. The functionality of the network can be formulated as:

$$g_{\alpha}\left(\boldsymbol{Z},\boldsymbol{P}\right) = \hat{\boldsymbol{z}} \to \boldsymbol{z}_{gt},\tag{6}$$

where  $g(\cdot)$  denotes the desired network with associated parameters  $\alpha$ ,  $\mathbf{Z} = \{\mathbf{z}_1, \cdots, \mathbf{z}_B\}$  denotes the set of latent codes, B represents the number of observations or latent codes,  $\mathbf{P} = \{\mathcal{P}_1, \cdots, \mathcal{P}_B\}$  denotes the set of multi-view point clouds, and  $\hat{\mathbf{z}} \in \mathbb{R}^{256}$  and  $\mathbf{z}_{gt} \in \mathbb{R}^{256}$  represent the estimated optimal and ground truth latent codes, respectively. Hence, the goal is converted into learning a network to make  $\hat{\mathbf{z}}$  approach  $\mathbf{z}_{gt}$  for subsequent reconstructions. To this end, we devise a new architecture to realize the functionality required by simultaneously learning global features and latent codes for generating predicted latent codes.

Specifically, inspired by PointNet [33], a pioneering work for 3D classification and segmentation, we transform the problem of shape estimation in the implicit space into an element-to-set feature extraction problem. Given a point cloud containing a set of 3D points, where all 3D points consistently describe a single object since they all belong to the same instance, we can learn a global feature for each point cloud that aggregates complementary features from each individual 3D point. However, only learning a global feature cannot be used for 3D reconstruction, so we propose to abstract the pipeline of PointNet (see Figure 4(a)) into a generalized process (see Figure 4(b)). As shown in Figure 4(a), PointNet first computes 3D point features using a shared MLP on each 3D point and then transforms these point features into a single global feature through max pooling. Now, suppose there exists a set composed of multiple elements. The goal of the generalized pipeline is to extract a feature for the entire set by aggregating all elements. To this



Figure 4. Comparison of three pipelines. (a) The pipeline of Point-Net [33]. (b) The pipeline of an abstract generalized process to extract the set feature from the elements. (c) The proposed pipeline to extract the instance feature from multi-sweep point clouds.

end, we propose two generalized operations. As shown in Figure 4(b), we define operation 1 as the computation of the element features (*e.g.*, the shared MLP in PointNet) and operation 2 as the transformation of the element features into a set feature (*e.g.*, the max pooling operation in PointNet).

Using this generalized process, we can now devise our architecture to simultaneously learn global features and latent codes for generating predicted latent codes. First, we construct a global feature extractor (see yellow block in Figure 2) for learning a global feature for each point cloud that aggregates complementary features from each individual 3D point. The architecture of this feature extractor is inspired by the PCN encoder [50], a variant of PointNet [33]. This block is a stack of four PointNet [33] encoders with 128, 256, 512, and 1024 units, respectively. Since this architecture is inspired by the PCN encoder [50] and its global feature extraction operation is applied to every individual point cloud, we name it shared PCN. Meanwhile, to exploit the implicit information of each element, we employ a pre-trained DeepSDF [32] model (as denoted by Eq. (1)) to generate a latent code for each partial point cloud. Next, we concatenate the global features and latent codes to generate the element-level representations. Then, we employ average pooling and mapping to aggregate the element-level representations into a set-level predicted latent code. As shown in Figure 4(c), operation 1 of our pipeline consists of a global feature extractor and a pre-trained DeepSDF model, while operation 2 consists of concatenation, average pooling, and mapping. This pipeline realizes a symmetric operation on unordered multi-sweep point clouds.

#### 3.5. Model Training of MV-DeepSDF

We adopt a curriculum learning strategy for model training, which consists of two stages.

**Stage One.** In stage one, we pre-train the DeepSDF decoder using watertight CAD models belonging to the car taxonomy of the ShapeNetV2 dataset [7]. The latent code z

and parameters of the decoder  $\theta$  are jointly optimized as:

$$\underset{\theta, \{\boldsymbol{z}_{j}\}_{j=1}^{J}}{\arg\min} \sum_{j=1}^{J} \left( \frac{1}{\sigma^{2}} \|\boldsymbol{z}_{j}\|_{2}^{2} + \sum_{k=1}^{K} \mathcal{L}_{1}\left(f_{\theta}\left(\boldsymbol{z}_{j}, \boldsymbol{x}_{k}\right), s_{k}\right) \right),$$
(7)

where J represents the number of 3D shapes used for training and K represents the number of points for each shape. During this stage, the decoder gains prior knowledge of vehicle shapes and once trained, the decoder remains fixed for the latent code generation step of our pipeline. Please refer to [32] for the details of the training at this step.

Training Dataset Preparation for Stage Two. As shown in Eq. (6), ground truth latent codes, partial point clouds, and their corresponding latent codes are required for training the reconstruction network. However, real-world datasets, such as Waymo [37] and KITTI [16], do not contain ground truth shapes, so it is necessary to use a synthetic dataset [7] for generating the training data. To resolve this problem, we utilize the latent codes of the training shapes (e.g., the complete dense point cloud shown in Case 1 of Figure 3) as ground truth latent codes in the second stage of model training. As for the partial point clouds, the domain gap between training instances and in-the-wild instances directly determines the ability of our network to generalize to real-world datasets. To reduce the domain gap, we adopt PCGen [24] as our partial point cloud generation method. PCGen places a virtual LiDAR with real-world parameters (resolution and sampling pattern) around the vehicle and simulates the point cloud captured by the virtual LiDAR. This differs from the method used by DeepSDF [32], where a simulated depth camera is used as the virtual sensor instead of a LiDAR. A visual comparison of the real-world point cloud and simulated partial point clouds obtained using various methods is given in Figure 5. As shown in Figure 5, the partial point clouds of PCGen [24] are visually closer to the in-the-wild LiDAR sweeps of Waymo [37] and KITTI [16] than those of the virtual depth camera approach [32], which only raycasts evenly in the inclination and azimuth directions. Apart from this visual superiority, the advantages of generating partial point clouds with PCGen can also be observed through the improved generalization ability of our network. See experiments for further details.

When generating training data from ShapeNetV2 [7], we randomly sample one side of the vehicle to accurately reflect the real-world scan of a LiDAR on an ego-vehicle. In this way, we produce six partial point clouds for each vehicle. The virtual LiDAR pose for each partial point cloud is generated using a set of restrictions. These restrictions are set by considering the possible relative poses of an ego vehicle's LiDAR in the coordinate system of the other vehicle. In particular,  $\theta \in [0^\circ, 180^\circ]$  or  $\theta \in [-180^\circ, 0^\circ]$ ,  $r \in [3, 15]$ , and  $h \in [0.8, 1.2]$  are used, where  $\theta$  represents the azimuth,



Figure 5. Visual comparison of in-the-wild partial point clouds and raycasted point clouds. (a) A partial point cloud from the Waymo tracking dataset [37]. (b) The raycasted point cloud generated by PCGen [24] using the LiDAR parameters of Waymo. (c) A partial point cloud from the KITTI tracking dataset [16]. (d) The raycasted point cloud generated by PCGen [24] using the LiDAR parameters of KITTI. (e) The raycasted point cloud obtained from the virtual depth camera approach used by DeepSDF [32].

r represents the distance between the LiDAR and the vehicle, and h represents the height of the LiDAR to the ground.  $\theta$ , r, and h describes the pose of the virtual LiDAR with respect to the frame of the other vehicle. Both r and h are expressed in the normalized space, where the size of the vehicle is normalized into the range [-1, 1]. Finally, we generate a latent code for each partial point cloud using Eq. (1).

**Stage Two.** In stage two, the network is trained to reconstruct 3D vehicles from multi-view partial point clouds and their corresponding latent codes. Consider the *c*th vehicle instance. Let  $P_{B,c}$  represent the set of partial point clouds,  $Z_{B,c}$  represent the set of latent codes corresponding to the multi-view point clouds, and  $z_{gt,c}$  represent the ground truth latent code acquired from stage one. Furthermore, let  $g(\cdot)$  represent the function of the implicit shape prediction network (refer to yellow and red blocks in Figure 2) and  $\alpha$  represent the parameters of the model. The training objective of stage two can be defined as:

$$\arg\min_{\alpha} \sum_{c=1}^{C} \mathcal{L}_2\left(g_{\alpha}(\boldsymbol{Z}_{B,c}, \boldsymbol{P}_{B,c}), \boldsymbol{z}_{gt,c}\right), \qquad (8)$$

where C represents the number of instances employed for training and  $\mathcal{L}_2$  is the Mean Squared Error loss. The training details are presented in Sec. 4.1.

# 4. Experiments

In this section, we present qualitative and quantitative results on two real-world autonomous driving datasets, namely Waymo [37] and KITTI [16]. The multi-sweep point clouds for Waymo are collected from 136 unique vehicle instances of Waymo Open Dataset's tracking data [37], while those for KITTI are extracted from 233 vehicle in-



Figure 6. Visual comparison with the state-of-the-art methods (DeepSDF [32], C-DeepSDF [11], MendNet [12], and AdaPoinTr [49]) on the Waymo [37] dataset. DeepSDF+MS and MendNet+MS indicate the models with multi-sweep input.

stances of KITTI's tracking dataset [16]. Due to page limitations, experiment details and some results (*e.g.*, the results on the synthetic dataset ShapeNetV2 [7]) are discussed in the supplementary material.

#### 4.1. Implementation Details

Architecture Details. Our global feature extractor, named *shared PCN*, is a variant of the PCN encoder [50]. The embedded shared MLP layers are identical to those of PointNet [33], which are implemented using 1D convolution layers. Moreover, since DeepSDF [32] normalizes the latent codes into the range [-1, 1], we add a tanh layer after the last shared MLP to normalize the global features into the same range as the latent codes. The decoder is identical to that of DeepSDF [32]. We report results based on an implementation with Python and PyTorch, but our approach also supports implementation with MindSpore [21].

Training Details. In the first stage of training, we follow the same method as presented in DeepSDF [32] to train the decoder using watertight CAD models from the car taxonomy of ShapeNetV2 [7]. Once trained, the decoder is fixed and projects a 256-dimensional latent code to an SDF in 3D space. In the second stage, we train the model for 20 epochs using the Adam optimizer [22] with a learning rate of 1e-5. The batch size is set to 1, allowing the model to simultaneously see all 6 frames of the given instance at each iteration. Since both the output of the model and the supervision are 1D vectors, the loss computation is very fast. It only takes 10 minutes to train our network on a single NVIDIA GeForce RTX 2080 GPU during stage two. However, the DeepSDF decoder [32] used to prepare the training dataset for stage two through latent code generation is a more timeconsuming process.

#### 4.2. Evaluation Protocol

**Metrics.** Since an in-the-wild instance does not have a ground truth 3D shape, we follow [12] to evaluate the recon-

struction results. In particular, the multiple LiDAR sweeps are stacked to construct the ground truth points set X. Moreover, we sample 30,000 points on the surface of each reconstructed mesh to generate the reconstruction points set Y. We employ Asymmetric Chamfer Distance (ACD) [12], which is the sum of the squared distance of each ground truth point to the nearest point in the reconstructed point set, to evaluate reconstruction results on real-world datasets:

$$\operatorname{ACD}\left(\boldsymbol{X},\boldsymbol{Y}\right) = \sum_{\boldsymbol{x}\in\boldsymbol{X}} \min_{\boldsymbol{y}\in\boldsymbol{Y}} \|\boldsymbol{x}-\boldsymbol{y}\|_{2}^{2}.$$
 (9)

For ACD, a smaller value is preferred. In addition, we compute the recall of the ground truth points from the reconstructed shape, which is defined as:

$$\operatorname{Recall}(\boldsymbol{X}, \boldsymbol{Y}) = \frac{1}{|\boldsymbol{X}|} \sum_{x \in \boldsymbol{X}} \left[ \min_{y \in \boldsymbol{Y}} \|x - y\|_2^2 <= t \right], \quad (10)$$

where the threshold t is set as 0.1, following [12].

**Competitors.** We compare our approach with four stateof-the-art methods, including three implicit modeling-based methods (DeepSDF [32], C-DeepSDF [11], and MendNet [12]) and AdaPoinTr [49], which is built upon PoinTr++ [48], the winner of Multi-View Partial Point Cloud Challenge 2021 on Completion and Registration [31]. Training data for DeepSDF and C-DeepSDF are generated following the original papers, while MendNet uses PCGen [24] to generate training data. The pipeline of AdaPoinTr [49] requires online generation of partial point clouds using the virtual depth camera approach and since it is not trivial to replace this process with PCGen, AdaPoinTr is trained using the default setup without PCGen. Note that with multi-sweep point clouds, competitors generate multiple reconstruction results (one for each partial point cloud), so we compare with the best single-shot reconstruction result, which is the mesh with the minimum ACD among the multiple single-shot reconstructed meshes. Furthermore, we report the results of DeepSDF and MendNet with multi-sweep

Method \ Metric	$ACD_{mean} \downarrow$	$ACD_{median} \downarrow$	Recall ↑
DeepSDF [32]	6.26	5.81	93.51
DeepSDF+MS	5.12	5.09	95.57
C-DeepSDF [11]	6.21	5.64	93.98
MendNet [12]	4.92	4.79	95.39
MendNet+MS	4.85	4.77	95.76
AdaPoinTr*[49]	4.79	4.74	95.95
Ours-VDC	4.76	4.55	96.05
Ours	3.36	2.26	96.84

Table 1. Comparison with the state-of-the-art methods on Waymo. ACD  $\downarrow$  is multiplied by 10<sup>3</sup>. Recall  $\uparrow$  is presented as a percentage. \*AdaPoinTr is trained with the default setup without PCGen on the large-scale ShapeNet-55 [48] following the original paper [49], and then fine-tuned on the car taxonomy of ShapeNetV2 in our experiments. Other methods are only trained with the car taxonomy of ShapeNetV2.

inputs (denoted as **DeepSDF+MS** and **MendNet+MS**), in which we stack the partial point clouds of multi-sweeps into a single point cloud and perform single-shot-based reconstruction on this stacked point cloud. For our approach, in addition to the default setup which generates the training dataset using PCGen (denoted as **Ours**), we also adapt our approach to use training data generated by the virtual depth camera approach [32] (denoted as **Ours-VDC**).

# 4.3. Results on Waymo

The qualitative and quantitative comparisons of our approach with the state-of-the-art methods are presented in Figure 6 and Table 1, respectively. Overall, our approach achieves the best reconstruction results both qualitatively and quantitatively. Specifically, from Figure 6 and Table 1, we can see that DeepSDF [32] and C-DeepSDF [11] are sensitive to noise, resulting in low fidelity. MendNet [12] generates more stable results than DeepSDF and C-DeepSDF, but its reconstruction meshes contain some irregular holes on the surface (see sixth column of Figure 6). Although AdaPoinTr [49] outputs point clouds with good fidelity, it expresses the shape with a limited resolution and thus fails to describe continuous local details (see seventh column of Figure 6).

Compared with vanilla DeepSDF and MendNet, DeepSDF+MS and MendNet+MS yield messier surfaces/structures (see fourth and eighth columns of Figure 6), despite their better metric scores. This shows that even though transforming all partial point clouds into a single frame does benefit information aggregation, it also accumulates noise and annotation errors in the process. Hence, due to noise in real-world data, it is not practical to geometrically stack multi-sweep point clouds and directly perform single-shot-based reconstruction. In comparison, our approach generates smooth watertight shapes (see last two columns of Figure 6) with high fidelity (see Table 1). Furthermore, Ours outperforms Ours-VDC, which indicates

Method \ Metric	$ACD_{mean} \downarrow$	$ACD_{median} \downarrow$	Recall ↑
DeepSDF [32]	6.81	6.17	80.65
DeepSDF+MS	6.11	5.83	82.73
C-DeepSDF [11]	6.73	5.99	80.77
MendNet [12]	5.94	5.64	83.84
MendNet+MS	5.83	5.61	84.26
AdaPoinTr [49]	5.89	5.67	84.20
Ours-VDC	5.75	5.24	84.39
Ours	4.27	3.01	85.88

Table 2. Comparison with the state-of-the-art methods on H	KITTI.
ACD $\downarrow$ is multiplied by 10 <sup>3</sup> . Recall $\uparrow$ is presented as a perce	entage.

that using PCGen to prepare training partial point clouds (refer to Figure 5) reduces the domain gap between training and in-the-wild instances.

#### 4.4. Results on KITTI

The qualitative and quantitative comparisons of our approach against the state-of-the-art methods are presented in Figure 7 and Table 2, respectively. From these results, we can also observe the superior performance of our approach over the state-of-the-art competitors. Specifically, qualitatively, as shown in Figure 7, our approaches (Ours and Ours-VDC) are more robust to noise when compared to DeepSDF [32], DeepSDF+MS, and C-DeepSDF [11] and generate smoother, watertight surfaces compared to Mend-Net [12], AdaPoinTr [49], and MendNet+MS. Quantitatively, as shown in Table 2, Ours yields the best ACD and Recall results, while Ours-VDC performs the second best.

#### 4.5. Ablation Study

To verify the effectiveness of the main components of our approach, we carry out a series of experiments on Waymo [37] as shown in Table 3. The first row of Table 3 refers to the baseline model of our approach, which is a shared PCN encoder (Enc.) (yellow block in Figure 2) followed by an average pooling layer (Avg.). This architecture is similar to that of MendNet [12], which adds two stacked PointNet encoders [33] to the DeepSDF decoder. In the second row, we directly use bitwise multiplication to merge the global features (B $\times$ 256, instead of B $\times$ 1024) and the latent codes ( $B \times 256$ ). However, this results in obvious performance degradation, which shows that directly merging the global features and latent codes obtained from multi-sweep information does not bring improvement to our baseline model. In comparison, in the third row, we add our proposed components, namely the input of latent codes (Dep.), concatenation (Con.), and mapping (Map.), to the baseline model, but with max pooling. This brings significant performance improvement compared to the first and second rows. Then, in the last row, we further replace max pooling with average pooling, which includes all proposed components of this work. This yields the best results.



Figure 7. Visual comparison with the state-of-the-art methods (DeepSDF [32], C-DeepSDF [11], MendNet [12], and AdaPoinTr [49]) on the KITTI [16] dataset. DeepSDF+MS and MendNet+MS indicate the models with multi-sweep input.

Opera	ation 1	Operation 2		Metric		
Enc.	Dep.	Mer.	Pool	Map.	ACD $\downarrow$	Recall ↑
1			Avg.		4.89	95.50
1	1	Mul.	Avg.		7.33	84.93
1	1	Con.	Max.	1	4.35	96.61
1	1	Con.	Avg.	1	3.36	96.84

Table 3. The ablation study of the proposed framework. Refer to Figure 4(b) to see the definitions of operations 1 and 2. Enc.: the shared PCN encoder; Dep.: latent code generation through DeepSDF as input; Mer.: the operation to merge global features and latent codes; Map.: the mapping conducted by the fullyconnected layer; Con.: the concatenation; Mul.: the bitwise multiplication; Avg.: average pooling; Max.: max pooling. ACD  $\downarrow$  is multiplied by 10<sup>3</sup>. Recall  $\uparrow$  is presented as a percentage.

#### 4.6. Extension to Other Taxonomies

It is straightforward to extend the proposed shape completion framework to other taxonomies. Suppose that we want to extend MV-DeepSDF to perform shape completion on sofas. First, a new DeepSDF decoder [32] needs to be pre-trained using watertight CAD models from the sofa taxonomy of ShapeNetV2 [7] or another synthetic dataset. Then, PCGen [24] or another sampling technique (e.g., the approaches proposed in [32, 48]) can be used to generate partial point clouds for each sofa instance. These partial point clouds are then passed into the pre-trained DeepSDF decoder [32] to produce their corresponding latent codes. Finally, the partial point clouds and latent codes can be used together to train MV-DeepSDF. Following the pipeline above, we evaluate on indoor sofas of ShapeNetV2 and outdoor cyclists of Waymo and show that our method is versatile and extendible to other taxonomies (see Figure 8).

While existing datasets, such as ScanNet [8] and Semantic3D [19], provide real-world point cloud data for many taxonomies, the availability of tracking data for multisweep scans is still fairly limited to the autonomous driving industry. Hence, it is necessary to first obtain the corresponding multi-sweep tracking labels for the desired class



Figure 8. Visual comparison of DeepSDF and MV-DeepSDF on the indoor sofa of ShapeNetV2 and the outdoor cyclist of Waymo.

when choosing to adopt MV-DeepSDF on a new taxonomy.

#### 5. Conclusion

In this work, we propose a new MV-DeepSDF framework to facilitate implicit modeling with multi-sweep point clouds for autonomous driving. The main idea is to abstract 3D vehicle reconstruction from multi-sweeps into an element-to-set feature extraction problem. Namely, we consider the multi-sweeps of a vehicle as elements composing a set and infer the set feature, which is an optimal estimation of the 3D shape described in the abstract implicit space. In particular, we present a theoretical analysis of multi-sweep consistency and complementarity in the latent feature space. Guided by this analysis, we design a new architecture to optimally estimate the Signed Distance Function shape of a vehicle from its in-the-wild multi-sweep point clouds. Qualitative and quantitative evaluations on both real-world and synthetic datasets show the superiority of our approach over the state-of-the-art methods.

**Limitation.** Despite promising results, our approach still relies on a synthetic 3D dataset to gain prior knowledge of 3D shapes for reconstruction. Exploring reconstruction by learning directly from real data is worthy of further study.

# Acknowledgments

The authors gratefully acknowledge the support of MindSpore, CANN, and Ascend AI Processor. The authors would also like to thank Andrew Yang and Hunter Schofield for their assistance.

# **Supplementary Material**

# A. Overview

This material provides quantitative and qualitative experimental results, dataset and implementation details, and discussions that are supplementary to the main paper.

### **B.** Dataset Details

The multi-sweep LiDAR point clouds for Waymo are collected from 136 unique vehicle instances of Waymo Open Dataset's tracking data [37], while those for KITTI are extracted from 233 unique vehicle instances of KITTI's tracking dataset [16]. In total, we extracted 3943 partial point clouds from the 136 vehicle instances of Waymo and 4235 partial point clouds from the 233 vehicle instances of KITTI. These raw multi-sweep point clouds are directly used as model input to obtain experimental results on our model and the state-of-the-art methods [32, 11, 12, 49]. When performing inference, all partial point cloud frames from the given instance are simultaneously passed into our model as input.

To construct the ground truth stacked point cloud, we first aggregate all partial point clouds within a multi-sweep to generate a dense stacked point cloud. Since this stacked point cloud contains unwanted noise, such as ground plane points and points lying on the exterior or interior of the vehicle surface, we perform statistical outlier removal on the stacked point cloud, which computes the average distance of a point from its neighbours and removes all points lying farther away from their neighbours than average. Denoising is essential for dataset processing since the presence of noise in the ground truth shape can result in false positives and false negatives in model performance, whereby a messy shape generated by a model that fits to the noise of the ground truth is deemed high fidelity and a smooth shape generated by a noise-robust model is deemed low fidelity.

# C. Visualization Results

Due to the page limitation of the main paper, we present more visualization results of our model in Figure 9. Additionally, to evidently present the significant improvement of our method over the baseline vanilla DeepSDF [32], a visual comparison on Waymo [37] is given in Figure 10.

# **D.** Results on ShapeNetV2

We randomly preserve 300 vehicles from the car taxonomy of ShapeNetV2 [7] as the test dataset. The remainder is used to train our network in stage one. Each vehicle instance is comprised of 6 partial point clouds generated by PCGen [24] under Waymo's LiDAR parameters. Note that the following results are generated solely using PCGen [24] to sample partial point clouds and the use of other sampling techniques (*e.g.*, the approaches proposed in [32, 48]) would yield different results on the same ShapeNetV2 dataset [7].

**Metrics.** Since ground truth shapes are readily available in synthetic datasets such as ShapeNetV2 [7], we use Chamfer Distance (CD) [14] to evaluate the 3D reconstruction results. Following DeepSDF [32], we sample 30,000 points on the surface of both the ground truth and reconstructed mesh. Given two point sets, the CD is the sum of the squared distance of each point to the nearest point in the other point set:

$$CD(\mathbf{X}, \mathbf{Y}) = \sum_{x \in \mathbf{X}} \min_{y \in \mathbf{Y}} \|x - y\|_{2}^{2} + \sum_{y \in \mathbf{Y}} \min_{x \in \mathbf{X}} \|x - y\|_{2}^{2}.$$
(11)

As outlined in the main paper, we only compare the results of our model with the best single-shot reconstruction result, which is the mesh with the minimum CD among the multiple single-shot reconstructed meshes.

**Oualitative and Ouantitative Comparison.** The qualitative and quantitative comparison of our approach against the state-of-the-art methods (DeepSDF [32], C-DeepSDF [11], MendNet [12], and AdaPoinTr [49]) are presented in Figure 11 and Table 4, respectively. Note that when testing on ShapeNetV2 [7], since PCGen [24] is used to generate both the training and test dataset, we only present Ours, in contrast with the comparison of Ours with Ours-VDC in the main paper. The meshes generated by DeepSDF [32], C-DeepSDF [11], and MendNet [12] show high fidelity compared to their performance on real-world datasets, but still show inferior performance to Ours. AdaPoinTr [49] also produces shapes with decent fidelity, but the reconstructed result is not watertight and expresses the shape with a limited resolution which fails to describe the continuous surface of the vehicle.

Method \ Metric	$CD_{mean}\downarrow$	$CD_{median} \downarrow$
DeepSDF [32]	5.47	5.15
C-DeepSDF [11]	5.31	5.03
MendNet [12]	4.22	3.65
AdaPoinTr [49]	4.10	3.36
Ours	3.17	2.54

Table 4. Comparison of the proposed network with the state-of-the-art approaches on ShapeNetV2 [7]. CD is multiplied by  $10^3$ .

# E. Comparison to a Non-Learning Approach

We now present an alternative non-learning approach, computing the mean latent code, for the task of multi-sweep



Figure 9. Additional visualization results of MV-DeepSDF on the Waymo [37] and KITTI [16] datasets.



Figure 10. Visual comparison of MV-DeepSDF and DeepSDF [32] on Waymo [37]. Individual point clouds are given in the first row and their corresponding reconstruction results from vanilla DeepSDF in the second row.



Figure 11. Visual comparison with the state-of-the-art methods (DeepSDF [32], C-DeepSDF [11], MendNet [12], and AdaPoinTr [49]) on the ShapeNetV2 [7] dataset.

3D vehicle reconstruction. As introduced in [32], linear interpolation between two latent codes in the latent space can also generate meaningful shape representations. Moreover, averaging is a common method of linear interpolation used for reducing error among multi-observation data. To this end, we investigate the effect of computing the mean latent code from the single-shot-based latent codes of a given multi-sweep and using this mean latent code for mesh reconstruction. We present the case shown in Figure 12, where two single-shot partial point clouds, PC 1 and PC 2, are used to generate two latent codes,  $z_1$  and  $z_2$ , and meshes, Mesh 1 and Mesh 2, respectively. We define the mean latent code as  $z_{mean} = 0.5(z_1+z_2)$  and generate the corresponding mesh, denoted by Mean. As shown, Mean is simply a uniform fusion of Mesh 1 and Mesh 2. Moreover, Mean is inferior to Mesh 2, the best single-shot in this



Figure 12. Comparison of the result of our approach to that of the mean latent code on ShapeNetV2 [7]. The proposed network is not fine-tuned.

example, which is also inferior to Ours, the result of our proposed model.

Num of PCs	$ACD_{mean}\downarrow$	$ACD_{median}\downarrow$
3	3.47	2.44
6	3.36	2.26
9	3.32	2.21

Table 5. Ablation study on Waymo [37] using different numbers of point clouds per instance. ACD is multiplied by  $10^3$ .

### F. Effect of Number of Point Clouds

The number of frames corresponding to an individual vehicle instance in Waymo [37] and KITTI [16] ranges up to 240 partial point clouds per instance. However, the vast majority of instances only contain between 3 to 9 partial point clouds. In this section, we investigate the relationship between the number of partial point clouds provided for each instance during stage two of training and overall model performance. Table 5 presents the experimental results of providing different numbers of partial point clouds increases. However, since generating the latent code for each partial point cloud with DeepSDF [32] is a timely process (around 10 seconds), we choose 6 observations per instance as a trade-off between performance and efficiency.

# G. Effect of Number of Points Per Point Cloud

The number of points captured in a single frame of Waymo [37] and KITTI [16] mostly falls into a range of 300 to 1000 points. Thus, we set the number of points per point cloud as 256 in our framework for performing FPS. In this section, we investigate the relationship between the number of points per point cloud during inference and model performance. Table 6 presents the experimental results of varying the number of points per point cloud on both DeepSDF and our model with Waymo [37]. As shown, when the number of points decreases, the performance of DeepSDF drops dramatically whereas our method holds steady.

Num of Points	256		128	
Metric	$ACD_{mean} \downarrow$	$ACD_{median} \downarrow$	$ACD_{mean} \downarrow$	$ACD_{median} \downarrow$
DeepSDF	6.26	5.81	12.52	8.51
Ours	3.36	2.26	3.47	2.64
<b>T11 ( 111</b>		117 505	11.00	1

Table 6. Ablation study on Waymo [37] using different numbers of points per point cloud. ACD is multiplied by  $10^3$ .

# References

- Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6290–6301, 2022.
- [2] Heli Ben-Hamu, Haggai Maron, Itay Kezurer, Gal Avineri, and Yaron Lipman. Multi-chart generative surface modeling. *ACM Transactions on Graphics*, 37(6):1–15, 2018.
- [3] Paul J Besl and Neil D McKay. Method for registration of 3-d shapes. In Sensor fusion IV: control paradigms and data structures, volume 1611, pages 586–606. Spie, 1992.
- [4] Hongrui Cai, Wanquan Feng, Xuetao Feng, Yan Wang, and Juyong Zhang. Neural surface reconstruction of dynamic scenes with monocular rgb-d camera. In *Proc. of Conference* on Neural Information Processing Systems, 2022.
- [5] Yingjie Cai, Kwan-Yee Lin, Chao Zhang, Qiang Wang, Xiaogang Wang, and Hongsheng Li. Learning a structured latent space for unsupervised point cloud completion. In Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5543–5553, 2022.
- [6] Rohan Chabra, Jan E Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *Proc. of European Conference on Computer Vision*, pages 608–625. Springer, 2020.
- [7] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012, 2015.
- [8] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In Proc. of the IEEE conference on computer vision and pattern recognition, pages 5828–5839, 2017.
- [9] Angela Dai, Christian Diller, and Matthias Nießner. Sg-nn: Sparse generative neural networks for self-supervised scene completion of rgb-d scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 849–858, 2020.
- [10] Angela Dai, Charles Ruizhongtai Qi, and Matthias Nießner. Shape completion using 3d-encoder-predictor cnns and shape synthesis. In *Proc. of the IEEE conference on computer vision and pattern recognition*, pages 5868–5877, 2017.
- [11] Yueqi Duan, Haidong Zhu, He Wang, Li Yi, Ram Nevatia, and Leonidas J Guibas. Curriculum deepsdf. In Proc. of European Conference on Computer Vision, pages 51–67. Springer, 2020.

- [12] Shivam Duggal, Zihao Wang, Wei-Chiu Ma, Sivabalan Manivasagam, Justin Liang, Shenlong Wang, and Raquel Urtasun. Mending neural implicit modeling for 3d vehicle reconstruction in the wild. In Proc. of IEEE/CVF Winter Conference on Applications of Computer Vision, pages 1900– 1909, 2022.
- [13] Francis Engelmann, Jörg Stückler, and Bastian Leibe. Samp: shape and motion priors for 4d vehicle reconstruction. In *Proc. of IEEE Winter Conference on Applications of Computer Vision*, pages 400–408. IEEE, 2017.
- [14] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proc. of the IEEE conference on computer vision and pattern recognition*, pages 605–613, 2017.
- [15] Qiaojun Feng, Yue Meng, Mo Shan, and Nikolay Atanasov. Localization and mapping using instance-specific mesh models. In Proc. of IEEE/RSJ International Conference on Intelligent Robots and Systems, pages 4985–4991. IEEE, 2019.
- [16] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- [17] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In Proc. of the IEEE conference on computer vision and pattern recognition, pages 216–224, 2018.
- [18] Jiayuan Gu, Wei-Chiu Ma, Sivabalan Manivasagam, Wenyuan Zeng, Zihao Wang, Yuwen Xiong, Hao Su, and Raquel Urtasun. Weakly-supervised 3d shape completion in the wild. In *Proc. of European Conference on Computer Vi*sion, pages 283–299. Springer, 2020.
- [19] Timo Hackel, Nikolay Savinov, Lubor Ladicky, Jan D Wegner, Konrad Schindler, and Marc Pollefeys. Semantic3d. net: A new large-scale point cloud classification benchmark. arXiv preprint arXiv:1704.03847, 2017.
- [20] Tao Hu, Zhizhong Han, Abhinav Shrivastava, and Matthias Zwicker. Render4completion: Synthesizing multi-view depth maps for 3d shape completion. In Proc. of the IEEE/CVF International Conference on Computer Vision Workshops, pages 0–0, 2019.
- [21] Huawei. Mindspore. https://www.mindspore.cn/, 2020.
- [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Proc. of the International Conference on Learning Representations (Poster), 2015.
- [23] Chen Kong, Chen-Hsuan Lin, and Simon Lucey. Using locally corresponding cad models for dense 3d reconstructions from a single image. In *Proc. the IEEE conference on computer vision and pattern recognition*, pages 4857–4865, 2017.
- [24] Chenqi Li, Yuan Ren, and Bingbing Liu. Pcgen: Point cloud generator for lidar simulation. In Proc. of International Conference on Robotics and Automation, 2023.
- [25] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. ACM siggraph computer graphics, 21(4):163–169, 1987.

- [26] Shirui Luo, Changqing Zhang, Wei Zhang, and Xiaochun Cao. Consistent and specific multi-view subspace clustering. In Proc. of the AAAI conference on artificial intelligence, volume 32, 2018.
- [27] Xiaoxu Meng, Weikai Chen, and Bo Yang. Neat: Learning neural implicit surfaces with arbitrary topologies from multi-view images. Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, June 2023.
- [28] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In Proc. of IEEE Conf. on Computer Vision and Pattern Recognition, 2019.
- [29] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [30] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *Proc. IEEE international symposium on mixed and augmented reality*, pages 127–136. Ieee, 2011.
- [31] Liang Pan, Tong Wu, Zhongang Cai, Ziwei Liu, Xumin Yu, Yongming Rao, Jiwen Lu, Jie Zhou, Mingye Xu, Xiaoyuan Luo, et al. Multi-view partial (mvp) point cloud challenge 2021 on completion and registration: Methods and results. arXiv preprint arXiv:2112.12053, 2021.
- [32] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In Proc. of the IEEE/CVF conference on computer vision and pattern recognition, pages 165–174, 2019.
- [33] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In Proc. of the IEEE conference on computer vision and pattern recognition, pages 652–660, 2017.
- [34] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. Advances in neural information processing systems, 30, 2017.
- [35] Radu Bogdan Rusu, Nico Blodow, and Michael Beetz. Fast point feature histograms (fpfh) for 3d registration. In *Proc. of IEEE international conference on robotics and automation*, pages 3212–3217. IEEE, 2009.
- [36] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5459– 5469, 2022.
- [37] 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 Proc. of the IEEE/CVF conference on computer vision and pattern recognition, pages 2446–2454, 2020.
- [38] Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Multi-view consistency as supervisory signal for learning

shape and pose prediction. In *Proc. of the IEEE conference on computer vision and pattern recognition*, pages 2897–2905, 2018.

- [39] Nitika Verma, Edmond Boyer, and Jakob Verbeek. Feastnet: Feature-steered graph convolutions for 3d shape analysis. In *Proc. of the IEEE conference on computer vision and pattern recognition*, pages 2598–2606, 2018.
- [40] Dan Wang, Xinrui Cui, Xun Chen, Zhengxia Zou, Tianyang Shi, Septimiu Salcudean, Z Jane Wang, and Rabab Ward. Multi-view 3d reconstruction with transformers. In Proc. of the IEEE/CVF International Conference on Computer Vision, pages 5722–5731, 2021.
- [41] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proc. of the European conference on computer vision*, pages 52–67, 2018.
- [42] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. arXiv preprint arXiv:2106.10689, 2021.
- [43] Xiaogang Wang, Marcelo H Ang, and Gim Hee Lee. Voxelbased network for shape completion by leveraging edge generation. In *Proc. of the IEEE/CVF international conference on computer vision*, pages 13189–13198, 2021.
- [44] Haozhe Xie, Hongxun Yao, Shangchen Zhou, Jiageng Mao, Shengping Zhang, and Wenxiu Sun. Grnet: Gridding residual network for dense point cloud completion. In *Proc. of the European conference on computer vision*, pages 365–381. Springer, 2020.
- [45] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Pointbased neural radiance fields. In Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5438–5448, 2022.
- [46] Jianglong Ye, Yuntao Chen, Naiyan Wang, and Xiaolong Wang. Gifs: Neural implicit function for general shape representation. In Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 12829–12839, 2022.
- [47] Jianglong Ye, Yuntao Chen, Naiyan Wang, and Xiaolong Wang. Online adaptation for implicit object tracking and shape reconstruction in the wild. *IEEE Robotics and Automation Letters*, 7(4):8909–8916, 2022.
- [48] Xumin Yu, Yongming Rao, Ziyi Wang, Zuyan Liu, Jiwen Lu, and Jie Zhou. Pointr: Diverse point cloud completion with geometry-aware transformers. In *Proc. of the IEEE/CVF international conference on computer vision*, pages 12498– 12507, 2021.
- [49] Xumin Yu, Yongming Rao, Ziyi Wang, Jiwen Lu, and Jie Zhou. Adapointr: Diverse point cloud completion with adaptive geometry-aware transformers. arXiv preprint arXiv:2301.04545, 2023.
- [50] Wentao Yuan, Tejas Khot, David Held, Christoph Mertz, and Martial Hebert. Pcn: Point completion network. In *Proc.* of International Conference on 3D Vision, pages 728–737, 2018.

[51] Changqing Zhang, Zongbo Han, Huazhu Fu, Joey Tianyi Zhou, Qinghua Hu, et al. Cpm-nets: Cross partial multiview networks. Advances in Neural Information Processing Systems, 32, 2019.