

OCBEV: Object-Centric BEV Transformer for Multi-View 3D Object Detection

Zhangyang Qi^{1†} Jiaqi Wang^{2*} Xiaoyang Wu¹ Hengshuang Zhao^{1*}

¹The University of Hong Kong, ²Shanghai Artificial Intelligence Laboratory

{zyqi, xywu3, hszhao}@cs.hku.hk, wangjiaqi@pjlab.org.cn

Abstract

Multi-view 3D object detection is becoming popular in autonomous driving due to its high effectiveness and low cost. Most of the current state-of-the-art detectors follow the query-based bird’s-eye-view (BEV) paradigm, which benefits from both BEV’s strong perception power and end-to-end pipeline. Despite achieving substantial progress, existing works model objects via globally leveraging temporal and spatial information of BEV features, resulting in problems when handling the challenging complex and dynamic autonomous driving scenarios. In this paper, we proposed an Object-Centric query-BEV detector OCBEV, which can carve the temporal and spatial cues of moving targets more effectively. OCBEV comprises three designs: 1) **Object Aligned Temporal Fusion** aligns the BEV feature based on ego-motion and estimated current locations of moving objects, leading to a precise instance-level feature fusion. 2) **Object Focused Multi-View Sampling** samples more 3D features from an adaptive local height ranges of objects for each scene to enrich foreground information. 3) **Object Informed Query Enhancement** replaces part of pre-defined decoder queries in common DETR-style decoders with positional features of objects on high-confidence locations, introducing more direct object positional priors. Extensive experimental evaluations are conducted on the challenging nuScenes dataset. Our approach achieves a state-of-the-art result, surpassing the traditional BEVFormer by 1.5 NDS points. Moreover, we have a faster convergence speed and only need half of the training iterations to get comparable performance, which further demonstrates its effectiveness.

1. Introduction

Multi-view 3D object detection is attracting increasing attention in autonomous driving. Compared to monocular detectors [49, 33, 36, 37], which lack depth information and have a limited field of view, multi-view detectors [39, 19, 20] provide a 360-degree view to fully ob-

*Corresponding Authors.

†The work is done during an internship at Shanghai AI Lab.

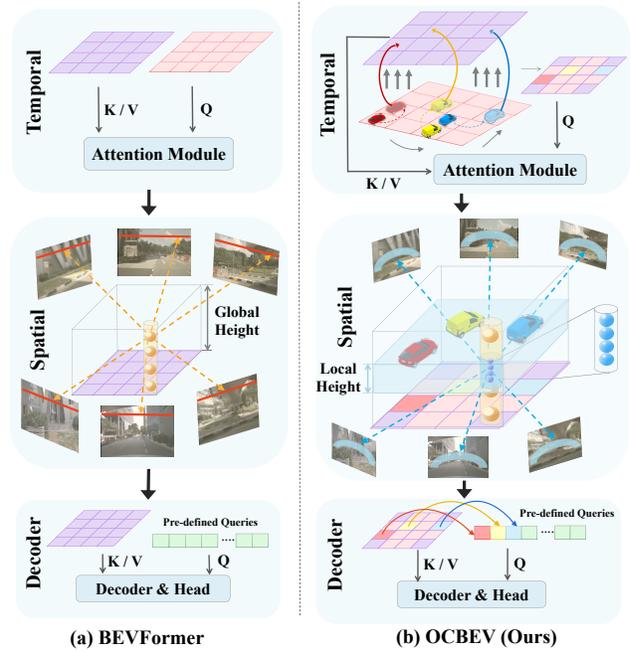


Figure 1: **Comparison between (a) BEVFormer and (b) our OCBEV.** The top shows the temporal modeling. BEVFormer fuses BEV features in two timestamps via global attention while we fuse them more precisely via ego-motion and object-motion. The middle part is the spatial extraction from multi-view image features. BEVFormer only samples points in the global height range evenly while we conduct it densely within an adaptive local height range with most objects for each scene. The bottom shows the decoder. BEVFormer adopts pre-defined queries, whereas we introduce object location priors to the decoder via replacing some queries with BEV features on high-confidence locations.

tain surrounding information. And unlike lidar-based methods [34, 16, 50], multi-view detectors have lower cost and rich semantic information such as traffic lights and lanes.

The multi-view detectors can be divided into *object-query-based methods* [39, 26, 23] and *bird’s-eye-view (BEV)-based solutions* [13, 19, 20]. The transformer series, like DETR3D [39], adopt a DETR-style [3] network, where one query is assigned to one object in 3D space. However,

since objects are not densely distributed, matching process become difficult, leading to low representation power. BEV-based methods employ the bird’s-eye view (BEV) as an intermediate state to model the positions of objects and can be applied to other tasks like lane segmentation.

BEV-based methods can be further categorized into *depth-BEV* [30, 13, 19] and *query-BEV methods* [20, 42], based on the way to construct BEV feature. Depth-BEV methods predict a latent depth distribution to build frusta, while the other use pre-defined BEV queries to encode information from image features. Query-BEV methods like BEVFormer [20] stand out again for their temporal and spatial encoding capabilities. However, the current temporal modeling is done globally, overlooking the fact that the scene is dynamic and objects are moving. For spatial exploitation, 3D reference points in the global height range are used for projection as in Fig. 1 (a), while this method falls short since the locations of moving objects is in a local height range area. Lastly, queries are all pre-defined for the decoder. It is known that object detection in 2D is difficult to optimize due to the scene’s sparsity, thus the matching process for queries and objects is even more challenging for 3D multi-view scenes.

We build an *object-centric query-based bird’s-eye view (BEV) detector* that incorporates the benefits of both object-query methods and query-BEV methods. *First*, to conduct temporal modeling, the query-BEV detector BEVFormer [20] uses deformable attention between historical and current BEV queries, as in Fig. 1(a) top. We propose a module called *Object Aligned Temporal Fusion* that predicts the current locations of objects by using their historical locations and velocity. This allows for better feature representation when fused with current BEV features, as shown in the top of Fig. 1(b). *Second*, to execute spatial exploitation, BEVFormer [20] only uses orange 3D reference points evenly distributed to the global height range, as in Fig. 1(a) middle. The top point hits the top of the image features like the orange line in the image features, which belong to the background. Contrarily, we develop *Object Focused Multi-View Sampling* module that predicts a local height for each scenario adaptively where most moving objects are located at this height, and we densely sample 3D points in this area to ensure more of the points hit objects in the 2D image feature like the blue bar area in the image feature, as demonstrated in the middle of Fig. 1(b). *Third*, the decoder of BEVFormer [20] uses pre-defined queries as in Fig. 1(a) bottom, which makes the object matching queries and objects more challenging. On the contrary, we design an *Object Informed Query Enhancement* module to replace part of the decoder queries with high-confidence locations of objects via a centerness heatmap supervision head after the encoder, as illustrated in the bottom of Fig. 1(b) which leads to a faster convergence speed for the whole network.

We name the whole algorithm as Object-Centric query-BEV detector OCBEV, and conduct extensive experimental evaluations on the challenging nuScenes [1] benchmark to validate its effectiveness. The proposed OCBEV achieves a state-of-the-art result, surpassing the classical BEVFormer by 1.5 NDS point. Further, it converges faster and only needs half of the training iterations to get comparable performance. Our code and models will be made publicly available. In summary, our contributions are as follows:

- We find that object-centric modeling is the key to the query-BEV 3D detectors. Existing methods only consider extracting temporal and spatial aspects globally, whereas we conduct in-depth design to overcome it.
- We present a powerful BEV detector named OCBEV, which consists of three newly proposed novel modules that can better conduct temporal modeling, execute spatial exploitation, and construct strong decoders.
- We evaluate our method on the challenging nuScenes dataset. Our OCBEV method can achieve state-of-the-art results in the vision-based methods. And our network converges faster and cut training iterations in half to get comparable results with existing methods.

2. Related Work

Multi-view vision-based 3D detection. For multi-view camera-only object detection, various object-query-based algorithms have been developed. DETR3D [39] follows the 2D detector DETR [3] and replaces the 2D queries with the 3D queries. PETR [3] develops position embedding transformation. Polarformer [15] and Polar-DETR [4] advocate exploiting the Polar coordinate system. Sparse4D [23] uses hierarchy feature fusion for sparse 3D detection.

Bird’s-eye-view (BEV) is a new trend for multi-view perception tasks. BEV is not originally from 3D detection tasks, and early works consider transforming perspective features to BEV features [32, 28] like the pseudo-lidar [38, 40]. Current BEV methods mainly extend to the 360-degree sensor to build the surrounding environment [30, 2, 48, 47]. LSS [30] predicts an implicit depth distribution to build the frusta from six cameras and construct the BEV feature by splatting frusta. While CVT [48] bridges perspective-view and BEV-view features by a camera-aware positional encoding module and dense cross-attention.

BEV is widely used in 3D object detection, which is originally a kind of lidar-based representation [10, 27]. TiG-BEV [14] and CMT [41] both use BEV feature to combine lidar-based methods and vision-based methods. For vision-based methods only, according to the way to build the BEV feature, these methods can be divided into depth-BEV-based methods [13, 18, 19, 35] and query-BEV-based methods [20, 42]. For depth-BEV-based methods, BEVDet [13] gives the paradigm of BEV 3D detection from

LSS [30]. BEVstereo [18] and Bevdepth [19] extend it to stereo-view input and depth supervision, respectively. Query-BEV-based methods like BEVFormer [20] do not need to build the latent depth but extract information from the perspective-view feature directly. Therefore, we focus on designing a query-BEV-based 3D detector. More details of *object-query-based*, *BEV-query-based*, and *BEV-depth-based methods* will be explained in the Appendix A.1.

Temporal and spatial exploration in 3D detection. Temporal modeling and spatial exploitation are extremely important in BEV detection. The former means modeling the history information since autonomous driving can be seen as a continuous video and history will give clues to the current frame. Actually, many 3D detectors have considered the temporal modeling [19, 18, 20]. There are also some works that only utilize temporal information. BEVDet4d [12], DETR4D [26] and PETR v2 [25] all do temporal modeling based on BEVDet [13], DETR3D [39], and PETR [24], respectively. Some works only study the effect of the temporal information, TWW [29] uses long and short temporal image feature fusion to get the BEV feature. However, all of them only consider BEV temporal modeling at the global level while we consider aligning moving objects of their displacement. Spatial information is to make the BEV feature extract information from the perspective view. And as we mentioned before, both depth-based [13, 18, 19, 35] and query-based methods [20, 42] try to bridge these two views. But they meet the same problem, *i.e.*, the spatial exploitation only bridges in a global range while most objects are densely in a local range. Some works like BEV-SAN [5] and OA-BEV [6] also try to pay more attention to these moving objects. We propose to predict an adaptive local height range for each scenario where most moving objects are located at this height range for better feature learning.

Decoder enhancement in object detection. Traditional CNN-based two-stage detectors [9, 8, 21] first give rough proposals, followed by precise predictions. In DETR-based methods [3, 52, 43, 17, 46], the encoder-decoder structure can be regarded as a two-stage pipeline. Two-stage deformable DETR [52] feeds the high score encoder embeddings to the decoder. Efficient DETR [43] enhances decoder queries by selecting top-k positions from the encoder’s output. Apart from giving positional clues to the decoder, DN-DETR [17] and DINO [46] add noisy labels and bounding boxes to the decoder for speeding up. In query-based 3D object detection, BEVFormer v2 [42] adopts a perspective view head after the image features to provide reference points and positional embeddings to the decoder. While we only need to provide positional instead of semantic information, so we add a centerness heatmap supervision which not only gives the decoder enough positional hints but makes the network end-to-end trainable as well. More details about decoders will be shown in the Appendix A.2.

3. Object-Centric BEV Transformer

3.1. Overall Architecture

The overall framework of our OCBEV is illustrated in Fig. 2. As most query-BEV methods like BEVFormer [20] and BEVFormer v2 [42], we start by encoding the pre-defined BEV queries with temporal modeling over historical BEV queries and spatial exploitation on perspective-view image features. Sec. 3.2 introduces Object Aligned Temporal Fusion, which aligns and fuses historical BEV, considering not only the background but also moving objects. And it can be divided into *ego-motion temporal fusion* and *object-motion temporal fusion* in consideration of ego-motion and object-motion, respectively. Besides, we utilize attention between these two BEV-feature timestamps after the fusion. For spatial information, we introduce Object Focused Multi-View Sampling in Sec. 3.3, which is based on spatial-cross attention in the BEVFormer [20]. We predict an adaptive local height within an object-focused area and sample the corresponding dense 3D points. Afterward, we project them to 2D image features and form a bar area where moving objects are densely distributed. After encoding, we construct Object Informed Query Enhancement for the decoder in Sec. 3.4, which includes a heatmap head that predicts the object centerness, and exchanges some pre-defined decoder queries with high-confidence locations that give positional hints. Finally, a detection head like Center-Head [44] is implemented to predict the final results.

3.2. Object Aligned Temporal Fusion

Recent BEV-based detectors incorporate temporal information by fusing the BEV features from previous timestamps, resulting in significant performance improvements. However, the existing modeling process meets problems when dealing with complex and dynamic autonomous driving scenarios and is incapable of discovering the dynamic characteristics. As an example, the depth-BEV series approaches utilize two timestamps as a stereo view to obtain the depth information, leaving the capturing of the moving status of objects untouched. Similarly, as representatives of query-BEV methods, BEVFormer [20] attempts to utilize deformable attention to capture history information, but also ignore the geometry modeling for moving objects.

Existing BEV temporal fusion strategies fail to address two critical problems: i) Fail to track the ego-motion information of the host vehicle, which is indispensable for modeling the driving scenario; ii) Fail to leverage geometry similarity, which is crucial to handle fast-moving vehicles. To tackle these challenges, we propose a novel solution Object Aligned Temporal Fusion. Specifically, we perform ego-motion temporal fusion and object-motion temporal fusion, resulting in a more comprehensive and effective fusion of temporal information.

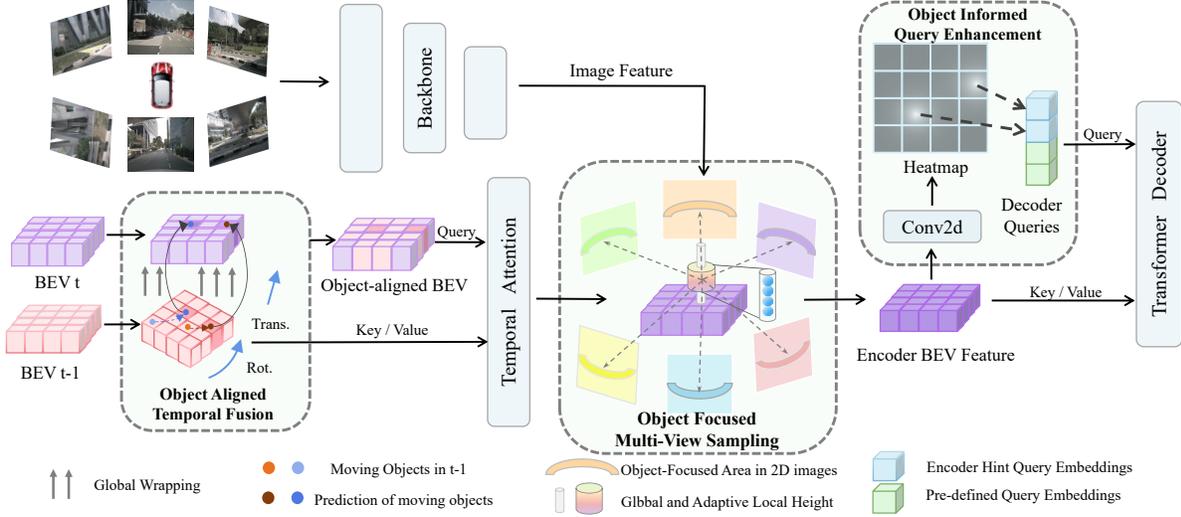


Figure 2: **Overall architecture of OCBEV.** (a) *Object Aligned Temporal Fusion*: warps BEV $t - 1$ to BEV t by ego motion transformation, then fuses moving objects by predicting their current locations via their velocity. (b) *Object Focused Multi-View Sampling*: pre-defines points in the 3D space and then projects them to perspective-view image features. Apart from sampling points in the global height range, we propose an adaptive local object-dense height with most objects for each autonomous driving scene and sample points densely in it. (c) *Object Informed Query Enhancement*: A heatmap supervision is added after the encoder output and replaces part of pre-defined decoder queries with high-confidence locations of objects.

Ego-motion temporal fusion. BEV ego-motion temporal fusion differs from points alignment in 3D space. To fuse a current query feature with its corresponding feature from a previous timestamp, we have to perform a bidirectional transformation to determine the mapping relationship which ensures precise alignment and fusion of features between different timestamps. For convenience, we define the previous timestamp as t' and the current timestamp as t . We also denote the BEV queries as $Q_t \in \mathbb{R}^{C_q \times H_q \times W_q}$ and its flatten version as $\bar{Q}_t \in \mathbb{R}^{C_q \times (H_q W_q)}$. Fig. 3 (a) demonstrates the process of aligning ego-motion. We aim to determine the mapping $\mathcal{P} = \{(i, j)\}_{t', t}$, where $\bar{Q}_{t'}(i) \rightarrow \bar{Q}_t(j)$. The set \mathcal{P} contains N_p pairs, indicating the number of overlapping queries. The indices of overlapping queries in timestamp t' and t are denoted as I and $J \in \mathbb{R}^{1 \times N_p}$, respectively.

The main idea to obtain I is to transform $Q_{t'}$ to the t reference frame for selecting the queries overlapped with Q_t . We simplify the problem by assuming that vehicles are moving on a plane, which means that the 3×3 rotation matrix $\mathcal{R}_{t' \rightarrow t}$ can be reduced to a simplified $\gamma_{t' \rightarrow t}$ representing the yaw angle along the Z axis. The 3×1 translation vector is denoted as $\mathcal{T}_{t' \rightarrow t}$, and we only use the xy translation because the motion is on a plane. The selected indices via the transformation from $Q_{t'}$ to Q_t are denoted as I and can be written as Eq. 1. Similarly, we obtain J as well.

$$\begin{aligned} I &= [g_t(g_r(Q_{t'}, \gamma_{t' \rightarrow t}), \mathcal{T}_{t' \rightarrow t}) \in Q_t] \\ J &= [g_t(g_r(Q_t, \gamma_{t \rightarrow t'}), \mathcal{T}_{t \rightarrow t'}) \in Q_{t'}] \end{aligned} \quad (1)$$

here g_t represents translation operation, whereas g_r represents rotation operation. For the BEV queries fusion, we replace overlapping BEV queries $\bar{Q}_{t'}(J)$ with corresponding previous BEV queries $\bar{Q}_{t'}(I)$ as $\bar{Q}_{t'}(I) \rightarrow \bar{Q}_t(J)$, and add the replaced BEV queries and current BEV queries Q_t . All alignment and fusion operations are expressed as:

$$Q_a^{ego} = f_a^{ego}(Q_{t'}, Q_t) + Q_t \quad (2)$$

Here f_a^{ego} is the function to align the BEV feature by ego-motion bidirectional transformation as mentioned above, we fuse it with current BEV embeddings Q_t by adding them together and Q_a^{ego} is the fused BEV queries.

Object-motion temporal fusion. Apart from the ego-motion temporal fusion, we also align and fuse moving objects as Fig. 3 (b). Suppose there are N_m objects in timestamp t' . $p_{t'}^m \in \mathbb{R}^{N_m \times 2}$ is their real xy coordinates with the corresponding indices $I_{t'}^m$ on the $Q_{t'}$. We predict velocity $v_{t'}^m \in \mathbb{R}^{N_m \times 2}$ to find the corresponding indexes J_t^m in Q_t and replace them as $\bar{Q}_{t'}(I_{t'}^m) \rightarrow \bar{Q}_t(J_t^m)$. Specifically, we first predict the locations after the motion in t reference frame as $p_t^m = p_{t'}^m + v_{t'}^m(t - t')$. Then convert coordinates into indices as I_t^m . To get its corresponding after-motion indices J_t^m in Q_t , we rotate and translate $\bar{Q}_{t'}(I_{t'}^m)$ as:

$$J_t^m = [g_t(g_r(\bar{Q}_{t'}(I_{t'}^m), \gamma_{t' \rightarrow t}), \mathcal{T}_{t \rightarrow t'}) \in Q_t] \quad (3)$$

I_t^m and J_t^m form the object-level mapping $\bar{Q}_{t'}(I_{t'}^m) \rightarrow \bar{Q}_t(J_t^m)$. Having the mapping relationship, we align and

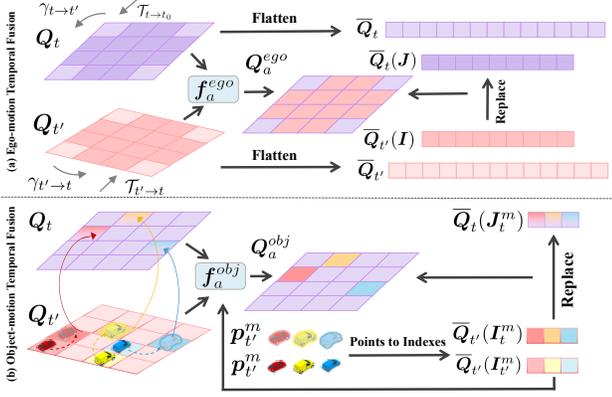


Figure 3: **Object Aligned Temporal Fusion.** (a) *Ego-motion temporal fusion* finds overlapping indices of adjacent BEV features by the ego-motion globally. Then we flatten BEV queries and replace current embeddings with corresponding previous embeddings. Original embeddings are added for fusion. (b) *Object-motion temporal fusion* is for moving objects. We predict their current locations and indices by velocity and use the above ego-motion fusion way to find and replace their corresponding current queries.

fuse the moving objects as before. The whole process is:

$$Q_a^{obj} = f_a^{obj}(Q_{t'}, Q_t, p_{t'}^m) \quad (4)$$

Here f_a^{obj} is the object-motion temporal alignment and fusion function locally between two timestamps. Apart from previous $Q_{t'}$ and current Q_t BEV queries, we have to input the original locations of objects $p_{t'}^m$. Overall, the Object Aligned Temporal Fusion can be written as:

$$\begin{aligned} Q_a &= Q_a^{ego} + Q_a^{obj} \\ &= f_a^{ego}(Q_{t'}, Q_t) + f_a^{obj}(Q_{t'}, Q_t, p_{t'}^m) + Q_t \end{aligned} \quad (5)$$

where Q_a is the aligned and fused BEV embeddings between t' and t timestamps. This operation contains ego-motion and object-motion temporal fusion.

After the alignment and fusion, to fully utilize the historical information, we adopt the deformable DETR attention [52] between aligned BEV feature Q_a and previous BEV feature $Q_{t'}$ as well. After finishing the temporal modeling, the BEV feature will be fed into the spatial module to extract information from perspective-view image features.

3.3. Object Focused Multi-View Sampling

Spatial exploitation for BEV queries is to extract information from perspective-view image features. All query-based methods depend on the deformable DETR [52], which makes the image features as keys and values. They pre-define 3D reference points for the BEV queries and project them to the 2D image features to sample keys and values. The question is whether the height of predefined

3D reference points is in the overall detection height range. Nevertheless, most vehicles are located at a local height. So our Object Focused Multi-View Sampling will predict an adaptive local height range that contains most objects with consideration of vehicles such as bus and truck and the height range can be adaptive for each unique scene.

We follow the spatial attention module based on the deformable attention [52] in the BEVFormer [20], which defines N_h pillar-like reference points along the z -axis for one specific BEV query $Q(u, v)$. The reference points are $p_{3d} \in \mathbb{R}^{N_h \times 3}$. All reference points share the same xy coordinates, which are the center of the BEV embedding corresponding area. Regarding the z coordinate, these points are distributed along the z -axis in the range Z_h evenly. So we denote $p_{3d}(Z_h)$ which means that the z coordinate of reference points p_{3d} are distributed evenly in the Z_h height range. Thus the spatial attention is shown below:

$$\text{SpaA}(Q_t, Z_h) = \text{DAttn}(Q_t, F_t, g_p(p_{3d}(Z_h))) \quad (6)$$

where SpaA is the spatial attention and inputs are BEV queries Q_t and height range Z_h . DAttn is the deformable attention [52]. Q_t is the queries, and multi-view image features F_t are the keys and values. g_p is the projection function for projecting 3D reference points to the 2D multi-view image features through extrinsic and intrinsic matrices.

In autonomous driving scenarios, most objects, including cars, buses, and pedestrians, are in a local height range. For the nuScenes dataset [1], the detection global range $Z_{h,g}$ is from $-5m$ to $3m$. Apart from the global range, we notice that most moving objects are concentrated in the height range of $[-2m, 2m]$ [5] because the camera sensor is usually located on the top of the data collection vehicle, so we define the local height $Z_{h,l}$ from $-2m$ to $2m$. We also predict an adaptive offset Δh for each scene to handle scenes with high vehicles like buses and trucks which is also common in autonomous driving scene. The adaptive local height is denoted as $\hat{Z}_{h,l} = Z_{h,l} + \Delta h$.

As shown in the middle right part of Fig. 1 and Fig. 2, these 3D points in the local height range are projected to the 2D image features and form bar areas that contain dense objects. So we sample 3D reference points in not only the global height range but also densely in the local height range. In this way, the spatial attention inside our Object Focused Multi-View Sampling is shown below:

$$\text{OFSpaA}(Q_t, Z_h) = \text{SpaA}(Q_t, Z_{h,g}) + \text{SpaA}(Q_t, \hat{Z}_{h,l}) \quad (7)$$

Here, OFSpaA means the spatial attention by Object Focused Multi-View Sampling, which is the sum of spatial attention with 3D reference points in global height range and adaptive local height range to focus the moving objects and make BEV's extraction power from image features stronger.

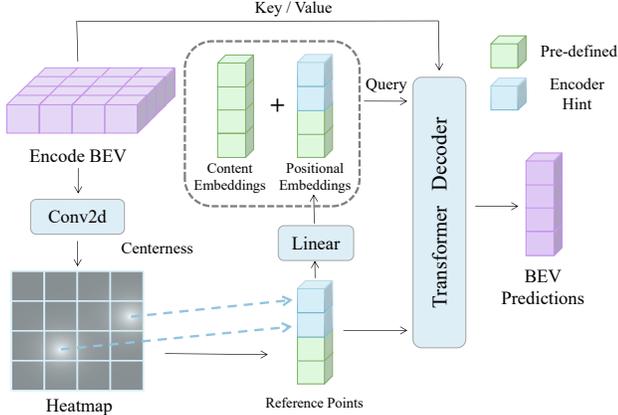


Figure 4: **Object Informed Query Enhancement.** Apart from feeding the BEV feature as keys and values, we use a heatmap supervision after the encoder BEV feature to get a confidence heatmap and replace the pre-defined reference points with high-confidence locations partially. The queries are divided into content embeddings that are predefined and positional embeddings that are from the reference points.

3.4. Object Informed Query Enhancement

As for the decoder, as mentioned above, DETR [3] style pre-defines queries randomly, which leads to difficulty in convergence and will be shared in all scenes. Inspired by two-stage CNN-based methods, we use encoder output to deliver a positional hint to the decoder. Our decoder is shown in Fig. 4, which divides the queries to content and positional embeddings following deformable DETR [52].

Unlike two-stage deformable DETR [52] that directly adds a detection head after the encoder, we add a heatmap head based on the encoder output. Since the outdoor 3D scene is within a wide range, we can focus more on the position in the BEV plane, rather than other properties like box size, velocity, attribute, etc. Here we introduce object centerness as defined in the FCOS3D [36] and CenterNet [49]. The centerness is represented by the 2D Gaussian distribution, which ranges from zero to one, and we utilize the binary cross entropy loss to conduct the training.

Apart from extra heatmap supervision, we replace some predefined queries with the BEV features on high-confidence locations. To be specific, for the decoder transformer, the keys and values are output features of the BEV encoder. We divide the queries into content and positional embeddings. The origin deformable DETR [52] initializes content embeddings positional embeddings and reference points randomly. Whereas in our framework, we use BEV points with high-confidence scores as reference points to partially replace the pre-defined part of reference points and positional embeddings through a linear layer whose input is replaced reference points, which brings beneficial positional priors to make the network converge quickly.

4. Experiments

4.1. Implementation Details

Dataset and metrics. To demonstrate the effectiveness of the proposed method, we conduct experimental evaluations on the challenging nuScenes 3D detection benchmark [1]. The dataset consists of 1000 scenes, each of which is about 20 seconds, and the sample rate is 2 HZ (40 annotation samples/scene). Each sample contains six images from six 360-degree cameras. It is split into 700/150/150 videos (28k/6k/6k samples) for training, validation, and testing, respectively. Ten categories, including car, bus, and pedestrian are adopted for evaluation. For the 3D object detection metrics, mean average precision (mAP) is computed over four thresholds using center distance on the ground plane. There are five other types of true positive metrics (TP metrics), including mean Average Translation Error (mATE), mean Average Scale Error (mASE), mean Average Orientation Error (mAOE), mean Average Velocity Error (mAVE), and mean Average Attribute Error (mAAE). Furthermore, the nuScenes detection score (NDS) is a comprehensive and the most important metric that combines all aspects above.

Experimental settings. The input image size is 1600×900 as default. Following previous practice [20, 39, 26], our model adopts the FCOS3D [36] pre-trained ResNet-101 [11] as backbone and initialized weights. After the backbone, we use a FPN [21] structure as the neck with the dimension of 256 and the size of $1/16, 1/32, 1/64$. As for the transformer layer, we utilize six layers for temporal, spatial, and decoder attention. For the BEV setting, we set the size of BEV queries as 300×300 (A minor variable, 200: 53.1% NDS vs. 300: 53.2% NDS), the perception range as $[-51.2m, 51.2m]$ for the XY plane, and the height range as $[-5m, 3m]$. For the training setting, we train our models for 12 and 24 epochs. The initial learning rate is set as 2×10^{-4} and a cosine learning rate decay policy is adopted.

4.2. Main Results

We compare our experimental results with other vision-based 3D detection methods on the nuScenes validation set, especially those query-based vision methods. The results are shown in Table 1. Compared to the current *state-of-the-art* query-based method BEVFormer [20], OCBEV outperforms it over 1.5 points on the validation set. Moreover, OCBEV outperforms DETR4D [26] over 2.3 points and PETRv2 [25] over 0.8 points with the same backbone and training schedule. For the depth-based methods, OCBEV is 1.7 points higher than the BEVDet4d [12]. All these results prove that focusing on objects can bring huge improvement, and our proposed OCBEV is a strong 3D outdoor detector and push the previous SOTA to 53.2%.

Moreover, OCBEV can converge faster than existing methods. With only half of the training schedule (12

| Style | Method | Backbone | Epoch | NDS \uparrow | mAP \uparrow | mATE \downarrow | mASE \downarrow | mAOE \downarrow | mAVE \downarrow | mAAE \downarrow |
|-------------|-------------------------|------------|-------|----------------|----------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| Monocular | CenterNet [49] | DLA [45] | - | 0.328 | 0.306 | 0.716 | 0.264 | 0.609 | 1.426 | 0.658 |
| | FCOS3D [36] | ResNet-101 | - | 0.415 | 0.343 | 0.725 | 0.263 | 0.422 | 1.292 | 0.153 |
| | PGD [37] | ResNet-101 | 24 | 0.428 | 0.369 | 0.683 | 0.260 | 0.439 | 1.268 | 0.185 |
| Depth-based | BEVDet \dagger [13] | Swin-T | 24 | 0.417 | 0.349 | 0.637 | 0.269 | 0.490 | 0.914 | 0.268 |
| | BEVDet4D \dagger [37] | Swin-B | 24 | 0.515 | 0.396 | 0.619 | 0.260 | 0.361 | 0.399 | 0.189 |
| Query-based | DETR3D \dagger [39] | ResNet-101 | 24 | 0.425 | 0.346 | 0.773 | 0.268 | 0.383 | 0.842 | 0.216 |
| | DETR4D \dagger [26] | ResNet-101 | 24 | 0.509 | 0.422 | 0.688 | 0.269 | 0.388 | 0.496 | 0.184 |
| | PETrv2 \dagger [25] | ResNet-101 | 24 | 0.524 | 0.421 | 0.681 | 0.267 | 0.357 | 0.377 | 0.186 |
| Query-BEV | BEVFormer-S [20] | ResNet-101 | 24 | 0.448 | 0.375 | 0.725 | 0.272 | 0.391 | 0.802 | 0.200 |
| | BEVFormer [20] | ResNet-101 | 24 | 0.517 | 0.416 | 0.673 | 0.274 | 0.372 | 0.394 | 0.198 |
| | OCBEV (Ours) | ResNet-101 | 12 | 0.523 | 0.408 | 0.633 | 0.270 | 0.377 | 0.333 | 0.194 |
| | OCBEV (Ours) | ResNet-101 | 24 | 0.532 | 0.417 | 0.629 | 0.273 | 0.339 | 0.342 | 0.187 |

Table 1: **Results comparison on the nuScenes detection validation set.** We conduct a comprehensive comparison of our OCBEV with four styles of vision-based 3D detectors. Notably, our model achieves SOTA and even outperforms latest models with data augmentation strategy CBGS [51] as indicated by \dagger . Also it achieves superior performance compared to BEVFormer [20] NDS result even with only 12 epochs of training, and extends this performance gap with a standard 24 epochs training schedule. Meanwhile it performs best on kinestate metrics such as mAOE and mAVE, underscoring our Object-Centric approach has the capability to model moving objects in the sparse autonomous driving scene.

| Backbone | Global | Num. Objects | NDS \uparrow | mAP \uparrow | mATE \downarrow | mASE \downarrow | mAOE \downarrow | mAVE \downarrow | mAAE \downarrow |
|------------|--------------|--------------|----------------|----------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| ResNet-101 | - | - | 0.371 | 0.316 | 0.798 | 0.292 | 0.605 | 0.919 | 0.252 |
| ResNet-101 | \checkmark | - | 0.401 | 0.308 | 0.801 | 0.288 | 0.555 | 0.669 | 0.231 |
| ResNet-101 | \checkmark | 10 | 0.403 | 0.303 | 0.826 | 0.289 | 0.533 | 0.622 | 0.220 |
| ResNet-101 | \checkmark | 30 | 0.413 | 0.323 | 0.801 | 0.285 | 0.533 | 0.640 | 0.232 |
| ResNet-101 | \checkmark | 50 | 0.409 | 0.318 | 0.809 | 0.293 | 0.524 | 0.647 | 0.230 |

(a) **Object Aligned Temporal Fusion.** We first enable ego-motion fusion and then object-motion fusion by increasing the number of aligned objects. Results indicate that Object Aligned Temporal Fusion works effectively with ego-motion fusion and achieves the best performance with 30 aligned objects.

| Backbone | Height Range | NDS \uparrow | mAP \uparrow |
|------------|--------------|----------------|----------------|
| ResNet-101 | - | 0.371 | 0.316 |
| ResNet-101 | [-5m, -2m] | 0.412 | 0.321 |
| ResNet-101 | [-2m, 2m] | 0.421 | 0.325 |
| ResNet-101 | [2m, 5m] | 0.410 | 0.325 |
| ResNet-101 | adaptive | 0.430 | 0.324 |

(b) **Object Focused Multi-View Sampling.** Our proposed model benefits significantly from the multi-view sampling strategy, particularly within the usual object height range of -2 meters to 2 meters. Moreover adaptive local height takes performance a step further.

| Backbone | Heatmap | Num. Points | NDS \uparrow | mAP \uparrow |
|------------|--------------|-------------|----------------|----------------|
| ResNet-101 | - | - | 0.371 | 0.316 |
| ResNet-101 | \checkmark | - | 0.414 | 0.292 |
| ResNet-101 | \checkmark | 30 | 0.412 | 0.327 |
| ResNet-101 | \checkmark | 50 | 0.422 | 0.338 |
| ResNet-101 | \checkmark | 70 | 0.416 | 0.320 |

(c) **Object Informed Query Enhancement.** Our model exhibits improved performance with the additional heatmap supervision. Further, the replacement of predefined queries of points with object positional information results in a further boost in performance until the replacement points number reaches 50.

Table 2: **Ablation study.** We conducted ablation experiments on a lighting version OFBEV model trained with a 12 epochs which remove temporal attention and reduce the spatial attention layer from 6 to 3. Our results indicate that each object-centric design has a positive impact on the overall performance of the model. More results are in the supplementary material.

epochs), the NDS of OCBEV is 0.6 points higher than BEVFormer [20] with 24 epochs (52.3% NDS vs. 51.7% NDS). And our 12-epoch result has already gone beyond most methods, which means the encoder positional hint to the decoder can help the network to converge faster largely. The benefits of the proposed Object-Centric modules are also revealed in the kinematic metrics. Metrics relying on the kinestate such as orientation (mAOE: 33.9%) and velocity (mAVE: 33.3% and 34.2%) outperform BEVFormer [20] and get best results among all existing methods. It also gets better results on the translation (mATE: 62.9%). This demonstrates that our method pays more attention to mov-

ing objects and can naturally improve overall performance. Some qualitative results are shown in Fig. 5 and show that our model can give more precise 3D bounding boxes.

4.3. Ablation Studies

To verify the effectiveness of the proposed three Object-Centric modules, we conduct ablation experiments on a lighting version OFBEV model trained with a 12 epochs training schedule and reduce the attention layer from 6 to 3 to save the memory. Results are reported in Table 2.

Object Aligned Temporal Fusion. To study the effect of the Object Aligned Temporal Fusion module, we remove the ego-motion and object ego-motion aligned fu-

| Backbone | Temporal | Spatial | Decoder | NDS↑ | mAP↑ |
|------------|----------|---------|---------|--------------|--------------|
| ResNet-101 | | | | 0.371 | 0.316 |
| ResNet-101 | ✓ | | | 0.413 | 0.323 |
| ResNet-101 | | ✓ | | 0.430 | 0.324 |
| ResNet-101 | | | ✓ | 0.422 | 0.338 |
| ResNet-101 | ✓ | ✓ | | 0.439 | 0.323 |
| ResNet-101 | ✓ | | ✓ | 0.439 | 0.339 |
| ResNet-101 | | ✓ | ✓ | 0.444 | 0.335 |
| ResNet-101 | ✓ | ✓ | ✓ | 0.448 | 0.333 |

Table 3: **Combination of object-centric components.** Temporal, Spatial and Decoder represents *Object Aligned Temporal Fusion*, *Object Focused Multi-View Sampling* and *Object Informed Query Enhancement* respectively. Performance increases with the number of modules increase.

sion. Note that we also remove the temporal attention between the current and previous BEV feature. Then, we add ego-motion and different numbers of objects for object ego-motion aligned fusion step by step. As shown in Table 2a, by adding the global alignment, the NDS increase by 0.5% points. Regarding the object ego-motion aligned fusion module, the gain increases from modeling 10 objects to 30 objects and drops from 30 to 50 objects. We owe this phenomenon to that one scene may not have so many objects. Thus in this setting, modeling 30 objects can get the best results and outperform by 1.7% points than that without object-focused alignment fusion.

Object Focused Multi-View Sampling. In the nuScenes validation set, we use the point cloud range of $[-5m, 3m]$ as the global height. In the spatial attention, apart from sample points in the global height range. We also sample points in the following local height ranges $[-5m, -2m]$, $[-2m, 2m]$ and $[2m, 3m]$. These three ranges are the bottom, middle, and top of the 3D scene, respectively. Also, we study the effect of making the height adaptive. The results are shown in Table 2b, which indicate that paying attention to object-dense local height $[-2m, 2m]$ (NDS: +5.0%) can have the most gain compared to other less-object-dense height ranges. While the adaptive height range will gain the most (NDS: +5.9%) and be much higher than other ranges.

Object Informed Query Enhancement. In order to validate the effectiveness of object-centric encoder output for the decoder, we remove the heatmap supervision and all encoder positional hints as the baseline. We then add the supervision and replace the predefined positional embeddings step by step. Results in Table 2c show that by adding heatmap supervision, the NDS increases by 4.3 points and achieves 41.4%. Moreover, replacing pre-defined reference points and positional embeddings will gain further improvements and reach a peak with a replacement points number of 50 and finally lead to a result of 42.2% NDS. Fig. 6 shows the visualization of our heatmaps from four different scenes.

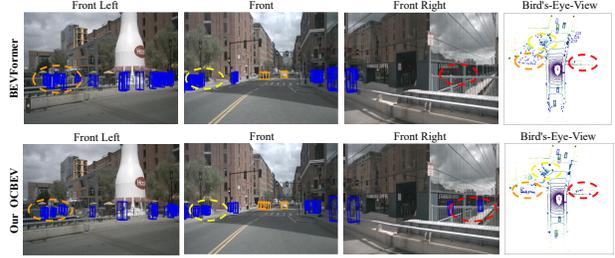


Figure 5: **Qualitive Results compared with BEVFormer.** Our model predicts more accurately than BEVFormer which successfully detects the pedestrian in the front right view and less false positive cases in the front right view.

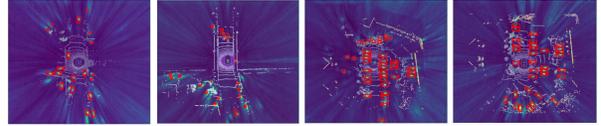


Figure 6: **Visualization of the heatmap.** It shows heatmaps of two scenes with ground truth boxes where red points are pedestrians and bounding boxes are cars. Bright spots (high confidence locations) correspond to the absolute position of objects, which gives the decoder a precise positional hint.

The high confidence position shown as bright spots corresponds to the absolute position of objects, which gives the decoder a precise positional hint. All prove the effectiveness of our Object Informed Query Enhancement.

Combination of object-centric components. In Table 3, we ablate the combination of three object-centric modules to confirm their contributions to the final results. The baseline is the same as mentioned before, we remove all our modules and the temporal attention. *Temp* means the Object Aligned Temporal Fusion, *Spatial* denotes Object-Focused Temporal Alignment, and *Decoder* represents Object Informed Query Enhancement for the decoder. Obviously, the performance will increase along with the number of modules increase. Finally, the complete version with all three modules outperforms 7.7 points compared to the baseline and reaches 44.8% NDS with only 12 epochs of training.

5. Conclusion

Query-based BEV paradigm benefits from both BEV’s perception power and end-to-end pipeline. Most 3D BEV detectors ignore moving objects in the scene and leverage temporal and spatial information globally. In this work, we propose a novel Object-Centric query-based BEV detector OCBEV. Three modules: Object Aligned Temporal Fusion, Object Focused Multi-View Sampling and Object Informed Query Enhancement focus on instance-level temporal modeling, spatial exploitation and decoder respectively. Our results achieve a state-of-the-art result on nuScenes benchmark and have a faster convergence speed as it only needs half training iterations to get comparable performance.

OCBEV: Object-Centric BEV Transformer for Multi-View 3D Object Detection

(Supplementary Material)

Content

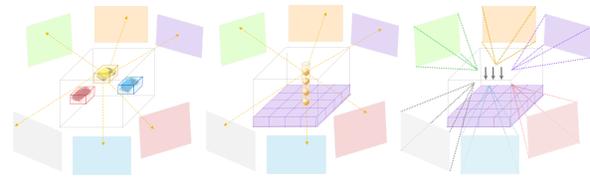
| | |
|---|----|
| A. Additional Related Work | 9 |
| A.1 Three Types of Multi-view 3D Object Detection | 9 |
| A.2 One-stage and Two-stage DETR Style Decoder | 10 |
| B. Implementation Details | 10 |
| B.1 Glossary of Notation | 10 |
| B.2 Network Architecture | 10 |
| B.3 Training Setting | 12 |
| B.4 Dataset Pre-processing | 12 |
| C. Extended Experiment Results | 12 |
| C.1. Analysis of Categories in Main Results | 12 |
| C.2. Additional Ablation Studies | 13 |
| D. Visualization | 13 |

A. Additional Related Work

A.1. Three Types of Multi-view 3D Object Detection

Here we will give a more detailed explanation on the classification of multi-view vision-based detectors. We first divided the multi-view vision-based detector into object-(query)-based methods [39] and BEV-based methods, and the classification is based on the format of features. Object-based approaches also called object-query-based methods use an object in 3D space as a query, and obtain temporal and spatial information by querying historical queries and 2D images as shown in Fig. S1. We can know that this type of method focuses on the object. However, due to the sparsity and scattered distribution of objects in the scene, accomplishing a meticulous one-to-one correlation between the physical entities and the related queries can be a daunting and demanding task. So BEV-based methods are proposed as BEV has strong representation capability.

Similar to the object-query-based methods, we also need to consider BEV to obtain temporal and spatial information. According to the way to extract information from perspective-view image features. Multi-view vision-based



(a) Object-query-based (b) BEV-query-based (c) BEV-depth-based

Figure S1: Three types of Multi-view 3D Detectors. The (a) Object-query-based and (b) BEV-query-based together form the query-based methods. The (b) BEV-query-based and (c) BEV-depth-based form the BEV-based methods.

detectors are divided into query-based methods and (BEV-)depth-based methods. BEV-depth-based methods [13] predict a latent depth to build the multi-view frusta and splat them to form the BEV while both object-query-based methods and BEV-query-based methods [20] pre-define learnable queries and using 3D points projecting to 2D features. In conclusion, there are three types of methods: *object-query-based*, *BEV-query-based*, and *BEV-depth-based methods*. The former two comprise query-based methods, and the latter two are BEV-based methods.

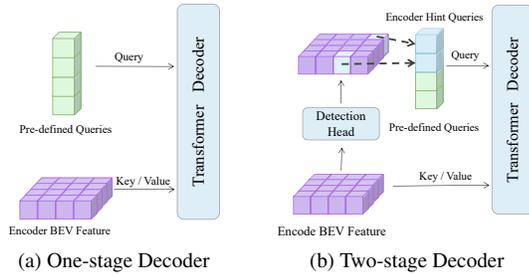


Figure S2: **One-stage and two-stage DETR style decoder.** (a) One-stage decoder uses the pre-defined queries to match the encoder output while (b) two-stage decoder use mixed queries while encoder hint queries are provided by the detection results after the encoder. These high-confidence locations are similar to proposals in two-stage CNN methods.

A.2. One-stage and Two-stage DETR Style Decoder

The original one-stage and two-stage DETR [3] style decoders are shown in Fig. S2. For the DETR [3] style decoder, the common one-stage approach [3] is to pre-define a certain number of queries as Fig. S2a. These queries correspond to the features of the encoder output. But it is difficult to find a direct match between the predefined queries and encoder output. Thus, inspired by the traditional CNN-based two-stage target detection, we can give rough results before giving final predictions. Two-stage DETR [3] style decoder [42] is such that, as shown in Fig. S2b, after encoder, we first use the detection head to make preliminary predictions on the encoder output and feed the detection high score results as queries of the decoder into the decoder.

B. Implementation Details

B.1. Glossary of Notation

All notations, including symbols and functions, are listed in Table S1 and in the order of appearance of the article. We explain all their definitions, and for symbols, we give their shapes while giving the inputs of all functions. The first part is the prerequisites to define the BEV queries. The two middle parts are the symbols that appear in Object Aligned Temporal Fusion. The second part is the notation used in the Ego-motion temporal fusion, and the third is for Object-motion fusion. The final part is the definition of Object Focused Multi-View Sampling and spatial exploitation.

B.2. Network Architecture

Backbone. OCBEV’s image backbone ResNet-101 [11] yields 3-level image feature maps of strides 16, 32, and 64. We replace the original convolution network with deformable convolutional networks v2 [7] for only the last two layers and the deformable groups equal to 1. The FPN [21]

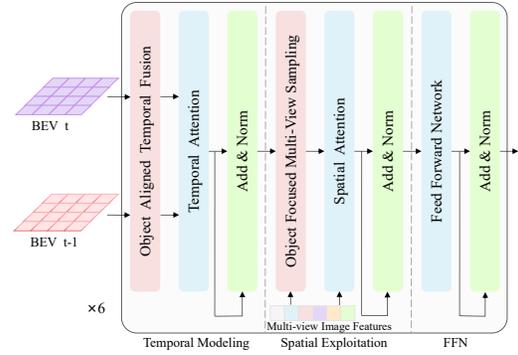


Figure S3: **The architecture of the encoder.** Red modules are added by ourselves. Blue modules are deformable attention, and green modules are *Add and Norm* modules. First and middle layers are temporal modeling and spatial exploitation, respectively. Final layers are Forward Network.

neck with input sizes 512, 1024, 2048 after the backbone produces 4-level features of stride 16, 32, 64, and 128.

Encoder layers. The architecture of the encoder is shown in Fig. S3. The dimension of BEV queries is 256. We use six layers of deformable attention [52] for the encoder to do temporal modeling and spatial exploitation. Note that we only consider the range with both length and width from $-51.2m$ to $51.2m$ and the height from $-5m$ to $3m$. As shown in Fig. S3, the encoder sequence is similar to BEVFormer [20]. For each sub-layer, the initial three sub-layers constitute temporal modeling. Before formal temporal attention, we add our Object Aligned Temporal Fusion to align and fuse the history information. The middle three sub-layers constitute spatial exploitation. We use Object Focused Multi-View Sampling to get object-dense 3D reference points and then feed these reference points to the spatial attention. For each height range, we sample 4 points along the z -axis evenly. All transformer layers have eight heads. A 2-layer Feed Forward Network with a 1024-dimension hidden layer is finally implemented after them.

Decoder layers. Decoder consists of 6 layers, similar to the encoder. The queries are 900 pre-defined queries, and keys and values are the encoder output BEV queries. So the attention for the decoder becomes regular 2D deformable attention [52]. Network parameters like eight heads and 0.1 dropout rate are identical to those in the encoder.

Heatmap branch. Before we feed the encoder output into the decoder, we first pass it through the heatmap module, which is similar to that in the FCOS3D [36]. The heatmap module is composed of simple convolutional layers, in order to downscale the 256 dimensions of the original BEV queries into a one-dimensional heatmap. As mentioned in the paper, centerness [49] value ranges from 0 to 1. During the training period, the centerness is 2D Gaussian distribu-

| Notation | Meaning | Shape |
|----------------------------------|---|----------------------------------|
| \mathbf{Q}_t | BEV queries at timestamp t . | $C_q \times H_q \times W_q$ |
| $\mathbf{Q}_{t'}$ | BEV queries at timestamp $t - 1$. | $C_q \times H_q \times W_q$ |
| $\bar{\mathbf{Q}}_{t'}$ | Flatten version of BEV queries at timestamp t . | $C_q \times (H_q W_q)$ |
| $\bar{\mathbf{Q}}_{t'}$ | Flatten version of BEV queries at timestamp $t - 1$. | $C_q \times (H_q W_q)$ |
| \mathcal{P} | Ego-motion overlapping part matching relationship between timestamp $t - 1$ and t . | N_p pairs |
| \mathbf{I} | Overlapping part matching indices at timestamp $t - 1$. | $1 \times N_p$ |
| \mathbf{J} | Overlapping part matching indices at timestamp t . | $1 \times N_p$ |
| $\mathcal{R}_{t' \rightarrow t}$ | Rotation matrix from timestamp $t - 1$ to t . | 3×3 |
| $\mathcal{R}_{t \rightarrow t'}$ | Rotation matrix from timestamp t to $t - 1$. | 3×3 |
| $\gamma_{t' \rightarrow t}$ | Yaw angle of rotation matrix from timestamp $t - 1$ to t . | \mathcal{R} |
| $\gamma_{t \rightarrow t'}$ | Yaw angle of rotation matrix from timestamp t to $t - 1$. | \mathcal{R} |
| $\mathcal{T}_{t' \rightarrow t}$ | Translation vector from timestamp $t - 1$ to t . | 3×1 |
| $\mathcal{T}_{t \rightarrow t'}$ | Translation vector from timestamp t to $t - 1$. | 3×1 |
| g_r | Function: Rotation operation. | - |
| g_t | Function: Translation operation. | - |
| f_a^{ego} | Function: Ego-motion temporal alignment and fusion. | - |
| \mathbf{Q}_a^{ego} | Aligned and fused BEV queries by Ego-motion temporal alignment and fusion. | $C_q \times H_q \times W_q$ |
| $\mathbf{p}_{t'}^m$ | xy coordinates of moving objects in timestamp $t - 1$ and $\mathbf{Q}_{t'}$ reference frame. | $N_m \times 2$ |
| $\mathbf{I}_{t'}^m$ | $\mathbf{Q}_{t'}$ indices of moving objects in timestamp $t - 1$ and $\mathbf{Q}_{t'}$ reference frame. | $1 \times N_m$ |
| $\mathbf{v}_{t'}^m$ | xy velocity of moving objects in timestamp $t - 1$. | $N_m \times 2$ |
| \mathbf{p}_t^m | xy coordinates of moving objects in timestamp t and $\mathbf{Q}_{t'}$ reference frame. | $N_m \times 2$ |
| \mathbf{I}_t^m | $\mathbf{Q}_{t'}$ indices of moving objects in timestamp t and $\mathbf{Q}_{t'}$ reference frame. | $1 \times N_m$ |
| \mathbf{J}_t^m | \mathbf{Q}_t indices of moving objects in timestamp t and \mathbf{Q}_t reference frame. | $1 \times N_m$ |
| f_a^{obj} | Function: Object-motion temporal alignment and fusion. | - |
| \mathbf{Q}_a^{obj} | Aligned and fused BEV queries by Object-motion temporal alignment and fusion. | $C_q \times H_q \times W_q$ |
| \mathbf{Q}_a | Aligned and fused BEV queries by Object Aligned Temporal Fusion. | $C_q \times H_q \times W_q$ |
| DAttn | Function: Deformable attention. | - |
| SpaA | Function: Spatial Attention between BEV and image features. | - |
| \mathbf{F}_t | Image features at timestamp t . | $C_f \times H_f \times W_f$ |
| \mathbf{Z}_h | Height range for the 3D reference points | $[\mathcal{R}_1, \mathcal{R}_2]$ |
| \mathbf{p}_{3d} | 3D reference points for Spatial Attention. | $N_h \times 3$ |
| g_p | Function: Projection 3D points to 2D image feature. | - |
| $\mathbf{Z}_{h,g}$ | Global height range for reference points | $[\mathcal{R}_1, \mathcal{R}_2]$ |
| $\mathbf{Z}_{h,l}$ | Local height range for reference points | $[\mathcal{R}_1, \mathcal{R}_2]$ |
| $\hat{\mathbf{Z}}_{h,l}$ | Adaptive local height range for reference points | $[\mathcal{R}_1, \mathcal{R}_2]$ |
| OFSpaA | Function: Spatial Attention by Object Focused Multi-View Sampling | - |

Table S1: Notation table: all symbols and functions. We show the shapes for symbols. Four parts are *BEV query definition*, *Ego-motion temporal fusion*, *Object-motion temporal fusion* and *Object Focused Multi-View Sampling*, respectively.

tion is written as:

$$c = e^{-\alpha((\Delta x)^2 + (\Delta y)^2)} \quad (8)$$

Here, c is the centerness [49] value. α adjusts the intensity attenuation from the center to the periphery, we set it to 2.5 for our experiments. Binary Cross Entropy (BCE) loss is used for its supervision, denoted as \mathcal{L}_c .

Head and loss. We adopt a DETR [52] style detection head after the decoder implemented by DETR3D [39], which consists of a classification head and a 3D bounding box

head. The classification head produces the logit number of each category. We use Focal loss [22] with $\gamma = 2$ and $\alpha = 0.25$ to supervise it, denoted as \mathcal{L}_{cls} . Besides, the 3D bounding box head predicts the boxes with the following coefficients: xyz coordinates of the center objects, the length, width, and height of the bounding boxes, the orientation of the boxes and the xy velocity of the objects. It adopts L1 loss \mathcal{L}_{bbox3d} for 3D bounding box regression. Here we have to notice these two losses have been implemented twice in the network. One of them is after the de-

| Hyperparameters | Value/Type |
|------------------------|------------------|
| batchsize | 1 |
| training epochs | 24 |
| optimizer | AdamW |
| learning rate | 2e-4 |
| weight decay | 0.01 |
| lr schedule | cosine annealing |
| warmup type | linear |
| warmup iters | 500 |
| warmup ratio | 1/3 |
| min lr ratio | 1e-3 |
| gradient clip | 35 |
| image size(main) | 1600 × 900 |
| image size(ablation) | 1280 × 720 |
| queue length(main) | 4 |
| queue length(ablation) | 3 |

Table S2: Training settings and hyperparameters.

coder as the final loss $\mathcal{L}_{\text{bbox3d}}$ which consists of the three losses in total:

$$\mathcal{L} = \lambda_c \mathcal{L}_c + \lambda_{cls} \mathcal{L}_{cls} + \lambda_{\text{bbox3d}} \mathcal{L}_{\text{bbox3d}} \quad (9)$$

Here we set $\lambda_{cls} = 2.0$ and $\lambda_{\text{bbox3d}} = 0.5$ following the DETR3D [39] and BEVFormer [20]. For the module we added, we make $\lambda_c = 1$ in practice. The other is that after the decoder, using a 3D Hungarian assigner to match the objects and decoder queries. To be more precise, the loss here should be called assigning cost.

B.3. Training Setting

We provide all the hyperparameters in our model in Table S2. Nearly all are the same as BEVFormer [20] and DETR3D [39]. To be specific, the first part is about some basic training parameters such as batchsize and epochs while the second part is about the optimizer. We use AdamW with an initial learning rate of 2e-4, and the schedule is Cosine Annealing Warm Restart. Also, we list some other parameters. The image size was set to 1600 × 900 for the main experimental results while 1280 × 720 for the ablation study. To get historical information, we take the previous four frames for the main experiments. For the ablation experiments, we take the preceding three frames. We illustrate how we use them for temporal modeling in Sec. B.4.

B.4. Dataset Pre-processing

How we use previous frames for temporal modeling. As mentioned earlier, we used the information from the previous frames to model the time. We model the temporal information slightly like the CLIP [31], which contains several sampling frames, and The last frame corresponds to the

current frame. So the first thing is to process the annotation information of the original nuScenes [1] so that the annotation of each frame contains the previous frames’ annotation as well. Similar to the BEVFormer [20], we use the ’Can bus’ information to record the ego-motion parameters in previous frames. The ’Can bus’ in the nuScenes [1] initially records the states of the ego-motion vehicle, including position, velocity, acceleration, steering, lights, battery, and so on. In temporal modeling, we only consider the position of the ego-motion vehicle. We set the position of the first frame of the sequence to zero, and the position of the subsequent frames to the distance moved relative to the first frame. This operation is offline and is a pre-processing of the overall nuScenes [1]. In the training dataloader, we save the various camera parameters of the historical frames and the state of the moving objects in them. Moreover, we align the history information between adjacent frames in the sequence. In the test mode, we test one scene like an extensive sequence and align them between adjacent frames.

Pre-processing and Transform. We have introduced the offline pre-processing of the entire nuScenes dataset [1] and the way to align and fuse historical frames. Now we introduce transform operations in the dataloader. In the training pipeline, Photo Metric Distortion was used for data augmentation and we also padded the multi-view images. Also we should note the image size is 1600 × 900 for the full version and 80% for the light version in some ablation studies.

C. Extended Experiment Results

In this section, we perform a series of additional experiments on top of the original ones. Here there are two main categories, and in Sec. C.1. we present a specific analysis of each category for the main experimental results. In Sec. C.2, we supplement the other ablation experiments.

C.1. Analysis of Categories in Main Results

The most complete experimental configuration was used for the main experimental results for comparison with other methods. Here we listed the results of each category. There are 10 evaluable categories in total as shown in Table S3. For each one, we show all its six metrics, including precision, translation, scale, orientation, velocity and attribute. For each small pane, there are three pieces of data. The data to the left of the slash is the BEVFormer [20] result, the data to the right of the slash is our OCBEV result, and the data in the right column is the result obtained by subtracting the BEVFormer [20] result from the OCBEV result. The green bar means that OCBEV is better than BEVFormer [20], and the gray bar means ours is worse than BEVFormer [20].

In terms of precision, BEVFormer [20] and OCBEV are almost the same, OCBEV performs poorly on large objects such as truck, bus and trailer and performs well on normal moving objects like car and motorcycle. Due to OCBEV is

| Object Class | AP(Precision) ↑ | ATE(Translation) ↓ | ASE(Scale) ↓ | AOE(Orientation) ↓ | AVE(Velocity) ↓ | AAE(Attribute) ↓ |
|----------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| car | 0.617 / 0.621 (+0.004) | 0.463 / 0.446 (-0.017) | 0.152 / 0.149 (-0.003) | 0.067 / 0.057 (-0.010) | 0.324 / 0.295 (-0.029) | 0.196 / 0.194 (-0.029) |
| truck | 0.370 / 0.345 (-0.025) | 0.726 / 0.697 (-0.029) | 0.212 / 0.204 (-0.008) | 0.094 / 0.079 (-0.015) | 0.341 / 0.327 (-0.014) | 0.194 / 0.195 (+0.001) |
| bus | 0.445 / 0.427 (-0.018) | 0.747 / 0.729 (-0.018) | 0.212 / 0.202 (-0.010) | 0.097 / 0.107 (+0.010) | 0.864 / 0.636 (-0.228) | 0.273 / 0.296 (+0.023) |
| trailer | 0.171 / 0.168 (-0.003) | 0.971 / 0.929 (-0.042) | 0.259 / 0.251 (-0.008) | 0.426 / 0.485 (+0.059) | 0.338 / 0.289 (-0.049) | 0.078 / 0.065 (-0.013) |
| construction vehicle | 0.129 / 0.115 (-0.019) | 0.999 / 0.952 (-0.047) | 0.451 / 0.455 (+0.004) | 1.157 / 0.954 (-0.203) | 0.131 / 0.133 (+0.002) | 0.373 / 0.340 (-0.034) |
| pedestrian | 0.494 / 0.502 (-0.008) | 0.642 / 0.624 (-0.018) | 0.294 / 0.296 (+0.002) | 0.432 / 0.410 (-0.022) | 0.361 / 0.345 (-0.016) | 0.174 / 0.157 (-0.017) |
| motorcycle | 0.433 / 0.434 (+0.001) | 0.643 / 0.589 (-0.054) | 0.257 / 0.261 (+0.004) | 0.441 / 0.342 (-0.099) | 0.441 / 0.466 (+0.025) | 0.266 / 0.232 (-0.034) |
| bicycle | 0.399 / 0.396 (-0.003) | 0.555 / 0.541 (-0.014) | 0.272 / 0.281 (+0.009) | 0.485 / 0.516 (+0.031) | 0.253 / 0.246 (-0.007) | 0.025 / 0.013 (-0.012) |
| traffic cone | 0.584 / 0.615 (+0.031) | 0.452 / 0.372 (-0.080) | 0.331 / 0.332 (+0.001) | - / - | - / - | - / - |
| barrier | 0.524 / 0.551 (+0.027) | 0.531 / 0.414 (-0.117) | 0.292 / 0.299 (+0.007) | 0.139 / 0.101 (-0.038) | - / - | - / - |

Table S3: **Categories details in the Main Result.** Here we show the main results comparison of BEVFormer [20] and our OCBEV. The data on the left of the slash is the result of BEVFormer [20], and the right is that of OCBEV. Data on the far right is the difference. Green means that OCBEV is better than BEVFormer [20], and gray means OCBEV is worse.

| Backbone | Loss Weight | NDS ↑ | mAP ↑ | Backbone | LR | NDS ↑ | mAP ↑ | Backbone | BEV Size | NDS ↑ | mAP ↑ |
|------------|-------------|-------|-------|------------|--------------------|-------|-------|------------|-----------|-------|-------|
| ResNet-101 | 2 : 0.4 | 0.528 | 0.413 | ResNet-101 | 1×10^{-4} | 0.509 | 0.316 | ResNet-101 | 200 × 200 | 0.531 | 0.414 |
| ResNet-101 | 2 : 0.5 | 0.532 | 0.417 | ResNet-101 | 2×10^{-4} | 0.532 | 0.417 | ResNet-101 | 300 × 300 | 0.532 | 0.417 |
| ResNet-101 | 2 : 0.6 | 0.531 | 0.414 | ResNet-101 | 4×10^{-4} | 0.487 | 0.375 | ResNet-101 | 350 × 350 | 0.529 | 0.415 |
| ResNet-101 | 2 : 0.7 | 0.530 | 0.413 | ResNet-101 | 6×10^{-4} | 0.493 | 0.356 | ResNet-101 | 400 × 400 | 0.528 | 0.415 |

(a) **Loss Weight.** The loss weight is weights of classification and 3D bounding boxes and we set the centerness loss as 1 in practice. Its fluctuation range becomes 0.2% in NDS and 0.4% in mAP which means that the effect of the loss weight on the main experimental results is not significant.

(b) **Learning Rate.** We change the learning rate for the whole network. All learning rates are in the 10^{-4} level. Its fluctuation range becomes 4.5% in NDS and 10.1% in mAP which means that the effect of the learning rate on the main experimental results is extremely significant.

(c) **BEV Size.** BEV size below 200 has been explored in BEVFormer so We set the BEV size larger than 200 respectively. Its fluctuation range becomes 0.4% in NDS and 0.3% in mAP which means that Our model is very robust to the parameter BEV size and its effect is very slight.

Table S4: **Ablation study.** We conducted additional ablation experiments on a full version OCBEV model trained with 12 epochs and the configuration of the model is the same as that in the main result experiments. Our results indicate that the model is robust to **loss weight** and **BEV size**, and **learning rate** have a significant impact on the performance of the model.

Object-Centric, OCBEV is particularly prominent in metrics about the state of motion of objects like translation, orientation and velocity, which means our Object-centric module model the motion characteristics of objects better. Also OCBEV is also better than BEVFormer [20] in terms of the attribute metric, but slightly worse in terms of scale. From the perspective of each category, OCBEV outperforms BEVFormer [20] in all aspects for the car category and has better performance in other kinds of objects as well.

C.2. Additional Ablation Studies

A series of additional experiments were carried out for the main results to exclude the interference of other factors. The setting of ablation experiments here differs from that of the main paper as the model configurations here are the same as those in the main results *i.e.* 6 layers of encoder and decoder. We did not disable the temporal attention as well.

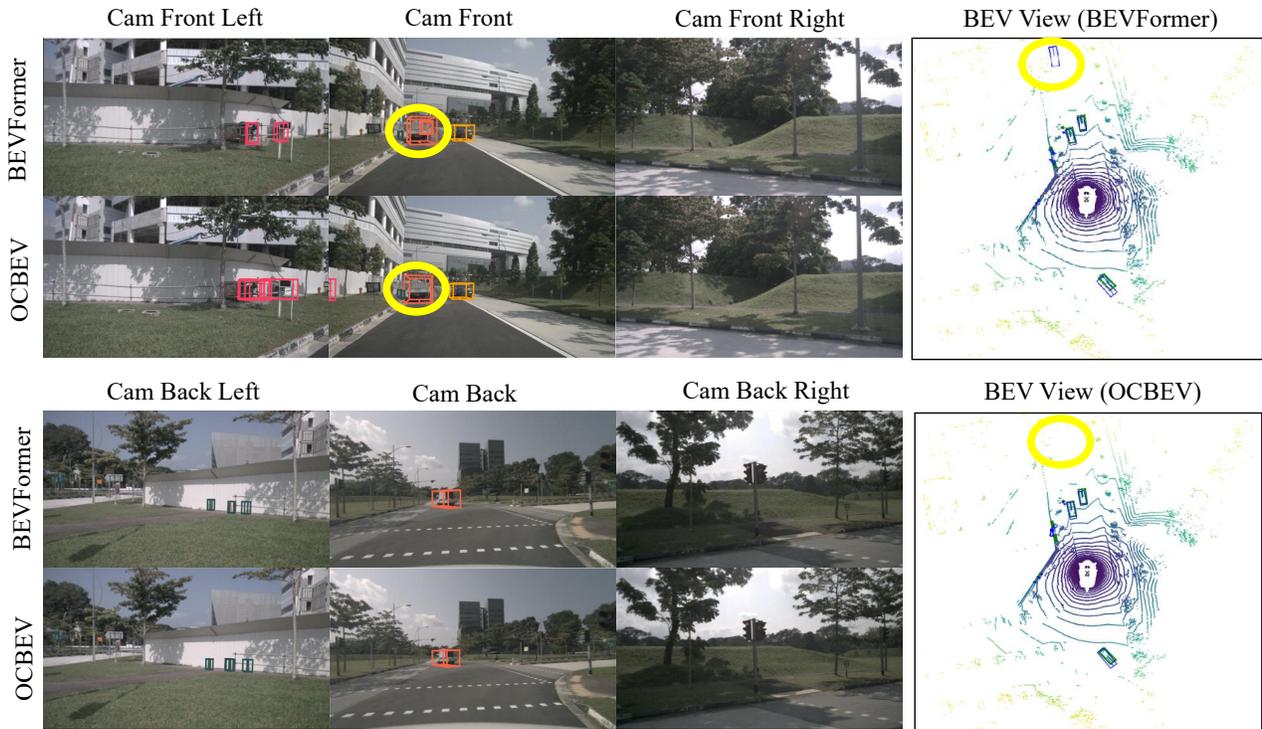
Loss weight. We do additional ablation experiments on the loss weight, learning rate and BEV size to exclude these factors. About the loss weights, the loss we added \mathcal{L}_c , and we set its weight $\lambda_c = 1$ as practice. But for the loss weights of classification and bounding boxes, we follow the DETR3D [39] and BEVFormer [20], and do ablation experiments on the loss weight. Table S4a shows that the effect

of loss weight on the experimental results is not significant. **Learning rate.** We also follow the DETR3D [39] and BEVFormer [20] and set the learning rate to 2×10^{-4} . Here we explore the effect of the learning rate on the experiments as well. Table S4b shows that the learning rate has a significant effect on the results. The fluctuation range reaches 4.5% in NDS (48.7% and 53.2%) and 10.1% in the mAP metric.

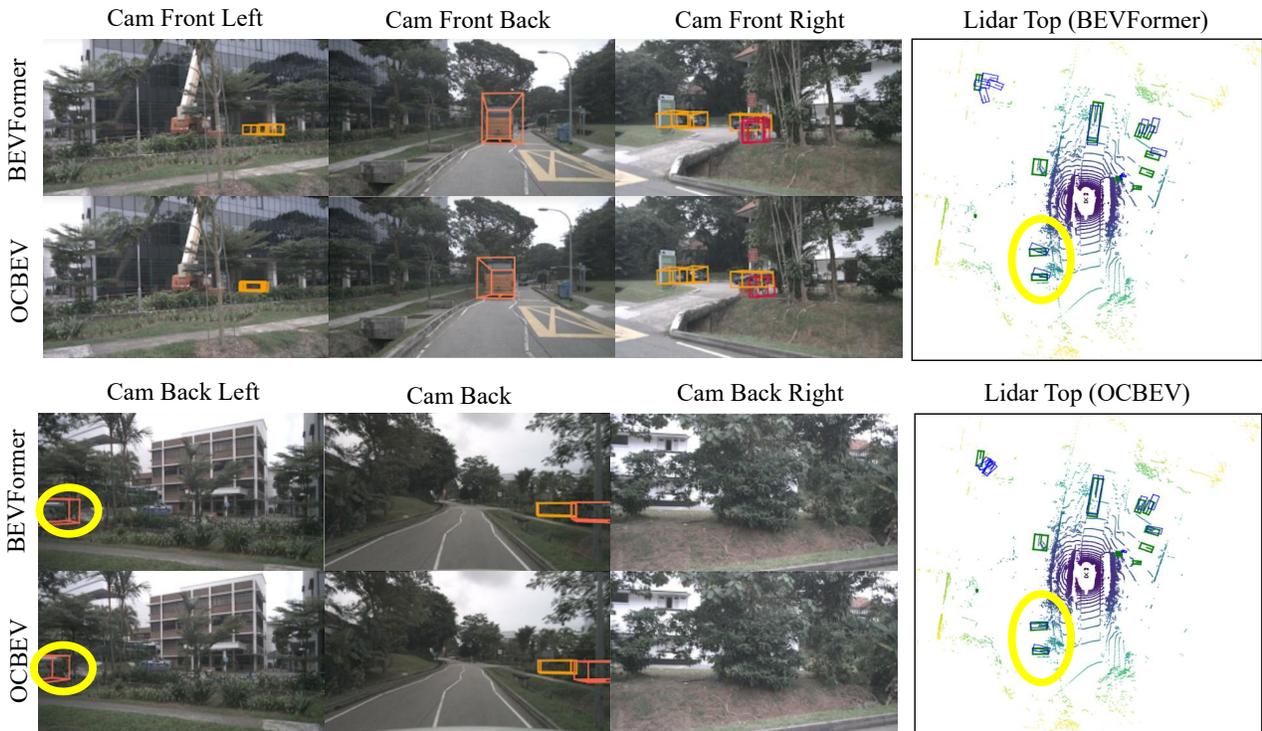
BEV size. As mentioned in the main paper. We did not follow the BEV size in DETR3D [39] and BEVFormer [20] which is 200×200 and we set the batchsize to 300×300 . We also performed additional ablation experiments for BEV size. BEV size below 200 has been explored in BEVFormer [20] so We set it larger than 200. Table S4c shows that BEV size is an irrelevant variable.

D. Visualization

We demonstrate the visualization and qualitative results of our 3D object detection in Fig. S1. We give the qualitative results of perspective views and the BEV view. We compare our results with BEVFormer. It is shown that our OCBEV predicts more accurate bounding boxes than that of BEVFormer. OCBEV can detect hard cases in the distance or with occlusion which prove OCBEV is more advanced.

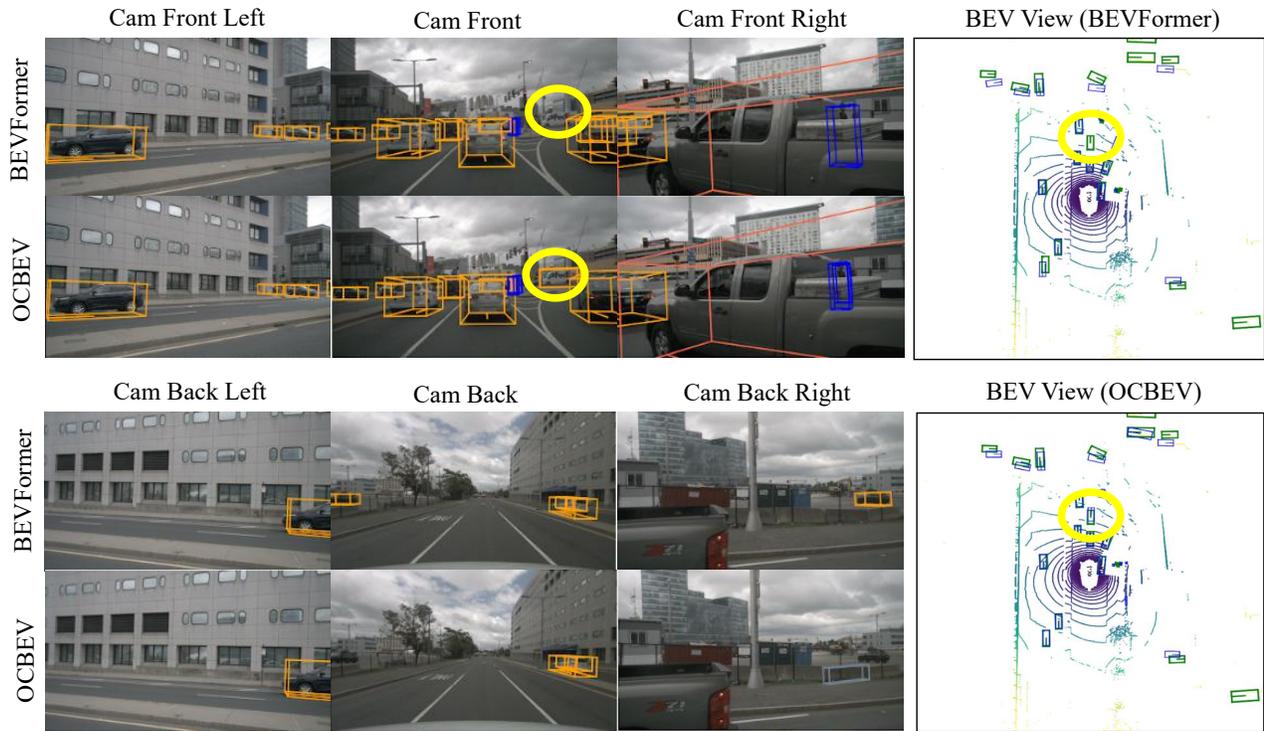


(a) OVBEV can reduce false positive cases. BEVFormer predicts a car in front view but actually there is no car there at all.

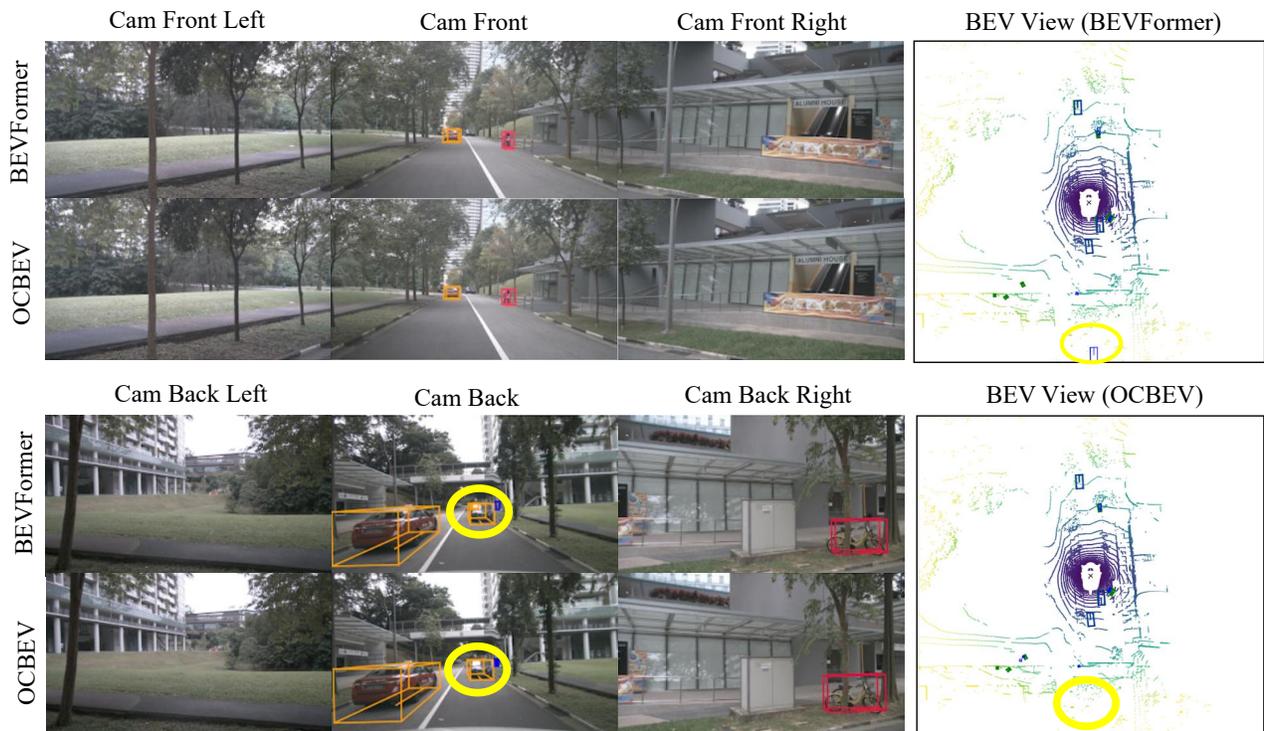


(b) OCBEV can give more precise locations and poses . BEVFormer predicts locations and poses of cars in back left view poorly than OCBEV.

Figure S4: **Comparison of BEVFormer and our OCBEV on the nuScenes val set.** In BEV view, the green boxes are the ground truth and blue boxes are predictions. And OCBEV can reduce false positive cases and give more accuracy boxes



(a) OCBEV detect the Fedex car in the front view while BEVFormer does not detect that successfully.



(b) For dense vehicle sceneslike car park, OCBEV can give better performance than BEVFormer.

Figure S5: **Comparison of BEVFormer and our OCBEV on the nuScenes val set.** In BEV view, the green boxes are the ground truth and blue boxes are predictions. OCBEV have higher precision and can handle dense vehicle scenes.

References

- [1] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, 2020. 2, 5, 6, 12
- [2] Yigit Baran Can, Alexander Liniger, Danda Pani Paudel, and Luc Van Gool. Structured bird’s-eye-view traffic scene understanding from onboard images. In *ICCV*, 2021. 2
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 1, 2, 3, 6, 10
- [4] Shaoyu Chen, Xinggang Wang, Tianheng Cheng, Qian Zhang, Chang Huang, and Wenyu Liu. Polar parametrization for vision-based surround-view 3d detection. *arXiv:2206.10965*, 2022. 2
- [5] Xiaowei Chi, Jiaming Liu, Ming Lu, Rongyu Zhang, Zhaoqing Wang, Yandong Guo, and Shanghang Zhang. Bev-san: Accurate bev 3d object detection via slice attention networks. *arXiv:2212.01231*, 2022. 3, 5
- [6] Xiaomeng Chu, Jiajun Deng, Yuan Zhao, Jianmin Ji, Yu Zhang, Houqiang Li, and Yanyong Zhang. Oa-bev: Bringing object awareness to bird’s-eye-view representation for multi-camera 3d object detection. *arXiv:2301.05711*, 2023. 3
- [7] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, 2017. 10
- [8] Ross Girshick. Fast r-cnn. In *ICCV*, 2015. 3
- [9] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 3
- [10] Adam W Harley, Zhaoyuan Fang, Jie Li, Rares Ambrus, and Katerina Fragkiadaki. Simple-bev: What really matters for multi-sensor bev perception? In *ICRA*, 2023. 2
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6, 10
- [12] Junjie Huang and Guan Huang. Bevdet4d: Exploit temporal cues in multi-camera 3d object detection. *arXiv:2203.17054*, 2022. 3, 6
- [13] Junjie Huang, Guan Huang, Zheng Zhu, and Dalong Du. Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. *arXiv:2112.11790*, 2021. 1, 2, 3, 7, 9
- [14] Peixiang Huang, Li Liu, Renrui Zhang, Song Zhang, Xinli Xu, Baichao Wang, and Guoyi Liu. Tig-bev: Multi-view bev 3d object detection via target inner-geometry learning. *arXiv:2212.13979*, 2022. 2
- [15] Yanqin Jiang, Li Zhang, Zhenwei Miao, Xiatian Zhu, Jin Gao, Weiming Hu, and Yu-Gang Jiang. Polarformer: Multi-camera 3d object detection with polar transformers. In *AAAI*, 2023. 2
- [16] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, 2019. 1
- [17] Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M Ni, and Lei Zhang. Dn-detr: Accelerate detr training by introducing query denoising. In *CVPR*, 2022. 3
- [18] Yinhao Li, Han Bao, Zheng Ge, Jinrong Yang, Jianjian Sun, and Zeming Li. Bevstereo: Enhancing depth estimation in multi-view 3d object detection with dynamic temporal stereo. *arXiv:2209.10248*, 2022. 2, 3
- [19] Yinhao Li, Zheng Ge, Guanyi Yu, Jinrong Yang, Zengran Wang, Yukang Shi, Jianjian Sun, and Zeming Li. Bevdepth: Acquisition of reliable depth for multi-view 3d object detection. *arXiv:2206.10092*, 2022. 1, 2, 3
- [20] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Qiao Yu, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *ECCV*, 2022. 1, 2, 3, 5, 6, 7, 9, 10, 12, 13
- [21] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 3, 6, 10
- [22] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 11
- [23] Xuewu Lin, Tianwei Lin, Zixiang Pei, Lichao Huang, and Zhizhong Su. Sparse4d: Multi-view 3d object detection with sparse spatial-temporal fusion. *arXiv:2211.10581*, 2022. 1, 2
- [24] Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr: Position embedding transformation for multi-view 3d object detection. In *ECCV*, 2022. 3
- [25] Yingfei Liu, Junjie Yan, Fan Jia, Shuailin Li, Qi Gao, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr2: A unified framework for 3d perception from multi-camera images. *arXiv:2206.01256*, 2022. 3, 6, 7

- [26] Zhipeng Luo, Changqing Zhou, Gongjie Zhang, and Shijian Lu. Detr4d: Direct multi-view 3d object detection with sparse attention. *arXiv:2212.07849*, 2022. [1](#), [3](#), [6](#), [7](#)
- [27] Sambit Mohapatra, Senthil Yogamani, Heinrich Gotzig, Stefan Milz, and Patrick Mader. Bevdetnet: bird’s eye view lidar point cloud based real-time 3d object detection for autonomous driving. In *ITSC*, 2021. [2](#)
- [28] Bowen Pan, Jiankai Sun, Ho Yin Tiga Leung, Alex Andonian, and Bolei Zhou. Cross-view semantic segmentation for sensing surroundings. *RAL*, 2020. [2](#)
- [29] Jinhyung Park, Chenfeng Xu, Shijia Yang, Kurt Keutzer, Kris Kitani, Masayoshi Tomizuka, and Wei Zhan. Time will tell: New outlooks and a baseline for temporal multi-view 3d object detection. In *ICLR*, 2023. [3](#)
- [30] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *ECCV*, 2020. [2](#), [3](#)
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. [12](#)
- [32] Thomas Roddick, Alex Kendall, and Roberto Cipolla. Orthographic feature transform for monocular 3d object detection. In *BMCV*, 2019. [2](#)
- [33] Andrea Simonelli, Samuel Rota Buló, Lorenzo Porzi, Manuel López-Antequera, and Peter Kotschieder. Disentangling monocular 3d object detection. In *ICCV*, 2019. [1](#)
- [34] Sourabh Vora, Alex H Lang, Bassam Helou, and Oscar Beijbom. Pointpainting: Sequential fusion for 3d object detection. In *CVPR*, 2020. [1](#)
- [35] Tai Wang, Qing Lian, Chenming Zhu, Xinge Zhu, and Wenwei Zhang. Mv-fcos3d++: Multi-view camera-only 4d object detection with pretrained monocular backbones. *arXiv:2207.12716*, 2022. [2](#), [3](#)
- [36] Tai Wang, Xinge Zhu, Jiangmiao Pang, and Dahua Lin. Fcos3d: Fully convolutional one-stage monocular 3d object detection. In *ICCV*, 2021. [1](#), [6](#), [7](#), [10](#)
- [37] Tai Wang, Xinge Zhu, Jiangmiao Pang, and Dahua Lin. Probabilistic and geometric depth: Detecting objects in perspective. In *CoRL*, 2022. [1](#), [7](#)
- [38] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *CVPR*, 2019. [2](#)
- [39] Yue Wang, Vitor Campagnolo Guizilini, Tianyuan Zhang, Yilun Wang, Hang Zhao, and Justin Solomon. Detr3d: 3d object detection from multi-view images via 3d-to-2d queries. In *CoRL*, 2022. [1](#), [2](#), [3](#), [6](#), [7](#), [9](#), [11](#), [12](#), [13](#)
- [40] Xinshuo Weng and Kris Kitani. Monocular 3d object detection with pseudo-lidar point cloud. In *ICCV Workshop*, 2019. [2](#)
- [41] Junjie Yan, Yingfei Liu, Jianjian Sun, Fan Jia, Shuailin Li, Tiancai Wang, and Xiangyu Zhang. Cross modal transformer via coordinates encoding for 3d object detection. *arXiv:2301.01283*, 2023. [2](#)
- [42] Chenyu Yang, Yuntao Chen, Hao Tian, Chenxin Tao, Xizhou Zhu, Zhaoxiang Zhang, Gao Huang, Hongyang Li, Yu Qiao, Lewei Lu, et al. Bevformer v2: Adapting modern image backbones to bird’s-eye-view recognition via perspective supervision. In *CVPR*, 2023. [2](#), [3](#), [10](#)
- [43] Zhuyu Yao, Jiangbo Ai, Boxun Li, and Chi Zhang. Efficient detr: improving end-to-end object detector with dense prior. *arXiv:2104.01318*, 2021. [3](#)
- [44] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-based 3d object detection and tracking. In *CVPR*, 2021. [3](#)
- [45] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In *CVPR*, 2018. [7](#)
- [46] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M. Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. In *ICLR*, 2023. [3](#)
- [47] Yunpeng Zhang, Zheng Zhu, Wenzhao Zheng, Junjie Huang, Guan Huang, Jie Zhou, and Jiwen Lu. Reverse: Unified perception and prediction in birds-eye-view for vision-centric autonomous driving. *arXiv:2205.09743*, 2022. [2](#)
- [48] Brady Zhou and Philipp Krähenbühl. Cross-view transformers for real-time map-view semantic segmentation. In *CVPR*, 2022. [2](#)
- [49] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. In *CVPR*, 2019. [1](#), [6](#), [7](#), [10](#), [11](#)
- [50] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *CVPR*, 2018. [1](#)
- [51] Benjin Zhu, Zhengkai Jiang, Xiangxin Zhou, Zeming Li, and Gang Yu. Class-balanced grouping and sampling for point cloud 3d object detection. *arXiv:1908.09492*, 2019. [7](#)
- [52] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *ICLR*, 2021. [3](#), [5](#), [6](#), [10](#), [11](#)