# Focal-PETR: Embracing Foreground for Efficient Multi-Camera 3D Object Detection

Shihao Wang*
wangshihao@bit.edu.cn

Xiaohui Jiang*
xhjiang@bit.edu.cn

Ying Li[†]
ying.li@bit.edu.cn

Beijing Institute of Technology

## Abstract

*The dominant multi-camera 3D detection paradigm is based on explicit 3D feature construction, which requires complicated indexing of local image-view features via 3D-to-2D projection. Other methods implicitly introduce geometric positional encoding and perform global attention (e.g., PETR) to build the relationship between image tokens and 3D objects. The 3D-to-2D perspective inconsistency and global attention lead to a weak correlation between foreground tokens and queries, resulting in slow convergence. We propose **Focal-PETR** with instance-guided supervision and spatial alignment module to adaptively focus object queries on discriminative foreground regions. **Focal-PETR** additionally introduces a down-sampling strategy to reduce the consumption of global attention. Due to the highly parallelized implementation and down-sampling strategy, our model, without depth supervision, achieves leading performance on the large-scale nuScenes benchmark and a superior speed of **30 FPS** on a single RTX3090 GPU. Extensive experiments show that our method outperforms PETR while consuming **3x** fewer training hours. The code will be made publicly available.*

## 1. Introduction

Camera-based 3D object detection, compared with LiDAR-based counterparts has attracted immense attention in recent years due to its low cost for deployment and dense semantic information [24, 25]. Multiple cameras with overlapping regions usually need to be setup for camera-based panoramic perception. To efficiently fuse information from overlapping regions, a unified representation for arbitrary camera rigs [10, 13, 16, 43, 46] has been widely concerned.

To explore a unified 3D representation, the dominant multi-camera 3D detection models are based on explicit 3D feature construction, which mainly relies on precise prediction of depth distribution or perspective projection
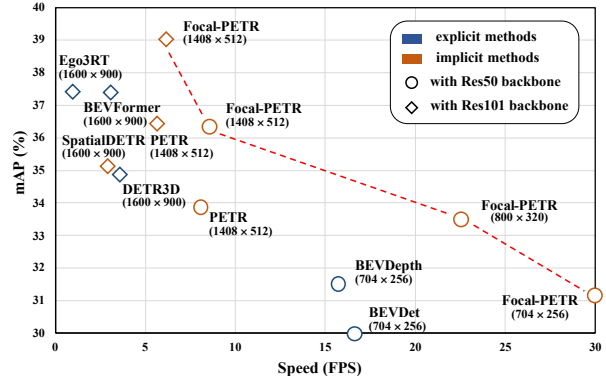
---

*equal contribution
[†]Corresponding author



Figure 1. Speed-accuracy trade-off of different models with a single time stamp input on nuScenes val set. The inference speed is calculated on a single RTX3090 GPU. Due to our focal sampling and high parallelism, we have achieved superior performance. The corresponding input resolutions are in brackets. More detailed analysis can be found in Table 1.

[15, 16, 28, 43]. These explicit models have achieved leading performance on benchmarks, while they are unfriendly for parallel computing on GPU devices due to complicated feature indexing operations. Another paradigm has been proposed by simply converting 2D image features into 3D position awareness through implicit position embedding. Such paradigm has achieved comparable performance with advantages of global modeling capabilities and parallel computing performance, as illustrated in Figure 1. While conciseness and effectiveness, the long training schedule and huge memory consumption limit its large-scale application in autonomous driving compared with the explicit paradigm. We suggest that the long training schedule is largely on account of the misalignment between the training objectives and features gathered by object queries. The explicit methods directly align features through auxiliary depth supervision or accurate 3D-to-2D projection. However, this mechanism in implicit methods has not been well explored.

We explore PETR [20, 21], which performs best among existing implicit methods, and suggest two limitations of
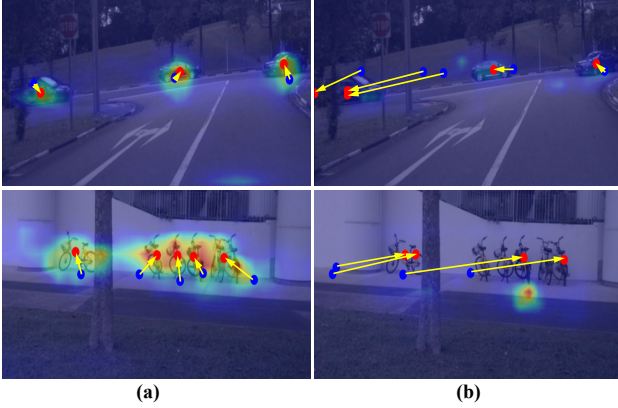
1

Figure 2. Comparison of attention weight maps for (a) Focal-PETR and (b) PETR with 24 training epochs without CBGS [48]. The reference points corresponding to different layers are projected onto the image plane. The blue and red dots represents initialized anchor points and final predictions respectively. It can be seen that the attention map of our Focal-PETR focuses on the foreground objects.

this paradigm: (i) semantic ambiguity and (ii) spatial misalignment. The reason for (i) is the similar content embedding of image tokens. As shown in Figure 2, the object queries give a nearly uniform attention weights on non-local regions and their initialized reference points are far from the predictions. We conclude that the lack of discrimination in foreground features makes queries difficult to focus on extremities. After subsequent attention layers, the obtained object queries for positive sample matching are not closely correlated with relevant foreground features during the training phase. The reason for (ii) is that the geometric cues only play the role of positional embedding and are not considered as a part of content collected by object queries (see Figure 4). Thus, the search mode of object queries is semantic-biased in the cross-attention modules, which means that the calculation of query-to-feature similarity is spatial insensitive.

We propose Focal-PETR, a semantic-aggregated and spatial-aligned framework to strengthen the discrimination and spatial sensitivity of foreground tokens. Specifically, three auxiliary tasks, namely class-aware, IoU-aware and centroid-aware strategies, are adopted to explore salient tokens in an instance-guided manner. Based on them we verify that the view transformer process can be considerably achieved by condensing the foreground tokens into a smaller set. This design greatly reduces the computing consumption of global cross attention. Additionally, a spatial alignment module is designed to embed the spatial information into image tokens, which enables the decoder to effectively search the content.

Figure 1 shows that the proposed Focal-PETR preserves and enhances the advantages of implicit 3D detec-

tion paradigm, which can be highly parallel implemented. On nuScenes [1] val set, the lightweight model with Res50 backbone achieves superior accuracy and speed trade-off (i.e., 31.1% mAP and 30.0 FPS on a single RTX3090). Extensive experiments and analysis on Section 4 justify the effectiveness and efficiency of the method.

To summarize, the contributions of this paper are:

- We first identify the semantic ambiguity and spatial misalignment in existing implicit 3D detection paradigms, causing suboptimal discriminant features extraction and slow convergence speed.

- We propose Focal-PETR by introducing focal sampling and spatial alignment modules. Our method mitigates the aforementioned problems and efficiently focuses on the foreground tokens. We also analyze the computational consumption and memory footprint to further validate the efficiency of our proposed method.

- Experiments on the large-scale nuScenes benchmark show superior efficiency and state-of-the-art performance (46.5% mAP and 51.6% NDS) of the proposed Focal-PETR with a single time stamp input.

## 2. Related Work

### 2.1. Multi-camera 3D Object Detection

3D object detection based on unified representation has attracted increasing attention for multi-sensors fusion. Previous works such as Pseudo-LiDAR [42], OFT [31], LSS [28], and BEVDet [10] explicitly predict the depth distribution to lift 2D features to 3D space. DETR3D [43] first projects the predefined queries onto images and then adopts attention mechanism to model the relation with the multiview features respectively. BEVFomrer [16] further extends this idea by using dense queries and temporal fusion. These methods explicitly index the local image feature from 2D perspective-view to 3D space, facilitating the alignment between training targets and image features. Other works [6, 20, 21, 46] model the view transformation by implicitly encoding geometric information, building the interaction between 3D queries and image tokens. Due to the powerful modeling capabilities of Transformer [37], models based on implicit positional encoding can extract global information in a parallel way but suffer from slow convergence and huge memory complexity. We analyze the mechanism of implicit paradigm and suggest that the weak ability to locate foreground information impairs the representation of object queries, resulting in the misalignment of training targets and semantic content. We perform auxiliary tasks to adaptively focus the attention on salient regions.
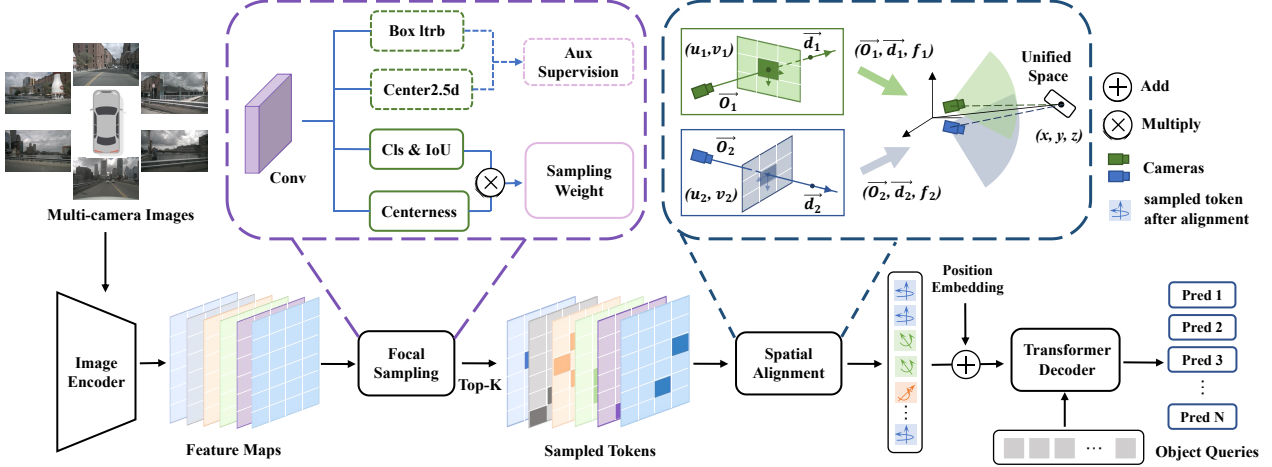
Figure 3. Overview of Focal-PETR. The multi-camera images are first fed into an image encoder to learn a high-level 2D representation. Then a focal sampling module is applied to distinguish discriminative foreground tokens. The auxiliary supervision which is turned off at inference is shown in dashed lines. The sampled tokens are further fed into the spatial alignment module to convert 2D features into a unified 3D space, according to parametric conical frustum. Finally, Transformer decoder is employed to generate 3D bounding boxes. $\vec{o}$ and $\vec{d}$ represent the optical center and direction of a specific pixel ray. $f$ is the camera intrinsic, additionally employed to describe frustum cone viewed by each ray.

## 2.2. Fast Converge of Query-based detectors

Query-based detectors have been widely studied due to their high performance and simplicity, while the slow training convergence limits the large-scale deployment [12]. Some works [4, 7, 12, 19, 26, 34, 44, 49] have tried to solve this problem. Several of them attempt to improve the network structure by taking local attention operation, such as Deformerable DETR [49], Dynamic DETR [4], AdaMixer [7]. Besides, initializing queries with meaningful information has been investigated. Anchor DETR [44] and DAB-DETR [19] interpret query as 2D reference points or 4D anchor boxes. Conditional DETR [26] combines the content and position information together to focus every query on a specific spatial space. DN-DETR [12] considers the influence of instable bipartite graph matching and introduces query denoising to mitigate slow convergence. The training objectives and features of 2D query-based detectors are all located in the same perspective view. While the 3D-to-2D query-based detectors are more difficult to converge due to the weak correspondence between object queries and the gathered features. We propose a spatial alignment module to enhance the spatial sensitivity of object queries on the basis of semantic correspondence.

## 2.3. Ranking and Sampling in Object Detection

Convolution based dense detectors [18, 29, 36] with one-to-many label assignment criteria have led the mainstream of object detection. Due to the non-uniformity and sparsity of foreground information in visual signal, duplicate predictions of objects are inevitable. Generally, additional NMS [29] and bounding box quality ranking strategies [14, 45],

such as classification, centerness [36] and IoU scores, are introduced to eliminate the redundant predictions. DETR [2] outperforms competitive Faster R-CNN [29] baseline by relying on the Transformer [37] architecture and one-to-one assignment strategy [2]. However, the huge computation and memory consumption of attention in Transformer limit its further development. Recent works conduct sparse sampling for attention mechanism [32, 40, 49]. Deformable DETR [49] performs learnable adjacent sampling instead of the entire image features. PnP-DETR [40] and Sparse DETR [32] sample a fraction of salient tokens in an unsupervised way for the subsequent encoder and decoder. In this paper, we propose focal sampling strategy based on the instance-guided supervision. Simply sampling foreground features via detection quality ranking achieves competitive performance and verify the training consistency between 2D and 3D object detection tasks.

## 3. Method

Existing 3D detection models based on implicit positional encoding [6, 20, 21] directly utilize pixel-level features extracted by an image encoder as the minimum unit to process view transformation. We suggest that the direct utilization of pixel representation makes object queries hard to focus on foreground features. Thus, we attempt to interpret the features as discriminative instances. In the following, we first revisit the detection pipeline proposed in PETR [20], which is high parallelism and suitable for fusing heterogeneous features between instances. Then we elaborate on our instance-guided down-sampling strategy and spatial alignment module.

## 3.1. Preliminaries: PETR

PETR is built upon the Transformer decoder architecture [37]. Its core components include image encoder, positional encoder and detection head. Combining sparse BEV queries with cross-attention mechanism, PETR refines the detection prediction using implicit features enhanced by 3D positional embedding in a cascade manner. A brief review of the PETR pipeline is made as follows.

**3D Positional Encoder.** Given an image $I_j \in \mathbb{R}^{3 \times H \times W}$ captured from one of N ($j \in \{1, 2, ..., N\}$) surround cameras with its corresponding intrinsic matrix $K_j \in \mathbb{R}^{3 \times 3}$ and extrinsic matrix $P_j \in \mathbb{R}^{4 \times 4}$, a tracing ray corresponding to each pixel center in a unified coordinate system $r_j^{u,v}(t) \epsilon \mathbb{R}^3 = \{o_j + t d_j^{u,v} | t \epsilon \mathbb{R}, o \epsilon \mathbb{R}^3, d \epsilon \mathbb{R}^3\}$ can be derived. Each pixel position $(u, v)$ on the specific camera emits a unique ray along the direction $d_j^{u,v}$, which passes through the optical center $o_j$ of the corresponding camera. Based on the aforementioned ray equation, the linear-increasing discretization (LID) [35] is adopted to approximately sample the rays at different distances $t_i$, where $i$ is the depth bin index in LID. Then the coordinates $C_j^{u,v}$ of sampled points are normalized and fed into a 2-layer multi-layer perception (MLP) $\varphi$. The resulting positional embedding is noted as $E_j^{u,v}$. This process can be abstracted as:

$$C_j^{u,v} = Concat[r_j^{u,v}(t_1), \ r_j^{u,v}(t_2), \ldots, r_j^{u,v}(t_i)] \quad (1)$$

$$E_j^{u,v} = \varphi(Norm(C_j^{u,v})) \quad (2)$$

**Detection Head.** The detection head is composed of conventional Transformer decoder layers [37], which interacts queries $q_L$ of layer $L$ with 3D positional enhanced features ($k$) from image encoder. To enable the features $F_j^{u,v}$ aggregated by queries be stably supervised, the queries are defined as a group of learnable spatial anchor points. The interaction mentioned above can be expressed as:

$$[k, v] = [F_j^{u,v} + E_j^{u,v}, \ F_j^{u,v}] \quad (3)$$

$$q_L = Softmax(q_{L-1} \cdot k^T) \cdot v \quad (4)$$

For simplicity, we ignore the scaling factor. It should be noticed that, the 3D implicit features are only used as keys ($k$) to interact with those queries. It is to say that the values ($v$) used for weighted sum are geometry-ignored. Finally, queries after multi-layer refinement are sent to finish classification and regression tasks respectively. For more details, please refer to the original paper of PETR [20].

## 3.2. Focal Sampling Strategy

The purpose of our focal sampling is to distinct features in instance-level and distill the representative tokens while ensuring high recall. Specifically, we first adopt an
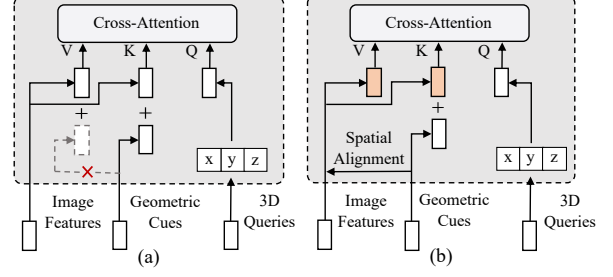


Figure 4. Content composition of query, key and value. (a) structure in PETR; (b) structure in Focal-PETR. The rectangle boxes in orange represent features after spatial alignment.

instance-guided way and decouple the foreground information into three categories, including semantic discriminability, object integrity and position sensitivity. To select corresponding features respectively, three sub-modules are conducted, which are class-aware, IoU-aware and centroid-aware modules. Several convolutional heads are appended to the image encoder to predict the quality score of attributes (see Figure 2). Focal sampling module is lightweight and has negligible impact on the inference time. The above three down-sampling strategies mainly emphasize the importance of positive samples assignment and loss definition.

**Class-aware Sampling.** In order to exploit the discriminative foreground information as soon as possible, additional 2D object detection task to infer the objectiveness is performed. Specifically, we follow FCOS [36] to supervise the classification score $c$ and the normalized distances $(l^*, t^*, r^*, b^*)$ from the location to the four sides of the bounding box. Beseides, Hugarian matching [2] that considers both classification and location costs is adopted for positive sample selection, which can supervise the network to generate high-recall predictions without post-processing [33,38]. For simplicity, Focal loss [18] and L1 loss are used for classification and regression supervision, respectively.

**IoU-aware Sampling.** Previous works [17] have indicated that the precise estimation of 2D pose attributes is highly related to 3D geometric information. Therefore, only sampling tokens with high classification scores will lead to performance degradation on pose prediction tasks. Therefore, an extra GIoU cost [30] is calculated when matching the positive samples. In addition, the network predicts the IoU quality, which is supervised by positive samples generally. However, the usage of one-to-one assignment leads to unfairness, which means that the background often has high confidence prediction. Inspired by Generalized Focal loss [14], we modify the classification branch to jointly estimate the classification-IoU quality $\mathcal{Q}$. In this way, the quality score of both positive and negative samples can be fairly evaluated, which can be formulated as follows:

$$\mathcal{L}_Q = -|y - \mathcal{Q}|^\beta ((1-y)log(1-\mathcal{Q}) + ylog(\mathcal{Q})) \quad (5)$$

where $y$ is the IoU for positive matching pairs while equals

0 for the negatives. The parameter $\beta$ is the modulating factor to down-weight easy examples [18] (we set it to 2.0 following common practice).

**Centroid-aware Sampling.** PETR [20] infers 3D attributes by learning the relationship between queries and 3D implicit features. This interaction is equivalent to learn the projection process from queries' reference point to image plane, as shown in Figure 1. We conclude that the accurate distillation of the projected 2.5D centers can help the network locate the objects. Therefore, we set an auxiliary task of 2.5D center offset for each location on feature map and supervise it with L1 loss. Besides, a key point prediction network is trained to give a high confidence on features around centroid. The Gaussian heatmap $\mathcal{H}$ is used to define the ground truth of 2.5D centers $c$:

$$\mathcal{H} = \exp(-\frac{(x - c_x)^2 + (y - c_y)^2}{2\delta_{\min(l^*, t^*, r^*, b^*)}}) \qquad (6)$$

where $c_x$ and $c_y$ are rounding coordinates of the 2.5D center on the feature map. $\delta_{\min(l^*, t^*, r^*, b^*)}$ is the size-adaptation coefficient controlled by normalized distances $(l^*, t^*, r^*, b^*)$ from the 2.5D center to the four sides of the bounding box. We use the variant of Focal loss proposed in CenterNet [47] to densely supervise the centerness $\rfloor$.

### 3.3. Spatial Alignment Module

In PETR [20], the implicit positon encoding is the only way to distinguish the camera poses, as shown in Equation 3, which causes ambiguity in overlapping areas. In addition, the instance-level down-sampling leads to the loss of global receptive field weakening the depth estimation. In order to align embedding space, the spatial alignment module (see Figure 4) is proposed to transform the 2D features $\boldsymbol{F}_j^{u,\,v}$ from image plane to a 3D unified space:

$$\boldsymbol{F}^*{}_j^{u,\,v} = \mathcal{T}_w(\boldsymbol{r}_j^{u,v},\,\boldsymbol{f}_j)\cdot\boldsymbol{F}_j^{u,\,v} + \mathcal{T}_b(\boldsymbol{r}_j^{u,v},\,\boldsymbol{f}_j) \qquad (7)$$

where $\mathcal{T}_w$ and $\mathcal{T}_b$ are two feedforward networks to encode intrinsic $\boldsymbol{f}_j$ and tracing ray $\boldsymbol{r}_j^{u,v}$ of a specific camera.

### 3.4. Training and Inference.

The priority of feature sampling $\mathcal{P}$ depends on both the quality score $\mathcal{Q}$ and the centerness $\mathcal{C}$ during the training or inference stage:

$$\mathcal{P} = \mathcal{Q}^{\,\alpha}\mathcal{C}^{1-\alpha} \qquad (8)$$

We use $\alpha$ to balance the weights in the sampling process. Following PnP-DETR [40], we dynamically select top $\rho$ ratio features in the training phase and set a fixed threshold $\rho^*$ for inference. The proposed focal sampling module is trained end-to-end and the auxiliary loss $\mathcal{L}_Q$ is plugged to the original 3D branch [20, 43].:

$$\mathcal{L}_Q = \frac{1}{N_{pos}}(\lambda_1\mathcal{L}_Q + \lambda_2\mathcal{L}_{2.5D} + \lambda_3\mathcal{L}_{GIoU} \\ + \lambda_4\mathcal{L}_{ltrb} + \lambda_5\mathcal{L}_{centerness}) \qquad (9)$$

where we adopt the weighted values $\lambda_{1-5}$ as 2, 10, 5, 2, 1, respectively.

## 4. Experiment

### 4.1. Dataset and Metircs

We validate our proposed Focal-PETR on the large-scale nuScenes dataset [1], which is the most frequently used dataset for vision-centric perception with 6 calibrated cameras covering a 360-degree horizontal FOV. This dataset consists of 1000 driving scenes, which are officially separated into 700/150/150 scenes for training, validation and testing. Specifically, each scene is of 20s duration and fully annotated every 0.5s. Following common practice, we report official metrics with NuScenes Detection Score (NDS), mean Average Precision (mAP) and 5 kinds of True Positive (TP) metrics including average translation error (ATE), average scale error (ASE), average orientation error (AOE), average velocity error (AVE), average attribute error (AAE).

### 4.2. Implementation Details

To verify the effectiveness of our method under different pre-training, three types of image encoders are employed: ResNet-50, ResNet-101 [8] and VoVNet-99 [11]. Note that, ResNet-50 is initialized from ImageNet [5] checkpoint, ResNet-101 is initialized from FCOS3D [41] checkpoint, and VoVNet-99 is initialized from DD3D [27] checkpoint. The down-sample stride of image encoder is set to 16. We adopt the same image and BEV data augmentation methods as PETR [10, 20]. The Transformer [37] decoder head consists of 6 layers with 900 object queries.

All models are trained using AdamW [22] optimizer with a total batch size of 16. The learning rate is initialized as 4e-4. There are two variants of our method, Focal-PETR without token sampling and Focal-PETR-H-$\rho^*$ with token sampling ratio $\rho^*$. For results on nuScenes val set, Focal-PETR and Focal-PETR-H-$\rho^*$ are trained for 24 and 60 epochs respectively. For experiments on test set, the training schedules are set up to 100 epochs. Query denoising [12] is not used in all experiments.

### 4.3. State-of-the-art Comparison

We firstly evaluate Focal-PETR on the nuScenes val set and compare it with the state-of-the-art methods listed in Table 1, including DETR3D [43], BEVDet [10], BEV-Former [16], BEVDepth [15], PETR [20], etc. As shown in Table 1, we set relative small input resolutions with ResNet-50 backbone to compare with light-weight models. Focal-PETR with $704 \times 256$ resolution input produces impressive inference speed at 30.0 FPS, which is 1.8 times faster than BEVDet. Remarkably, Focal-PETR with size $800 \times 320$ reaches a desirable trade-off between accuracy and speed. With ResNet-101 backbone, Focal-PETR is

Table 1. Comparison of 3D object detection performance on nuScenes val set. ∗ indicates that the models are trained with CBGS [48] strategy. † notes using the pre-trained FCOS3D backbone. "S" indicates model with a single time stamp input.

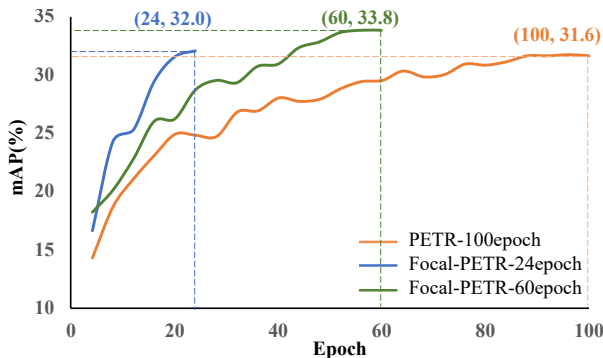| Methods | Backbone | Resolution | mAP↑ | NDS↑ | mATE↓ | mASE↓ | mAOE↓ | mAVE↓ | mAAE↓ | FPS↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| BEVDet∗ [10] | Res-50 | $704 \times 256$ | 0.298 | 0.379 | 0.725 | 0.279 | 0.559 | 0.860 | 0.245 | 16.7 |
| BEVDepth-S [15] | Res-50 | $704 \times 256$ | 0.315 | 0.367 | **0.702** | 0.271 | 0.621 | 1.042 | 0.315 | 15.7 |
| PETR∗ [20] | Res-50 | $1408 \times 512$ | 0.339 | 0.403 | 0.748 | 0.273 | 0.539 | 0.907 | **0.203** | 8.1 |
| Focal-PETR-H-0.33 | Res-50 | $704 \times 256$ | 0.311 | 0.388 | 0.780 | 0.281 | 0.545 | 0.867 | 0.204 | **30.0** |
| Focal-PETR-H-0.5 | Res-50 | $800 \times 320$ | 0.335 | 0.406 | 0.746 | **0.270** | **0.523** | **0.862** | 0.220 | 22.9 |
| Focal-PETR | Res-50 | $800 \times 320$ | 0.320 | 0.381 | 0.788 | 0.278 | 0.595 | 0.893 | 0.228 | 20.2 |
| Focal-PETR | Res-50 | $1408 \times 512$ | **0.363** | **0.414** | 0.760 | 0.279 | 0.538 | 0.876 | 0.224 | 8.5 |
| FCOS3D [41] | Res-101 | $1600 \times 900$ | 0.295 | 0.372 | 0.806 | 0.268 | 0.511 | 1.315 | **0.170** | 1.7 |
| PGD [39] | Res-101 | $1600 \times 900$ | 0.335 | 0.409 | 0.732 | 0.263 | 0.423 | 1.285 | 0.172 | 1.4 |
| EPro-PnP [3] | Res-101 | $1600 \times 900$ | 0.352 | 0.430 | 0.667 | **0.258** | **0.337** | 1.031 | 0.193 | 3.3 |
| DETR3D∗ † [43] | Res-101 | $1600 \times 900$ | 0.349 | 0.434 | 0.716 | 0.268 | 0.379 | 0.842 | 0.200 | 3.7 |
| BEVFormer-S† [16] | Res-101 | $1600 \times 900$ | 0.375 | 0.448 | 0.725 | 0.272 | 0.391 | **0.802** | 0.200 | 3.0 |
| Ego3RT† [23] | Res-101 | $1600 \times 900$ | 0.375 | 0.450 | **0.657** | 0.268 | 0.391 | 0.850 | 0.206 | 1.7 |
| SpatialDETR† [6] | Res-101 | $1600 \times 900$ | 0.351 | 0.425 | 0.772 | 0.274 | 0.395 | 0.847 | 0.217 | 3.5 |
| PETR∗ † [20] | Res-101 | $1408 \times 512$ | 0.366 | 0.441 | 0.717 | 0.267 | 0.412 | 0.834 | 0.190 | 5.7 |
| Focal-PETR-H-0.5† | Res-101 | $1408 \times 512$ | 0.390 | **0.461** | 0.678 | 0.263 | 0.395 | 0.804 | 0.202 | **6.6** |
| Focal-PETR† | Res-101 | $1600 \times 640$ | 0.385 | 0.448 | 0.737 | 0.265 | 0.404 | 0.831 | 0.207 | 4.4 |
| Focal-PETR-H-0.5† | Res-101 | $1600 \times 640$ | **0.393** | 0.457 | 0.695 | 0.264 | 0.383 | 0.850 | 0.206 | 4.9 |



Figure 5. Convergence curves of Focal-PETR and PETR. Specifically, we denote the mAP of the last epoch for each curve.

trained with $1408 \times 512$ resolution to fairly compare with PETR. We can see that our method exceeds PETR by 2.4% in mAP and 2.0% in NDS though PETR is trained with CBGS strategy. With a larger ($1600 \times 640$) resolution input, Focal-PETR outperforms Ego3RT, the state-of-the-art method with ($1600 \times 900$) resolution input, by 1.8% in mAP and 0.7% in NDS respectively.

As shown in Table 2, we also conduct experiments on nuScenes test set and Focal-PETR produces the outstanding results both on mAP and NDS. On the ResNet-101 backbone, Focal-PETR outperforms BEVFormer, the state-of-the-art method, by 1.7% in mAP and 2.4% NDS. It is also noteworthy that with VoVNet-99 [11] backbone, Focal-PETR exceeds PETR by 2.4% in mAP and 1.2% in NDS.

For convergence speed, we train Focal-PETR with 24, 60 epochs and PETR with 100 epochs respectively. Both two

methods are evaluated every 4 epochs on nuScenes val set visualized in Figure 5. The results show that Focal-PETR achieces higher mAP (32.0% vs 31.6%) than PETR even with 3x fewer training epochs. Notably, when comparing the results in Table 1 and Table 2, Focal-PETR achieves better performance than PETR with different backbone pre-training.

## 4.4. Ablation Study

In this section, we conduct experimental analysis on important components of our methods. All models are trained for 24 epochs without CBGS [48] and evaluated on nuScenes val set. More experimental results are presented in the supplementary.

**Analysis of Focal Sampling.** The results in Table 7 show that each component in focal sampling contributes to the performance improvement. Notably, our method with additional instance-guided supervision, outperforms the unsupervised method PnP-DETR, which is inconsistent with the experience in 2D paradigm [32, 40]. One possible reason is that the decoder-only architecture is weak in casting potential objects. Under this design, the foreground token is more required. In addition, the centroid-aware sampling brings much more improvement than IoU-aware sampling, which is not consistent with the conclusion drawn in Generalized Focal Loss [14]: IoU always performs better than centerness as a measurement of localization quality. This implies that the accurate estimation of the centroid in image plane is crucial for view transformer to learn the implicit 3D-to-2D projection.

We also conduct experiments on the influence of differ-

6

Table 2. Comparison of 3D object detection performance on nuScenes test set.

| Methods | Backbone | mAP↑ | NDS↑ | mATE↓ | mASE↓ | mAOE↓ | mAVE↓ | mAAE↓ |
|---|---|---|---|---|---|---|---|---|
| FCOS3D [41] | Res-101 | 0.358 | 0.428 | 0.690 | 0.249 | 0.452 | 1.434 | **0.124** |
| EPro-PnP [3] | Res-101 | 0.373 | 0.453 | 0.605 | **0.243** | **0.359** | 1.067 | **0.124** |
| PGD [39] | Res-101 | 0.386 | 0.448 | 0.626 | 0.245 | 0.451 | 1.509 | 0.127 |
| Ego3RT [23] | Res-101 | 0.389 | 0.443 | **0.599** | 0.268 | 0.470 | 1.169 | 0.172 |
| BEVFormer-S [16] | Res-101 | 0.409 | 0.462 | 0.650 | 0.261 | 0.439 | 0.925 | 0.147 |
| PETR [20] | Res-101 | 0.391 | 0.455 | 0.647 | 0.251 | 0.433 | 0.933 | 0.143 |
| Focal-PETR | Res-101 | **0.426** | **0.486** | 0.617 | 0.250 | 0.398 | **0.862** | 0.146 |
| DD3D [27] | VoV-99 | 0.418 | 0.477 | 0.572 | 0.249 | 0.368 | 1.014 | **0.124** |
| DETR3D [43] | VoV-99 | 0.412 | 0.479 | 0.641 | 0.255 | 0.394 | 0.845 | 0.133 |
| Ego3RT [23] | VoV-99 | 0.425 | 0.473 | 0.549 | 0.264 | 0.433 | 1.014 | 0.145 |
| BEVDet [10] | VoV-99 | 0.424 | 0.488 | **0.524** | **0.242** | **0.373** | 0.950 | 0.148 |
| BEVFormer-S [16] | VoV-99 | 0.435 | 0.495 | 0.589 | 0.254 | 0.402 | 0.842 | 0.131 |
| SpatialDETR [6] | VoV-99 | 0.424 | 0.486 | 0.613 | 0.253 | 0.402 | 0.857 | 0.131 |
| PETR [20] | VoV-99 | 0.441 | 0.504 | 0.593 | 0.249 | 0.383 | **0.808** | 0.132 |
| Focal-PETR | VoV-99 | **0.465** | **0.516** | 0.578 | 0.247 | 0.390 | 0.817 | 0.135 |

Table 3. Ablation study of Focal-PETR with different sampling strategies as the sampling ratio is 0.25. Correspondingly, We report the mAP and DNS metrics here. Cls and Ctr denote classification and centerness scores respectively.

| | Cls | IoU | Ctr | PnP [40] | mAP↑ | NDS↑ |
|---|---|---|---|---|---|---|
| (1) | - | - | - | ✓ | 0.273 | 0.335 |
| (2) | ✓ | - | - | - | 0.301 | 0.361 |
| (3) | ✓ | ✓ | - | - | 0.303 | 0.356 |
| (4) | ✓ | ✓ | ✓ | - | **0.313** | **0.372** |

Table 4. Ablation study with different sampling ratio. Mem. indicates consumption of GPU memory here. To be noted, only computational costs (FLOPs) belonging to detection head are counted. Results with ∗ is obtained with vanillia PETR [20].

| ratio | mAP↑ | NDS↑ | FLOPs(G) | Mem.(G) | FPS |
|---|---|---|---|---|---|
| 0.25 | 0.313 | 0.372 | **24.1(-44.0%)** | **3.6(-43.8%)** | **24.0** |
| 0.5 | 0.319 | 0.378 | 30.5(-29.1%) | 4.5(-29.7%) | 22.9 |
| 0.75 | **0.320** | **0.379** | 34.8(-14.1%) | 5.5(-14.3%) | 21.2 |
| 1.0 | **0.320** | **0.379** | 40.1(+10.2%) | 6.4(+1.6%) | 20.7 |
| 1.0∗ | 0.286 | 0.339 | 36.5 | 6.3 | 20.7 |

Table 5. Different designs of spatial alignment module. We analyze various network designs (module) and feature selection (content). The ray representation only considers the direction of pixel rays, while cone includes the volume viewed by each ray.

| module | content | mAP↑ | NDS↑ |
|---|---|---|---|
| - | - | 0.314 | 0.362 |
| pos. | ray | 0.303 | 0.350 |
| ours. | ray | 0.319 | 0.371 |
| SE [9] | cone | 0.316 | 0.380 |
| ours. | cone | **0.320** | **0.381** |

Table 6. Performance Comparison of different sampling ratio on previous and current frame for temporal extension. We additionally report the results of PETRv2 [21] and marked with ∗.

| previous | current | mAP↑ | NDS↑ | mAVE↓ |
|---|---|---|---|---|
| 0.2 | 0.5 | 0.337 | 0.437 | 0.443 |
| 0.2 | 1.0 | 0.339 | 0.439 | 0.444 |
| 0.5 | 1.0 | 0.340 | **0.442** | 0.426 |
| 1.0 | 1.0 | **0.341** | **0.442** | **0.425** |
| 1.0∗ | 1.0∗ | 0.306 | 0.412 | 0.471 |

ent sampling ratios in Table 4. We firstly compare the results of our method and PETR without sampling (1.0 sampling ratio). It shows that our method achieves 3.4% and 4.0% improvement in mAP and NDS. Note that, the extra 3.7 GFLOPs introduced by the proposed sampling module can be discarded when using 1.0 sampling ratio. Further, when conducting experiments with different sampling ratios in Focal-PETR, they have a more significant influence on computation metrics. The results show that with 0.25 sampling ratio, the model reduces the FLOPs and memory

consumption by absolutely 16.0 G and 2.8 G compared with 1.0 sampling ratio. That is to say, it only sacrifices 0.7% mAP and 0.7% NDS for nearly 43.8% fewer memory costs and 44.0% fewer decoder FLOPs.

**Effectiveness of Spatial Alignment.** As shown in Table 9, different designs for spatial alignment module lead to large variance in detection accuracy. Simply applying the same position encoding scheme for key-value pairs will impair the performance. It can be seen that our design performs better than the commonly used SE-like architec-
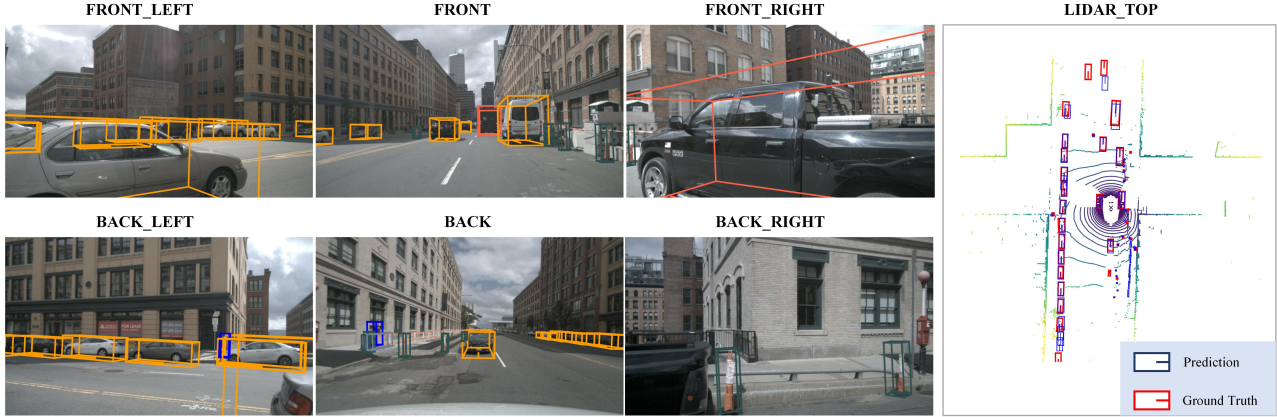
Figure 6. Qualitative detection results of Focal-PETR on nuScenes val set. We show 3D bounding boxes predicted both in multi-camera images and bird's eye view. In multi-camera images, 3D boxes in different colors note different classifications. While in bird's eye view, the ground-truthes are drawn in red to be distinguishable with our predictions in blue.
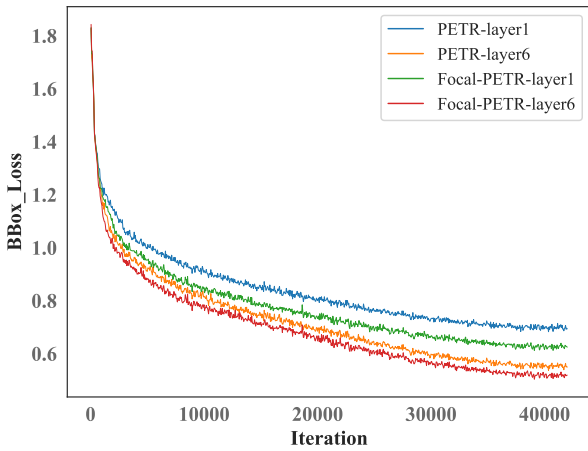


Figure 7. Comparison of L1 losses in the 1st and the 6th layer.

ture [9], indicating that the modulation of both weight and bias (see Equation 7) can achieve more flexible spatial transformation. Furthermore, frustum cone provides a more precise geometric prior than pixel ray to perceive the 3D scene, and it brings significantly 0.9% NDS improvement.

**Extension in Temporal Modeling.** Existing work [15, 16] has demonstrated that temporal clues can assist network in velocity estiamtion and highly occluded objects detection. To verify the scalability of our method in temporal modeling, we conduct a detail analysis on various sampling ratio of different time stamps. As shown in Table 6, the sampling of past image features have a little impact on the final predictions. Only with 20% tokens of historical features reaches almost equivalent performance (mAP of 33.9% vs. 34.1%, NDS of 43.9% vs 44.2%). Remark that, our semantics and spatial alignment strategies still bring great performance improvement. Compared with the PETRv2 baseline, the imporvement in mAP is noticeable of 3.5%, and the performance gains in NDS is 3.0%. This phenomenon also im-

plies the potential of our method in temporal modeling.

### 4.5. Visualization and Analysis

We visualize the detection results of Focal-PETR, as illustrated in Figure 14. We employ Res101 as backbone with $1600 \times 640$ resolution input. The sampling ratio is set as 0.5. Our model is capable of capturing small or dense objects.

To quantitatively prove that Focal-PETR can efficiently focus on the foreground tokens, we show the L1 loss of bounding boxes predicted by the 1st and the 6th decoder layers in Figure 7. Note that, in the 1st decoder layer, the loss of Focal-PETR is significantly less than PETR, indicating that bounding boxes predicted by our method is closer to the objects in the shallow decoder layer.

## 5. Conclusion

We propose Focal-PETR, a multi-camera 3D detection method that mitigates the semantic ambiguity and spatial misalignment of the implicit paradigm. Considering that object detection inherently focuses on foreground information, Focal-PETR takes instance-guided supervision to select discriminative image tokens. These tokens are semantically centralized, which is helpful for the detection head to quickly locate foreground instances. The proposed spatial alignment module enhances the search sensitivity of object query by introducing precise geometric representation. Extensive experiments on the large-scale nuScenes benchmark demonstrate that Focal-PETR achieves state-of-the-art performance and superior efficiency.

Although the proposed Focal-PETR achieves high efficiency in multi-camera 3D object detection, our sampling strategy ignores the representation of map elements, which may hinder the joint learning of object detection and high-quality map segmentation. For future work, we will explore advanced techniques to mitigate this limitation.

## 6. Acknowledgement

## 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 2, 5, 13

[2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 3

[3] Hansheng Chen, Pichao Wang, Fan Wang, Wei Tian, Lu Xiong, and Hao Li. Epro-pnp: Generalized end-to-end probabilistic perspective-n-points for monocular object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2781–2790, 2022. 6, 7

[4] Xiyang Dai, Yinpeng Chen, Jianwei Yang, Pengchuan Zhang, Lu Yuan, and Lei Zhang. Dynamic detr: End-to-end object detection with dynamic attention. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2988–2997, 2021. 3

[5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 5

[6] Simon Doll, Richard Schulz, Lukas Schneider, Viviane Benzin, Markus Enzweiler, and Hendrik PA Lensch. Spatialdetr: Robust scalable transformer-based 3d object detection from multi-view camera images with global cross-sensor attention. 2, 3, 6, 7

[7] Ziteng Gao, Limin Wang, Bing Han, and Sheng Guo. Adamixer: A fast-converging query-based object detector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5364–5373, 2022. 3

[8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5

[9] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. 7, 8, 12

[10] Junjie Huang, Guan Huang, Zheng Zhu, and Dalong Du. Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. *arXiv preprint arXiv:2112.11790*, 2021. 1, 2, 5, 6, 7

[11] Youngwan Lee, Joong-won Hwang, Sangrok Lee, Yuseok Bae, and Jongyoul Park. An energy and gpu-computation efficient backbone network for real-time object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 0–0, 2019. 5, 6

[12] Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M Ni, and Lei Zhang. Dn-detr: Accelerate detr training by introducing query denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13619–13627, 2022. 3, 5, 11, 12

[13] Qi Li, Yue Wang, Yilun Wang, and Hang Zhao. Hdmapnet: An online hd map construction and evaluation framework. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 4628–4634. IEEE, 2022. 1

[14] Xiang Li, Wenhai Wang, Lijun Wu, Shuo Chen, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss: Learning qualified and distributed bounding boxes for dense object detection. *Advances in Neural Information Processing Systems*, 33:21002–21012, 2020. 3, 4, 6

[15] 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 preprint arXiv:2206.10092*, 2022. 1, 5, 6, 8

[16] 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. *arXiv preprint arXiv:2203.17270*, 2022. 1, 2, 5, 6, 7, 8

[17] Qing Lian, Peiliang Li, and Xiaozhi Chen. Monojsg: Joint semantic and geometric cost volume for monocular 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1070–1079, 2022. 4

[18] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 3, 4, 5

[19] Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang. Dab-detr: Dynamic anchor boxes are better queries for detr. *arXiv preprint arXiv:2201.12329*, 2022. 3

[20] Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr: Position embedding transformation for multi-view 3d object detection. *arXiv preprint arXiv:2203.05625*, 2022. 1, 2, 3, 4, 5, 6, 7, 11, 13

[21] Yingfei Liu, Junjie Yan, Fan Jia, Shuailin Li, Qi Gao, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petrv2: A unified framework for 3d perception from multi-camera images. *arXiv preprint arXiv:2206.01256*, 2022. 1, 2, 3, 7

[22] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5

[23] Jiachen Lu, Zheyuan Zhou, Xiatian Zhu, Hang Xu, and Li Zhang. Learning ego 3d representation as ray tracing. *arXiv preprint arXiv:2206.04042*, 2022. 6, 7

[24] Xinzhu Ma, Wanli Ouyang, Andrea Simonelli, and Elisa Ricci. 3d object detection from images for autonomous driving: a survey. *arXiv preprint arXiv:2202.02980*, 2022. 1

[25] Yuexin Ma, Tai Wang, Xuyang Bai, Huitong Yang, Yuenan Hou, Yaming Wang, Yu Qiao, Ruigang Yang, Dinesh Manocha, and Xinge Zhu. Vision-centric bev perception: A survey. *arXiv preprint arXiv:2208.02797*, 2022. 1

[26] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3651–3660, 2021. 3

[27] Dennis Park, Rares Ambrus, Vitor Guizilini, Jie Li, and Adrien Gaidon. Is pseudo-lidar needed for monocular 3d object detection? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3142–3152, 2021. 5, 7

[28] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *European Conference on Computer Vision*, pages 194–210. Springer, 2020. 1, 2

[29] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 3

[30] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 658–666, 2019. 4

[31] Thomas Roddick, Alex Kendall, and Roberto Cipolla. Orthographic feature transform for monocular 3d object detection. *arXiv preprint arXiv:1811.08188*, 2018. 2

[32] Byungseok Roh, JaeWoong Shin, Wuhyun Shin, and Saehoon Kim. Sparse detr: Efficient end-to-end object detection with learnable sparsity. *arXiv preprint arXiv:2111.14330*, 2021. 3, 6

[33] Peize Sun, Yi Jiang, Enze Xie, Wenqi Shao, Zehuan Yuan, Changhu Wang, and Ping Luo. What makes for end-to-end object detection? In *International Conference on Machine Learning*, pages 9934–9944. PMLR, 2021. 4

[34] Zhiqing Sun, Shengcao Cao, Yiming Yang, and Kris M Kitani. Rethinking transformer-based set prediction for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3611–3620, 2021. 3

[35] Yunlei Tang, Sebastian Dorn, and Chiragkumar Savani. Center3d: Center-based monocular 3d object detection with joint depth understanding. In *DAGM German Conference on Pattern Recognition*, pages 289–302. Springer, 2020. 4

[36] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9627–9636, 2019. 3, 4

[37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2, 3, 4, 5

[38] Jianfeng Wang, Lin Song, Zeming Li, Hongbin Sun, Jian Sun, and Nanning Zheng. End-to-end object detection with fully convolutional network. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15849–15858, 2021. 4

[39] Tai Wang, ZHU Xinge, Jiangmiao Pang, and Dahua Lin. Probabilistic and geometric depth: Detecting objects in perspective. In *Conference on Robot Learning*, pages 1475–1485. PMLR, 2022. 6, 7

[40] Tao Wang, Li Yuan, Yunpeng Chen, Jiashi Feng, and Shuicheng Yan. Pnp-detr: Towards efficient visual analysis with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4661–4670, 2021. 3, 5, 6, 7

[41] Tai Wang, Xinge Zhu, Jiangmiao Pang, and Dahua Lin. Fcos3d: Fully convolutional one-stage monocular 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 913–922, 2021. 5, 6, 7

[42] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8445–8453, 2019. 2

[43] 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 *Conference on Robot Learning*, pages 180–191. PMLR, 2022. 1, 2, 5, 6, 7

[44] Yingming Wang, Xiangyu Zhang, Tong Yang, and Jian Sun. Anchor detr: Query design for transformer-based object detection. *arXiv preprint arXiv:2109.07107*, 2021. 3

[45] Shengkai Wu, Xiaoping Li, and Xinggang Wang. Iou-aware single-stage object detector for accurate localization. *Image and Vision Computing*, 97:103911, 2020. 3

[46] Brady Zhou and Philipp Krähenbühl. Cross-view transformers for real-time map-view semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13760–13769, 2022. 1, 2

[47] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. 5

[48] Benjin Zhu, Zhengkai Jiang, Xiangxin Zhou, Zeming Li, and Gang Yu. Class-balanced grouping and sampling for point cloud 3d object detection. *arXiv preprint arXiv:1908.09492*, 2019. 2, 6, 12

[49] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 3

# Supplementary Material

## A. Overview

This document provides more details of network architecture, additional experimental results and qualitative results of ablation study and visualization to the main paper.

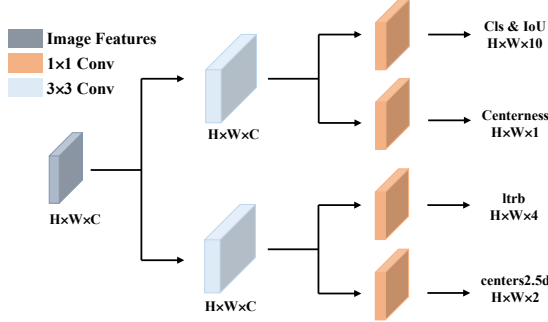## B. Detailed Network Architecture

### B.1. Focal Sampling



Figure 8. The network design of sampling module. Each image is fed into two light-weight convolution layers to predict the quality scores and attributes. For the supervision of Focal sampling module, please refer to Method in main text.
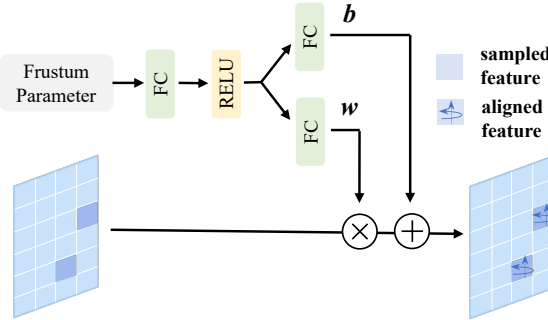
### B.2. Spatial Alignment Module



Figure 9. The network design of spatial alignment module. In order to embed spatial information into image features, we use two MLPs to encode parameterized frustum cones. The obtained weight $w$ and bias $b$ are applied to re-weight the image features.

## C. Additional Experimental Results

We train PETR [20] and Focal PETR with query denosing [12] under the same setting. The training loss and evaluation results are shown in Figure 10 and Figure 11. It can be seen that our method is superior to query denoising in
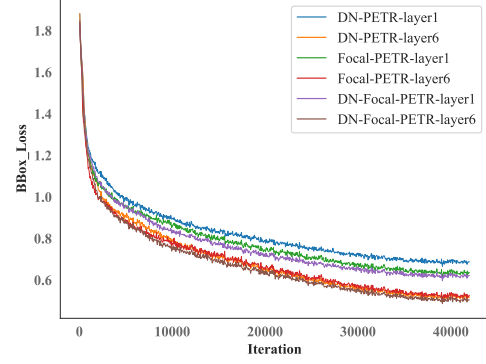


Figure 10. Comparison of regression losses in the 1st and the 6th layers. Note that "DN" indicates models trained with query denoising.
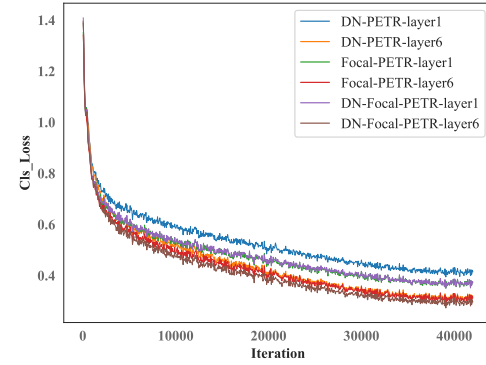


Figure 11. Comparison of classification losses in the 1st and the 6th layers. Note that "DN" indicates models trained with query denoising.

multi camera 3D detection. Focal PETR can still benefit from query denosing.

More experimental and visualization results are as follows.

Table 7. Complete results of Focal-PETR with different sampling strategies as the sampling ratio is 0.25. Cls and Ctr denote classification and centerness scores respectively. We also provide the evaluation results of the model only using C5 feature without sampling. It has worse performance and comparable computational cost compared with the models using P4 feature with 0.25 sampling ratio.

|     | Cls | IoU | Ctr | feature | mAP↑ | NDS↑ | mATE↓ | mASE↓ | mAOE↓ | mAVE↓ | mAAE↓ |
|-----|-----|-----|-----|---------|------|------|-------|-------|-------|-------|-------|
| (1) | ✓ | - | - | P4 | 0.301 | 0.361 | 0.831 | 0.279 | 0.634 | **0.908** | 0.239 |
| (2) | ✓ | ✓ | - | P4 | 0.303 | 0.356 | 0.833 | 0.279 | 0.640 | 0.965 | 0.239 |
| (3) | ✓ | ✓ | ✓ | P4 | **0.313** | **0.372** | **0.807** | **0.277** | 0.584 | 0.943 | **0.236** |
| (4) | ✓ | ✓ | ✓ | C5 | 0.308 | 0.371 | 0.808 | 0.280 | **0.576** | 0.911 | 0.247 |

Table 8. Complete results of different designs for spatial alignment module. We analyze various network designs (module) and feature selection (content). The ray representation only considers the direction of pixel rays, while cone additionally includes the volume viewed by each ray.

| module | content | mAP↑ | NDS↑ | mATE↓ | mASE↓ | mAOE↓ | mAVE↓ | mAAE↓ |
|--------|---------|------|------|-------|-------|-------|-------|-------|
| - | - | 0.314 | 0.362 | 0.808 | 0.281 | 0.601 | 1.049 | 0.262 |
| pos. | ray | 0.303 | 0.350 | 0.842 | 0.281 | 0.686 | 0.957 | 0.245 |
| ours. | ray | 0.319 | 0.371 | 0.800 | 0.278 | 0.587 | 0.984 | 0.241 |
| SE [9] | cone | 0.316 | 0.380 | 0.803 | 0.278 | **0.575** | **0.885** | 0.242 |
| ours. | cone | **0.320** | **0.381** | **0.791** | **0.276** | 0.607 | **0.885** | **0.232** |

Table 9. Influence of query denoising [12] on Focal-PETR. We compare the results of Focal-PETR and PETR with query denoising under the same setting.

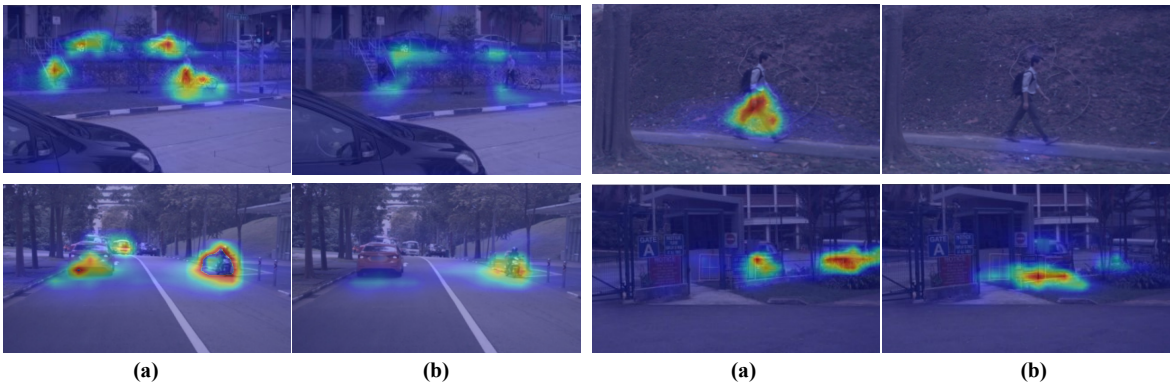| Focal | Denosing | mAP↑ | NDS↑ | mATE↓ | mASE↓ | mAOE↓ | mAVE↓ | mAAE↓ |
|-------|----------|------|------|-------|-------|-------|-------|-------|
| ✓ | - | 0.321 | 0.381 | 0.800 | 0.280 | **0.589** | **0.895** | **0.235** |
| - | ✓ | 0.307 | 0.355 | 0.810 | 0.278 | 0.716 | 0.944 | 0.240 |
| ✓. | ✓ | **0.328** | **0.382** | **0.779** | **0.275** | 0.631 | **0.895** | 0.240 |



(a)　　　　　(b)　　　　　(a)　　　　　(b)

Figure 12. Comparison of attention weight maps for (a) Focal-PETR and (b) PETR with 24 training epochs without CBGS [48]. It can be seen that the attention map of Focal-PETR focuses more on the foreground objects.
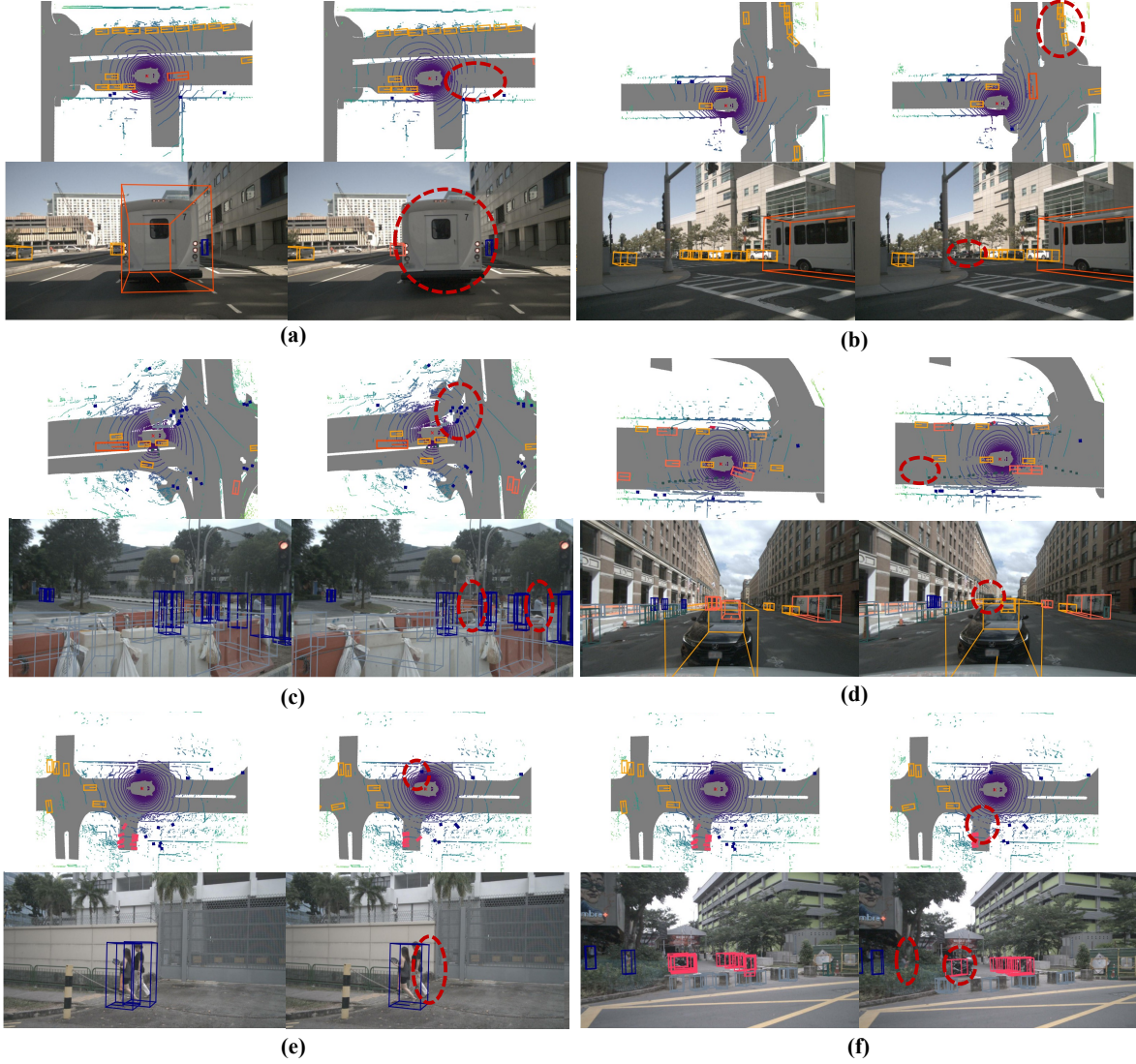
Figure 13. Qualitative detection results of Focal-PETR and PETR [20] on nuScenes [1] val set. We visualize one of six camera views in each scene and compare the results of two methods. The 3D bounding boxes are projected into BEV and image view. Boxes in different colors note different classifications. In particular, We note the sub-optimal predictions of PETR in red circles.



Figure 14. Qualitative detection results of (a) Focal-PETR and (b) PETR on nuScenes val set. We visualize multi-camera views in each scene and compare the results of two methods. Boxes in different colors note different classifications. In particular, we note the sub-optimal predictions of PETR in red circles.

Figure 15. Visualization of sampling locations in multi-camera views with different sampling ratios, including (a) 0.25 sampling ratio; (b) 0.33 sampling raito; (c) 0.5 sampling ratio; (d) 0.75 sampling raito. The yellow dots indicate the sampled locations.