

Sparse Dense Fusion for 3D Object Detection

Yulu Gao, Chonghao Sima, Shaoshuai Shi, Shangzhe Di, Si Liu and Hongyang Li

Abstract—With the prevalence of multimodal learning, camera-LiDAR fusion has gained popularity in 3D object detection. Although multiple fusion approaches have been proposed, they can be classified into either sparse-only or dense-only fashion based on the feature representation in the fusion module. In this paper, we analyze them in a common taxonomy and thereafter observe two challenges: 1) sparse-only solutions preserve 3D geometric prior and yet lose rich semantic information from the camera, and 2) dense-only alternatives retain the semantic continuity but miss the accurate geometric information from LiDAR. By analyzing these two formulations, we conclude that the information loss is inevitable due to their design scheme. To compensate for the information loss in either manner, we propose Sparse Dense Fusion (SDF), a complementary framework that incorporates both sparse-fusion and dense-fusion modules via the Transformer architecture. Such a simple yet effective sparse-dense fusion structure enriches semantic texture and exploits spatial structure information simultaneously. Through our SDF strategy, we assemble two popular methods with moderate performance and outperform baseline by 4.3% in *mAP* and 2.5% in NDS, ranking first on the nuScenes benchmark. Extensive ablations demonstrate the effectiveness of our method and empirically align our analysis.

I. INTRODUCTION

The human experience of the world is multimodal - through vision, auditory, thermoception, and other sensory systems. Modality denotes how something happens or is experienced and multimodal learning is characterized as a research problem when a machine learning system includes multiple such modalities [1]. This problem comes of vital importance in the scenario of autonomous driving as sensor configurations get more complex on vehicles, providing a suitable playground [2]. Within the realm of 3D object detection, point cloud from LiDAR and RGB images from camera are two general perception sources with complementarity and redundancy to each other. While point cloud can provide accurate 3D spatial structure around the ego vehicle, RGB images perceive the scene with rich semantic information.

Given such a complementary setting of sensors, multimodal learning naturally raises the audience’s interest in information interaction between those sensors. However, the design is non-trivial due to the considerable difference between the sparsity in LiDAR data and the density in camera data. Distinguished by the representation of feature in fusion modules, we categorize existing methods into two fashion: *sparse-only* fusion and *dense-only* fusion. To alleviate the sparsity in LiDAR data, *sparse-only* methods

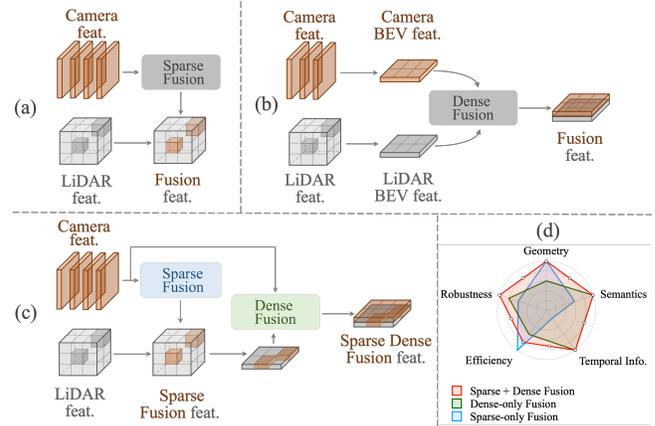


Fig. 1: **Design scheme** of (a) sparse-only fusion, preserving 3D geometric prior and yet losing rich semantic information, (b) dense-only fusion, retaining the semantic continuity but missing the accurate geometric information, and (c) our sparse-dense fusion, preserving 3D geometric prior and semantic information simultaneously. (d) illustrates the strengths and weaknesses of each fashion and the corresponding design scheme. The difference in geometry, semantics, and temporal information is demonstrated via *mAP* and NDS in Tab. I, while the difference in robustness explained in Tab. IV and efficiency is explained in Tab. V.

usually transform image semantic label [3], [4] or feature [5], [6], [7] to point cloud representation via LiDAR-camera sampling [8], [9] or predefined depth bins [10], [11]. To address the density in camera data, intuitive approaches in *dense-only* fashion are fusing feature from both modalities in bird’s-eye-view (BEV) representation [12], [13], [14], [15], [16], [17]. Besides BEV, voxel representation also draws some interest in dense-only fusion. The camera-voxel transformation utilizes either camera parameters [18] or estimated depth per pixel [19] to establish the projection mapping. As we observe this vast difference between the two fashions, some natural but non-trivial questions are raised - *what is the core issue left unexplored in sensor fusion? what are the key factors to improve the upper bound of performance?*

We conduct a caveat analysis of the prevailing sensor fusion methods, and it reveals that either fashion has inherent advantages and disadvantages due to its design scheme. Fig. 1 (d) (a) (b) illustrates the strengths and weaknesses of each fashion and corresponding design scheme in a nutshell.

In the *sparse-only* fashion, point cloud is projected to image plane and gather semantic label or image feature around the projection position, thus fusing image information

Yulu Gao, Shangzhe Di, and Si Liu are with Institute of Artificial Intelligence, Beihang University, Beijing, China. Chonghao Sima and Hongyang Li are with OpenDriveLab, Shanghai AI Laboratory, Shanghai, China. Shaoshuai Shi is with the Department of Electronic Engineering at The Chinese University of Hong Kong, Hong Kong, China.

to point cloud representation. Though this can preserve 3D geometric information in the point cloud feature, the rich semantic information provided by the image feature is broken. This is because a 3D object, described by several points in point cloud data, corresponds to hundreds of pixels in image data. This sparsity-density inconsistency leads to a problem: when model appends image feature to voxel representation, around 1 to 2 magnitude of image feature is deprecated via LiDAR-camera projection.

In *dense-only* fashion, image feature and point cloud feature are transformed to the BEV space, respectively, then fused grid by grid there. While transforming image feature to BEV space remains semantic information, squeezing point cloud feature to BEV space genuinely loses the geometric prior in 3D space. Also, projecting image feature to BEV space is ill-posed since camera does not capture any 3D geometry. This results in the grid inconsistency of the same object between image-BEV and point-cloud-BEV feature, increasing the difficulty of fusion. Therefore, information loss is inevitable in either manner, and we are motivated to design a fusion mechanism that sparse-only and dense-only fusion can jointly contribute to the network.

Given that the two fashions operate feature map in different representations (voxel, perspective and BEV), Transformer architecture [20] meets our demands in aggregating valuable information from different modalities conceptually. Therefore, we propose **Sparse-Dense Fusion (SDF)**, a complementary structure connecting sparse-fusion and dense-fusion modules to enjoy the best of both worlds. Conceptually, SDF extracts image and point cloud feature, respectively, then performs sparse-fusion to aggregate geometric prior and dense-fusion to exploit semantic information. The sparse-fusion module is responsible for gathering image feature around non-empty voxels, aiming to preserve the geometric prior provided by point cloud feature. Meanwhile, the dense-fusion module projects point cloud feature to bird’s-eye-view and queries image feature in perspective view, in the goal of utilizing rich semantic information in image. Illustrated by Fig. 1 (c), such a simple yet effective design can cover the disadvantage of sparse-only or dense-only fashion and enjoy the best of both worlds. To demonstrate the effectiveness of our SDF design, we assemble two moderate methods, BEVFormer [21] for camera branch and TransFusion [17] for LiDAR branch to construct the baseline network. Then the SDF framework improves their joint performance by a large margin, outperforms baseline by **4.3%** in *mAP* and **2.5%** in NDS and reaches the state of the art on the nuScenes benchmark [2].

Extensive ablations respond to our analysis in the two fashions as well. Dense-only fashion has the advantage in texture-intensive categories such as truck and car, while sparse-only fashion is capable of localization-aware categories such as trailer and pedestrian.

To sum up, our work has three-fold contributions: (a) We propose a caveat analysis on existing sensor-fusion methods and point out their inevitable information loss respectively. (b) We devise a simple yet effective framework, named

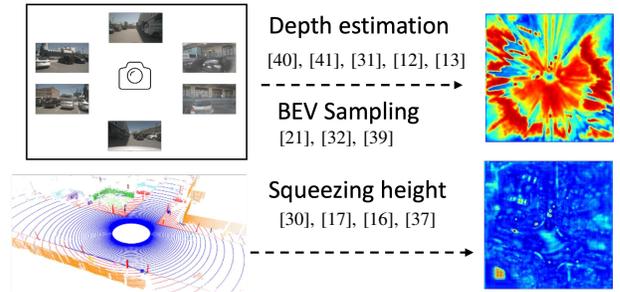


Fig. 2: Visualization of *dense-only* manner. Camera feature is transformed to BEV via depth estimation or BEV sampling, while LiDAR feature is via squeezing height in voxel space. It can be observed that there exists “blurring long tail” in the shape of ray-tracing dilation. Since transforming camera feature to BEV is ill-posed, feature from two modalities cannot align well grid by grid, consequently hindering dense feature fusion.

Sparse-Dense Fusion, which shares the advantage of sparse and dense fusion via a complementary structure. (c) The proposed framework improves the baseline method by a large margin and ranks first on the nuScenes benchmark with extensive ablations.

II. RELATED WORK

A. 3D Object Detection in General

As autonomous driving system is a hierarchical design of perception [22], prediction [23] and planning [24], [25], [26], [27], 3D bounding boxes usually play the role of connecting the pipeline. Based on the input modality, we divide 3D object detection research into three parts mainly - camera-based [28], [29], LiDAR-based [30] and sensor-fusion-based [12] methods. Compared to the perspective view originating from camera, BEV representation has taken camera-based 3D perception by storm in terms of objection detection [21], [31], lane detection [32], [33] and BEV map segmentation [34]. As BEV representation is a natural and straightforward candidate view to serve as a unified representation [22] across different modalities, representing image feature in BEV set up the foundation of sensor-fusion methods [12]. In LiDAR-based methods, point cloud data is first processed to voxel [35] or pillar [36] space, then go through 3D/2D feature extractor [37], and finally, perform detection in BEV space [30], [38]. Different from representing image feature in BEV, LiDAR-to-BEV transformation is more intuitive as point cloud data naturally preserve 3D geometric information.

Given such progress in camera-based and LiDAR-based methods, fusing them to enjoy the best of both worlds is a natural idea, and many attempts did happen in the sensor-fusion domain. In general, we present in detail various early- or mid-fusion methods in a taxonomy and categorize them into two fashions, termed *dense-only* and *sparse-only*.

B. Dense-only Fusion

In dense-only fusion, image feature is either sampled to voxel space or projected to BEV space, and fused with point

cloud feature in that representation. Inspired by Lift-Splat-Shoot [40], projecting image feature to BEV space via depth estimation and simply concatenate with point cloud feature in BEV space have demonstrate its effectiveness by two BEVFusion [12], [13]. Similar approaches follow the idea of fusing feature in BEV space with additional tricks such as region of interest pooling in DeepInteraction [14] and 3D-CVF [15], BEV space augmentation in DeepFusion [16] and two-stage refinement in TransFusion [17]. Another line of works attempts to fuse feature in voxel space by expanding image feature to frustum and aligning it to voxel representation. UVTR [18] transforms image feature to voxel space via predefined 3D point sampling, and aligns it with point cloud voxel feature. Unlike BEV representation, voxel space preserve more 3D geometric information.

C. Dense-only Fusion

In sparse-only fusion, point cloud is projected to image plane and gather semantic label or image feature around the projection position. In terms of gathering semantic labels, point-painting [3] and point-augmenting [4] project point cloud to image plane and append semantic labels around projection location to 3D points. Recently the trend has been switched to fusing feature under this philosophy, such as AutoAlign [5], its V2 [6] and EPNet++ [7] focusing on the pixel-instance feature, MMF [8] with multi-scale fusion, SFD [10] constructing pseudo point cloud from image, and VPFNet [11] equipped with stereo data augmentation. Another line of work applies object-level sparse fusion such as FUTR3D [42] and CLOCS [43].

III. ANALYSIS

Though plenty of work has explored the potential of sensor-fusion following either sparse-only or dense-only fashion, both of them have inevitable information loss inherited from the design philosophy, respectively. In this section, we will first make an ideal analysis of how a 3D object is perceived via different sensors. Later, we highlight how the information loss happened in current sensor-fusion methods and finally introduce the motivation of our framework.

A. Representation Matters in 3D Perception

An object in the real world is captured by LiDAR into point cloud representation and by camera into image representation. LiDAR is often regarded as a critical sensor for distance sensing in autonomous driving due to its ability to perceive and capture the 3D geometry of the environment. The prevailing LiDAR configuration is equipped with a fixed number of beams, a default scanning frequency, and finally provides a stable number of scanned points per ring with slight fluctuation. Consequently, the size of the point cloud data per frame is usually limited to tens of thousands, e.g., in nuScenes [2], approximately 35,000. The camera is often seen as a complementary sensor for semantic information as it can capture the texture appearance of objects in the field of view (FOV). Predominant camera configuration comprises several cameras with different locations and orientations

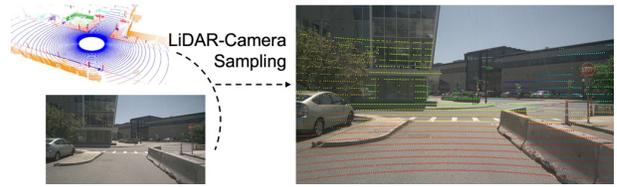


Fig. 3: Visualization of *sparse-only* manner. Point cloud is projected to image plane and gather semantic label or image feature around the projection position. Since the projection positions are sparse on image plane, the process of sampling feature will break the texture continuity in image.

on the autonomous vehicle. Each camera can perceive the scene in image with a resolution of, e.g., 900×1600 in nuScenes [2]. Therefore, in each frame, the camera system receives approximately 8,000,000 pixels in total. The huge quantitative and qualitative difference between point cloud and image data implies the difficulty in designing sensor-fusion algorithms.

B. Achilles Heel of Sensor Fusion

In a nutshell, the core problem of sensor fusion is to align the feature from different modalities, so transforming feature into the exact representation plays a vital role in feature. Distinguished by where and how to transform feature, two intuitive methodologies have taken camera-LiDAR fusion by storm, the *sparse-only* and the *dense-only* fashion.

In sparse-only fusion, image feature is transformed to point cloud representation via LiDAR-camera sampling, and feature fusion operates only on those sampled positions. The sampling process can be formulated as:

$$z_c \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} [\mathbf{R} \quad \mathbf{T}] \begin{bmatrix} x_w & y_w & z_w & 1 \end{bmatrix}^T, \quad (1)$$

where (u, v) denotes 2D position in image coordinate, z_c denotes the depth per pixel, (x_w, y_w, z_w) denotes the 3D position in the world coordinate, (f_x, f_y, c_x, c_y) denotes camera intrinsic, $(\mathbf{R} \quad \mathbf{T})$ denotes camera-LiDAR rotation and translation coefficient. Such a transformation, however, breaks the texture continuity in image feature, and consequently cannot utilize the rich semantic information extracted via image backbone [44]. Numerically, because of sparsity in point cloud, the sampled positions only count for around 1/100 to 1/1000 of the object appearance on the image. It's more evident through visualization in Fig. 3. Therefore, sparse-only fashion has information loss in rich semantic feature provided by image.

In dense-only fusion, the image and point cloud features are transformed to BEV space, respectively, and are fused grid by grid. Because BEV is a unified representation of different modalities, fusing features in BEV space is prevailing [12], [13], [14].

However, through compressing point cloud feature to BEV space, its 3D geometric structure is broken inevitably. Meanwhile, projecting image feature to BEV space is ill-

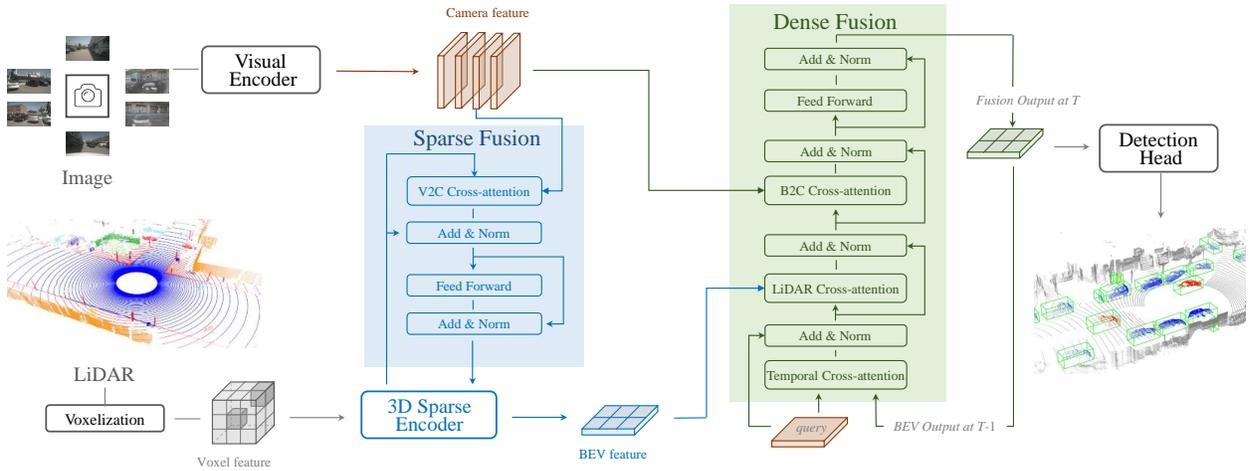


Fig. 4: **Pipeline of our proposed SDF method.** SDF includes two sub-parts: Sparse Fusion Module and Dense Fusion Module. In Sparse Fusion Module, the non-empty voxel query interacts with image features. In Dense Fusion Module, each BEV query interacts alternately with the history BEV feature, lidar BEV feature, and image feature.

posed, resulting in location inconsistency. Through visualization in Fig. 2, it’s observed that there exists “blurring long tail” in camera-BEV feature because image does not preserve any 3D geometric prior, making the projection ill-posed. Consequently, an object represented by several BEV grids in LiDAR-projected BEV feature may not align well with the corresponding BEV grids in camera-projected BEV feature. Therefore, dense-only fashion has information loss in accurate 3D geometric feature provided by point cloud.

To this point, we have witnessed the strengths and weaknesses of each fashion; they are complementary to each other. A straightforward idea is to incorporate both design schemes into a single network and enjoy the benefits from both worlds.

IV. METHODOLOGY

This paper proposes a unified camera-lidar fusion framework based on attention. Our main contribution is fusing camera and lidar from both sparse and dense perspectives.

A. Overall Structure

The core idea of our SDF is fusing camera and lidar from both sparse and dense views. The overall architecture is illustrated in Fig. 4. SDF includes two sub-parts: Sparse Fusion Module and Dense Fusion Module. Considering that deformable transformer [45] has linear complexity and fast convergence speed, our module uses a deformable transformer [45] to gather and fuse features in each part. Specifically, the Sparse Fusion Module fuses image features to valid voxels. Here, *valid voxels* denote the voxels with lidar points inside. The fused voxel features maintain the original 3D spatial structure and are further sent to the subsequent network. As for dense fusion, we adopt BEV queries to integrate camera, lidar, and past BEV features, respectively. The BEV queries are grid-shaped parameters and are learned during training.

B. Sparse Fusion in Voxel Space

The purpose of sparse fusion is to preserve 3D geometric prior in the fusion process. In the pipeline of the entire lidar branch, voxel space can maintain geometric information. Since the voxel space is large and sparse, taking all voxel as queries for fusion is useless and inefficient. Thus, we propose a modified deformable attention module to establish valid connections between the voxel and camera features.

Specifically, we focus on valid voxels to improve efficiency and accuracy. We take the centers of all valid voxels as 3D reference points and project them onto image planes to get 2D reference points. Then, we use V2C(voxel-to-camera) cross-attention to gather the camera feature, defined as:

$$\text{V2C}(Q_p, F) = \frac{1}{|\mathcal{V}_{\text{hit}}|} \sum_{i \in \mathcal{V}_{\text{hit}}} \sum_{j=1}^{N_{\text{ref}}} \text{DeAttn}(Q_p, \mathcal{P}(p, i, j), F^i), \quad (2)$$

where \mathcal{V}_{hit} indexes the number of images hit by 2D reference points, i indexes the image view, j indexes the reference points, and N_{ref} is the total number of reference points for each voxel query. F^i is the features of the i -th camera view. For each voxel query Q_p , we utilize a project function $\mathcal{P}(p, i, j)$ to get the j -th reference point on the i -th view image.

Following that, we show how to build 2D reference points on a view image using the projection function \mathcal{P} and how to adapt to lidar geometric augmentation. We calculate the real world position (x', y', z') corresponding to the query Q_p located at $p = (x, y, z)$ of Q as follows:

$$x' = (x - \frac{W}{2}) \times s_x; \quad y' = (y - \frac{H}{2}) \times s_y; \quad z' = (z - Z) \times s_z + z_{\text{base}}, \quad (3)$$

where H, W, Z are the spatial shapes of the entire voxel feature, s_x, s_y, s_z are the size of resolution of voxel, z_{base} is the base height, and (x', y', z') are the coordinate of the ego vehicle coordinate system. As a result, for each query

	S-Fusion	D-Fusion	Temporal	mAP	NDS	Car	Truck	C.V.	Bus	Trailer	Barrier	Motor.	Bike	Ped.	T.C.
1				64.9	69.9	86.9	60.9	26.1	73.7	42.8	68.8	71.3	55.9	86.8	74.4
2	✓			68.1(+3.2)	70.9(+1.0)	88.4	64.4	30.0	77.2	45.8	69.4	77.0	63.1	88.1	77.5
3		✓		67.9(+3.0)	71.5(+1.6)	89.4	65.8	28.6	77.4	43.8	70.6	76.9	61.2	87.7	77.7
4	✓	✓		68.8(+3.9)	72.0(+2.1)	89.5	65.9	30.7	77.6	45.2	70.6	78.2	63.5	88.3	77.9
5	✓	✓	✓	69.2(+4.3)	72.4(+2.5)	89.5	67.0	32.4	77.3	43.4	71.5	79.2	65.7	88.5	77.8

TABLE I: **Component analysis.** We verify the effectiveness of the Sparse Fusion Module and Dense Fusion Module. “S-Fusion” and “D-Fusion” denotes the Sparse and Dense Fusion Module, respectively. “Temporal” means the Dense Fusion Module with temporal information. “C.V.,” “Ped.,” and “T.C.” are abbreviations of construction vehicle, pedestrian, and traffic cone, respectively. Compared with each other, sparse-only fashion performs better on small objects such as bikes and pedestrians, while dense-only fashion performs better on objects with large scale and regular locations such as cars, trunks, and buses. Sparse Dense Fusion (SDF) benefits from both.

Q_p , we obtain 3D reference points (x', y', z') . Finally, we project the 3D reference points to different image views using Eqn. 1. When we apply data augmentation to a point cloud, we reverse all 3D reference points according to data augmentation, and then use the projection matrix to project the 3D reference points to 2D reference points.

C. Dense Fusion in BEV Space

Our Dense Fusion Module is composed of 6 transformer layers, and each transformer layer is composed of Temporal Cross-attention, LiDAR Cross-attention, B2C(Bev to Camera) Cross-attention module, and Feed Forward module. The purpose of these cross-attention modules is to gather temporary information, lidar information, and image information, respectively. In the beginning, we initialize the BEV query following BEVFormer [21].

Different from the sparse fusion, which only performs fusion on valid voxel, Dense Fusion Module treats each BEV query equally. First, we initialize the BEV query with a set of learnable parameters. Then, the BEV queries are used to collect and integrate features from different modalities. Directly fusion in BEV space will cause semantic ambiguity because of the height compression. Following BEVFormer [21], we predefine a set of anchor heights for each BEV query to alleviate this problem.

When fusing camera features, we set 3D reference points for the BEV query at different heights and use the same projection method described in section IV-B to generate the 2D reference points. As for lidar fusion, we directly set reference points to gather lidar features at the corresponding locations.

As shown in Fig. 4, our Dense Fusion Module contains a modified transformer layer that alternately fuses the lidar, camera, and history BEV features. The BEV features generated by the Dense Fusion Module will be sent to the 3D detection head. In contrast to Temporal Cross-attention, we first align B_{t-1} to Q according to ego-motion, then update the 2D reference point location to keep the features at the same grid corresponding to the exact real-world location. With the above modification, we get the final BEV feature equipped with temporal and multimodal information.

V. EXPERIMENTS

In this section, we conduct extensive comparisons and ablation studies to examine the effectiveness of our method.

We use the gray background to indicate our method in the following tables.

A. Setup and Implementation Details

Dataset. We conduct experiments on the nuScenes dataset [2], which contains 1000 scenes, with each sample consisting of RGB images from 6 cameras with 360° horizontal FOV and one LiDAR point cloud. There are 1.4 million 3D bounding boxes with annotations from ten categories. We use mean average precision (mAP) and nuScenes detection score (NDS) evaluation metrics provided by the dataset to evaluate our results.

Implementation details. We implement SDF with Python 3.8.13, PyTorch 1.9.0, and CUDA 11.3. All models are trained with 8 Nvidia A100 GPUs. Following BEVFormer [21], we use ResNet101-DCN [44] as our image backbone. For lidar, we use the VoxelNet [35] as our LiDAR backbone, and TransFusion-L [17] is chosen as our 3D detection head. And we also conduct a two-stage training process following BEVFusion [13]. Like BEVFusion, we do not use data augmentation when training the camera and fusion models. And we only use data augmentation for the fusion model when comparing it with the state-of-the-art methods. And we employ a single model during testing with no test-time augmentation.

B. Components of SDF

We examine the effectiveness of each proposed module, and results are listed in Tab. I. The 1st row is our baseline Transfusion-L [17]. We individually add Sparse Fusion Module and Dense Fusion Module to the baseline model, and the result is shown in the 2nd and 3rd rows. Each component improves the performance, indicating their effectiveness in the 3D detection task. In particular, we can observe that sparse fusion has better performance on small objects such as bikes and pedestrians, which proves that sparse fusion can effectively leverage the spatial structure and has a more accurate mapping from 3D space to camera view than dense fusion. And dense fusion performs better on objects with large scale and regular location such as car, trunk, and bus, because such objects rely less on spatial structure and dense fusion provides richer semantic information. The 4th row is the result of combining both fusion modules, achieving 4.0% mAP and 2.1% NDS absolute improvements compared to the baseline. Besides, using both fusion modules is better than



Fig. 5: **Visualization results of SDF on nuScenes val set.** We show the 3D bounding box predictions across each camera’s field of view and the bird-eye view. Except for a few errors in crowded areas and remote objects, our method yields remarkable results.

using only one, showing that our sparse and dense modules can complement each other. Additionally, in the 5th row, we extend our model in the temporal dimension, and the result shows that our model can effectively utilize the temporal information.

	Loc	Method	mAP	NDS
1	C_1	deformable attention	66.6	70.4
2	C_2	deformable attention	66.9	70.0
3	C_3	deformable attention	67.3	70.8
4	C_4	deformable attention	68.1	70.9
5	C_5	deformable attention	67.4	70.4
6	C_4	extended deformable attention	68.1	71.2

TABLE II: Sparse Fusion Module analysis. We explore the influence of different fusion locations and different fusion methods in the sparse module. “Loc” denotes the location of the Sparse Fusion Module. C_x means insert the deformable attention after layer x of the sparse encode.

C. Ablation Study

Robustness. Here, we simulate LiDAR sensor failures to test the robustness of different methods. Following BEVFusion [13], we set a limited field-of-view (FOV) and remove the out-of-view LiDAR points. Tab. IV shows the influence of limited LiDAR field-of-view. One can observe that the Sparse-only method is slightly better than the LiDAR-only method. The Sparse-only method relies heavily on the LiDAR point, so the fusion process can only bring slight improvement compared to the LiDAR-only method. The Dense-only method is much better than the LiDAR-only method, which achieves 21.3% and 12.7% absolute improvements on mAP and NDS, respectively, under the robustness setting.

Efficiency. We compare the efficiency of the Sparse Fusion Module and Dense Fusion Module in the fusion process by calculating the number of fusion blocks (voxels or BEV grids). Because the number of valid voxels is related to the input point cloud, here we count the average of valid voxels after the fourth layer of the sparse encoder in Val set. The dense Fusion Module fuses on all grids, so the number of fusion grids depends on the size of the BEV.

In our experiment, the BEV size is set to 180×180 , so the number of fusion grids is 32400. As shown in Tab. V, the number of fusion blocks for our Sparse Fusion Method is approximately 13k, fewer than the Dense Fusion Method, which proves that sparse fusion is more efficient than dense fusion.

Sparse Fusion Module. We then explore different settings of the sparse module and summarize the results in Tab. II. From the 1st to 5th row, we insert our V2C deformable attention into the different locations of the sparse encoder. It is effective to insert the deformable attention after the fourth layer of the sparse encoder. In other words, deeper voxel features are more beneficial for aggregating image information. Furthermore, we conduct experiments on the use of reference points. In the 6th row, we use offset sampling while setting multiple reference points to voxels adjacent to the query. The modification (Row 6) performs slightly better than the original deformable attention (Row 4) because of its wider perception field. Nevertheless, we still use the original V2C deformable attention to keep it simple.

Dense Fusion Module. We compare our fusion method in the dense module with dynamic fusion from BEVFusion [13]. Dynamic fusion first concatenates multimodal features and fuses them with learnable static weights, then applies channel attention to select fused features. This module is simple but lacks the ability for iterative interaction between features. To apply the dynamic fusion module in our framework, we remove the lidar input in the dense fusion encoder to get the pure BEV camera feature and then fuse the lidar and camera BEV features with the dynamic fusion module. As shown in Tab. VI, our dense fusion method is better than dynamic fusion when using the same number of encoders. With iterative interaction, our fusion encoder achieves more sufficient fusion. We also test the scalability of our Dense Fusion Module and conclude that more encoder layers can lead to better results.

Different Detection Methods. Tab. VII shows the result of different detection methods with our fusion module. Our fusion method performs feature fusion at the voxel or BEV space, which can be easily extended to existing detection

Method	Modality	Fusion Type	mAP	NDS	Car	Truck	C.V.	Bus	Trailer	Barrier	Motor.	Bike	Ped.	T.C.
<i>val</i>														
FUTR3D [42]	LC	D	64.2	68.0	86.3	61.5	26.0	71.9	42.1	64.4	73.6	63.3	82.6	70.1
BEVFusion* [13]	LC	D	69.6	72.1	89.1	66.7	30.9	77.7	42.6	73.5	79.0	67.5	89.4	79.3
BEVFusion [12]	LC	D	68.5	71.4	-	-	-	-	-	-	-	-	-	-
Deepinteraction [14]	LC	D	69.9	72.6	-	-	-	-	-	-	-	-	-	-
Ours	LC	SD	68.8	72.0	89.5	65.9	30.7	77.6	45.2	70.6	78.2	63.5	88.3	77.9
Ours-T	LC	SD	69.2	72.4	89.5	67.0	32.4	77.3	43.4	71.5	79.2	65.7	88.5	77.8
Ours*	LC	SD	70.4	72.8	90.0	68.2	33.5	79.5	48.3	70.8	78.4	66.5	89.1	79.4
<i>test</i>														
PointPillars [36]	L	-	30.5	45.3	68.4	23.0	4.1	28.2	23.4	38.9	27.4	1.1	59.7	30.8
CenterPoint [†] [46]	L	-	60.3	67.3	85.2	53.5	20.0	63.6	56.0	71.1	59.5	30.7	84.6	78.4
TransFusion-L [17]	L	-	65.5	70.2	86.2	56.7	28.2	66.3	58.8	78.2	68.3	44.2	86.1	82.0
PointPainting [3]	LC	S	46.4	58.1	77.9	35.8	15.8	36.2	37.3	60.2	41.5	24.1	73.3	62.4
3D-CVF [15]	LC	D	52.7	62.3	83.0	45.0	15.9	48.8	49.6	65.9	51.2	30.4	74.2	62.9
TransFusion [17]	LC	D	68.9	71.7	87.1	60.0	33.1	68.3	60.8	78.1	73.6	52.9	88.4	86.7
Autoalignv2 [6]	LC	S	68.4	72.4	87.0	59.0	33.1	69.3	59.3	78.0	72.9	52.1	87.6	85.1
BEVFusion* [13]	LC	D	71.3	73.3	88.5	63.1	38.1	72.0	64.7	78.3	75.2	56.5	90.0	86.5
BEVFusion [12]	LC	D	70.2	72.9	88.6	60.1	39.3	69.8	63.8	80.0	74.1	51.0	89.2	86.5
DeepInteraction [14]	LC	D	70.8	73.4	87.9	60.2	37.5	70.8	63.8	80.4	75.4	54.5	91.7	87.2
Ours*	LC	SD	70.9	73.4	89.2	63.2	38.6	73.2	65.5	80.9	75.3	48.1	89.2	86.0

TABLE III: Performance on the nuScenes *val* and *test* sets. ‘‘C.V.’’, ‘‘Ped.’’, and ‘‘T.C.’’ are abbreviations of construction vehicle, pedestrian, and traffic cone, respectively. ‘L’ and ‘C’ represent LiDAR and camera, respectively. ‘D’, ‘S’, and ‘SD’ represent ‘Dense-only’, ‘Sparse-only’, and ‘Sparse Dense’. † indicates the results with double-flip during the test. * indicates method with data augmentation during training. Even *without* temporal information, our SDF achieves new state-of-the-art on both *val* and *test* sets. Particularly, the augmentation version of SDF achieves 0.8% and 0.6% absolute improvements compared to BEVFusion with augmentation during training.

	Fusion type	FOV	mAP	NDS
1	-	-	64.9	69.9
2	-	$(-\pi/2, \pi/2)$	28.1	50.6
3	S	-	68.1	70.9
4	S	$(-\pi/2, \pi/2)$	30.5	51.4
5	D	-	67.9	71.5
6	D	$(-\pi/2, \pi/2)$	49.4	58.2
7	SD	-	68.8	72.0
8	SD	$(-\pi/2, \pi/2)$	50.2	59.0

TABLE IV: The result of limited LiDAR field-of-view. Dense-only Fusion outperforms baseline by 21.3% mAP and 12.7% NDS under the robustness setting, which proves its robustness.

	Fusion type	NF	mAP	NDS
1	-	-	64.9	69.9
2	S	13k	68.1	70.9
3	D	32k	67.9	71.5
4	SD	55k	68.8	72.0

TABLE V: The number of fusion blocks. By calculating the number of fusion blocks, we can find that Sparse-only Fusion has fewer fusion blocks than Dense-only Fusion.

methods. With the SD-Fusion, Centerpoint and Transfusion-L get 6.8% 3.8% and 3.9%, 2.8% absolute improvements on mAP and NDS respectively, indicating that our approach has the general ability to apply to different detection methods.

Different frame number during training. Tab. VIII shows the effect of the frame number during training. We can observe that introducing temporal information could bring improvements to our model. Particularly, our method with 3 frames during training achieves 0.4% and 0.4% absolute improvements compared to the version that does not leverage

	Method	# Encoder Layer	mAP	NDS
1	dynamic fusion	6	66.9	70.9
2	fusion encoder	1	66.8	70.8
3	fusion encoder	6	67.9	71.5

TABLE VI: Dense Fusion Module analysis. We explore the impact of different fusion methods and different numbers of encoder layers.

	Method	SD-Fusion	mAP	NDS
1	Centerpoint		57.1	65.4
2	Centerpoint	✓	63.9	69.2
3	Transfusion-L		64.9	69.9
4	Transfusion-L	✓	68.8	72.0

TABLE VII: Expandability analysis. With our Sparse Dense Fusion, Centerpoint and Transfusion-L get 6.8% 3.8% and 3.9%, 2.8% absolute improvements on mAP and NDS respectively.

temporal information.

D. Comparison to State-of-the-arts

In Tab. III, we report our main results on nuScenes *val* and *test* splits. On the *val* set, our method achieves the state-of-the-art performance of 70.4% mAP and 72.7% NDS, outperforming all previous methods. Compared to BEVFusion, which is a dense-only fashion, our method can surpass it in terms of both non-augmentation and augmentation versions. In particular, the augmentation version of our method achieves 0.8% and 0.6% absolute improvements compared to BEVFusion. On the *test* set, our method achieves state-of-the-art with 70.9% mAP and 73.4% NDS.

	#Frame	mAP	NDS	mAVE
1	1	68.8	72.0	26.2
2	2	68.7	72.0	25.8
3	3	69.2	72.4	25.5

TABLE VIII: Temporal analysis. We explore the influence of different frame numbers during training. “#Frame” denotes the frame number during training.

Visualization We show the detection results of a complex scene in Fig. 5. Except for a few errors in crowded areas and remote objects, our method yields remarkable results.

VI. CONCLUSION

In this paper, we analyze current fusion methods and point out their inevitable information loss. Thus, we propose a unified framework termed Sparse-Dense Fusion, which enriches semantic texture and exploits spatial structure. The framework significantly outperforms the baseline method and ranks first on the nuScenes benchmark.

REFERENCES

- [1] T. Baltrušaitis, C. Ahuja, and L.-P. Morency, “Multimodal machine learning: A survey and taxonomy,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 2, pp. 423–443, 2018.
- [2] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, “nusenes: A multimodal dataset for autonomous driving,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 621–11 631.
- [3] S. Vora, A. H. Lang, B. Helou, and O. Beijbom, “Pointpainting: Sequential fusion for 3d object detection,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 4604–4612.
- [4] C. Wang, C. Ma, M. Zhu, and X. Yang, “Pointaugmenting: Cross-modal augmentation for 3d object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 794–11 803.
- [5] Z. Chen, Z. Li, S. Zhang, L. Fang, Q. Jiang, F. Zhao, B. Zhou, and H. Zhao, “Autoalign: Pixel-instance feature aggregation for multi-modal 3d object detection,” *arXiv preprint arXiv:2201.06493*, 2022.
- [6] Z. Chen, Z. Li, S. Zhang, L. Fang, Q. Jiang, and F. Zhao, “Autoalignv2: Deformable feature aggregation for dynamic multi-modal 3d object detection,” *ECCV*, 2022.
- [7] Z. Liu, B. Li, X. Chen, X. Wang, X. Bai, *et al.*, “Epnnet++: Cascade bi-directional fusion for multi-modal 3d object detection,” *arXiv preprint arXiv:2112.11088*, 2021.
- [8] M. Liang, B. Yang, Y. Chen, R. Hu, and R. Urtasun, “Multi-task multi-sensor fusion for 3d object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7345–7353.
- [9] T. Yin, X. Zhou, and P. Krähenbühl, “Multimodal virtual point 3d detection,” 2021.
- [10] X. Wu, L. Peng, H. Yang, L. Xie, C. Huang, C. Deng, H. Liu, and D. Cai, “Sparse fuse dense: Towards high quality 3d detection with depth completion,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5418–5427.
- [11] H. Zhu, J. Deng, Y. Zhang, J. Ji, Q. Mao, H. Li, and Y. Zhang, “Vpfnet: Improving 3d object detection with virtual point based lidar and stereo data fusion,” *IEEE Transactions on Multimedia*, 2022.
- [12] Z. Liu, H. Tang, A. Amini, X. Yang, H. Mao, D. Rus, and S. Han, “Befusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation,” *arXiv*, 2022.
- [13] T. Liang, H. Xie, K. Yu, Z. Xia, Z. Lin, Y. Wang, T. Tang, B. Wang, and Z. Tang, “Befusion: A simple and robust lidar-camera fusion framework,” *arXiv preprint arXiv:2205.13790*, 2022.
- [14] Z. Yang, J. Chen, Z. Miao, W. Li, X. Zhu, and L. Zhang, “Deepinteraction: 3d object detection via modality interaction,” in *NeurIPS*, 2022.

- [15] J. H. Yoo, Y. Kim, J. Kim, and J. W. Choi, “3d-cvf: Generating joint camera and lidar features using cross-view spatial feature fusion for 3d object detection,” in *European Conference on Computer Vision*. Springer, 2020, pp. 720–736.
- [16] Y. Li, A. W. Yu, T. Meng, B. Caine, J. Ngiam, D. Peng, J. Shen, B. Wu, Y. Lu, D. Zhou, *et al.*, “Deepfusion: Lidar-camera deep fusion for multi-modal 3d object detection,” *arXiv preprint arXiv:2203.08195*, 2022.
- [17] X. Bai, Z. Hu, X. Zhu, Q. Huang, Y. Chen, H. Fu, and C.-L. Tai, “Transfusion: Robust lidar-camera fusion for 3d object detection with transformers,” *arXiv preprint arXiv:2203.11496*, 2022.
- [18] Y. Li, Y. Chen, X. Qi, Z. Li, J. Sun, and J. Jia, “Unifying voxel-based representation with transformer for 3d object detection,” in *Advances in Neural Information Processing Systems*, 2022.
- [19] X. Li, B. Shi, Y. Hou, X. Wu, T. Ma, Y. Li, and L. He, “Homogeneous multi-modal feature fusion and interaction for 3d object detection,” in *European Conference on Computer Vision*. Springer, 2022, pp. 691–707.
- [20] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [21] Z. Li, W. Wang, H. Li, E. Xie, C. Sima, T. Lu, Q. Yu, and J. Dai, “Befformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers,” *arXiv preprint arXiv:2203.17270*, 2022.
- [22] H. Li, C. Sima, J. Dai, W. Wang, L. Lu, H. Wang, E. Xie, Z. Li, H. Deng, H. Tian, *et al.*, “Delving into the devils of bird’s-eye-view perception: A review, evaluation and recipe,” *arXiv preprint arXiv:2209.05324*, 2022.
- [23] X. Jia, L. Chen, P. Wu, J. Zeng, J. Yan, H. Li, and Y. Qiao, “Towards capturing the temporal dynamics for trajectory prediction: a coarse-to-fine approach,” in *Conference on Robot Learning*. PMLR, 2023, pp. 910–920.
- [24] Y. Hu, J. Yang, L. Chen, K. Li, C. Sima, X. Zhu, S. Chai, S. Du, T. Lin, W. Wang, L. Lu, X. Jia, Q. Liu, J. Dai, Y. Qiao, and H. Li, “Planning-oriented autonomous driving,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- [25] P. Wu, L. Chen, H. Li, X. Jia, J. Yan, and Y. Qiao, “Policy pre-training for autonomous driving via self-supervised geometric modeling,” in *The Eleventh International Conference on Learning Representations*.
- [26] S. Hu, L. Chen, P. Wu, H. Li, J. Yan, and D. Tao, “St-p3: End-to-end vision-based autonomous driving via spatial-temporal feature learning,” in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVIII*. Springer, 2022, pp. 533–549.
- [27] P. Wu, X. Jia, L. Chen, J. Yan, H. Li, and Y. Qiao, “Trajectory-guided control prediction for end-to-end autonomous driving: A simple yet strong baseline,” *arXiv preprint arXiv:2206.08129*, 2022.
- [28] T. Wang, X. Zhu, J. Pang, and D. Lin, “Fcos3d: Fully convolutional one-stage monocular 3d object detection,” 2021.
- [29] Z. Wang, Z. Huang, J. Fu, N. Wang, and S. Liu, “Object as query: Equipping any 2d object detector with 3d detection ability,” *arXiv preprint arXiv:2301.02364*, 2023.
- [30] T. Yin, X. Zhou, and P. Krahenbuhl, “Center-based 3d object detection and tracking,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 11 784–11 793.
- [31] J. Huang, G. Huang, Z. Zhu, and D. Du, “Bevdet: High-performance multi-camera 3d object detection in bird-eye-view,” *arXiv preprint arXiv:2112.11790*, 2021.
- [32] L. Chen, C. Sima, Y. Li, Z. Zheng, J. Xu, X. Geng, H. Li, C. He, J. Shi, Y. Qiao, *et al.*, “Persformer: 3d lane detection via perspective transformer and the openlane benchmark,” *arXiv preprint arXiv:2203.11089*, 2022.
- [33] S. Huang, Z. Shen, Z. Huang, Z. Ding, J. Dai, J. Han, N. Wang, and S. Liu, “Anchor3dlane: Learning to regress 3d anchors for monocular 3d lane detection,” *arXiv preprint arXiv:2301.02371*, 2023.
- [34] B. Zhou and P. Krähenbühl, “Cross-view transformers for real-time map-view semantic segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 760–13 769.
- [35] Y. Zhou and O. Tuzel, “Voxelnet: End-to-end learning for point cloud based 3d object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4490–4499.

- [36] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, "Pointpillars: Fast encoders for object detection from point clouds," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 12 697–12 705.
- [37] Y. Yan, Y. Mao, and B. Li, "Second: Sparsely embedded convolutional detection," *Sensors*, vol. 18, no. 10, p. 3337, 2018.
- [38] J. Fu, G. Ren, Y. Chen, and S. Liu, "Improved pillar with fine-grained feature for 3d object detection," *arXiv preprint arXiv:2110.06049*, 2021.
- [39] C. Yang, Y. Chen, H. Tian, C. Tao, X. Zhu, Z. Zhang, G. Huang, H. Li, Y. Qiao, L. Lu, *et al.*, "Bevformer v2: Adapting modern image backbones to bird's-eye-view recognition via perspective supervision," *arXiv preprint arXiv:2211.10439*, 2022.
- [40] J. Phillion and S. Fidler, "Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d," 2020.
- [41] Y. Li, Z. Ge, G. Yu, J. Yang, Z. Wang, Y. Shi, J. Sun, and Z. Li, "Bevdepth: Acquisition of reliable depth for multi-view 3d object detection," *arXiv preprint arXiv:2206.10092*, 2022.
- [42] X. Chen, T. Zhang, Y. Wang, Y. Wang, and H. Zhao, "Futr3d: A unified sensor fusion framework for 3d detection," *arXiv preprint arXiv:2203.10642*, 2022.
- [43] S. Pang, D. Morris, and H. Radha, "Clocs: Camera-lidar object candidates fusion for 3d object detection," 2020.
- [44] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [45] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable DETR: Deformable transformers for end-to-end object detection," 2021.
- [46] T. Yin, X. Zhou, and P. Krähenbühl, "Center-based 3d object detection and tracking," 2021.