

Polar Parametrization for Vision-based Surround-View 3D Detection

Shaoyu Chen^{1,2}, Xinggang Wang^{1*}, Tianheng Cheng^{1,2},
Qian Zhang², Chang Huang², and Wenyu Liu¹

¹ Huazhong University of Science & Technology

² Horizon Robotics

{shaoyuchen, xgwang, thch, liuwu}@hust.edu.cn

{qian01.zhang, chang.huang}@horizon.ai

Abstract. 3D detection based on surround-view camera system is a critical technique in autopilot. In this work, we present Polar Parametrization for 3D detection, which reformulates position parametrization, velocity decomposition, perception range, label assignment and loss function in polar coordinate system. Polar Parametrization establishes explicit associations between image patterns and prediction targets, exploiting the view symmetry of surround-view cameras as inductive bias to ease optimization and boost performance. Based on Polar Parametrization, we propose a surround-view 3D DEtection TRansformer, named PolarDETR. PolarDETR achieves promising performance-speed trade-off on different backbone configurations. Besides, PolarDETR ranks 1st on the leaderboard of nuScenes benchmark in terms of both 3D detection and 3D tracking at the submission time (Mar. 4th, 2022). Code will be released at <https://github.com/hustvl/PolarDETR>.

Keywords: 3D Detection, Autopilot, Surround-View Perception, View Symmetry, Polar Parametrization

1 Introduction

In the field of autopilot, surround-view camera system has been attached great importance and popularized by both industry and academia, given its low assembly cost and rich semantic information. And 3D detection based on such vision system has become the critical technique of environmental perception. The passed several years have witnessed tremendous progress in vision-based 3D detection [13,12,32,24,16,25,28,7]. Most methods parameterize object’s position in two manners, 1) Image-based Parametrization (Image-based Param. for short) and 2) Cartesian Parametrization (Cartesian Param. for short).

Image-based Parametrization. For Image-based Param. (Fig. 1 left), object’s position is defined in the pixel coordinate frame, parameterized by a 3-tuple (u, v, d) , where (u, v) is the pixel coordinate and d is the object’s depth relative

* Xinggang Wang is the corresponding author.

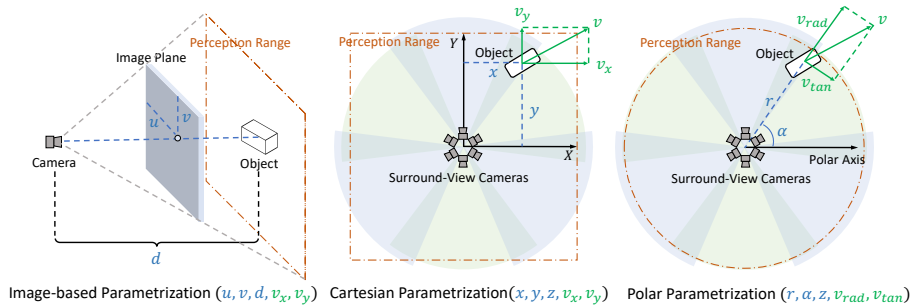


Fig. 1. Illustration of parametrization of objects. Object’s position is parameterized by pixel coordinate with depth (u, v, d) in Image-based Parametrization, by cartesian coordinate (x, y, z) in Cartesian Parametrization, and by polar (cylindrical) coordinate (r, α, z) in Polar Parametrization. Besides, Polar Parametrization decomposes velocity along radial and tangential direction (to v_{rad} and v_{tan}), and reformulates perception range, label assignment and loss function in polar coordinate system. Sector denotes field of view (FoV) of each camera.

to camera. With camera’s extrinsic and intrinsic parameters, we can transform (u, v, d) to a 3D point to localize the object in 3D space. Image-based Param. is widely used in monocular methods [24,16,25,11]. To process surround-view images, they independently take each view’s image as input to regress object’s (u, v, d) . Predicted objects of different views are then projected to the same 3D space for inter-camera merging. Post-processing among views (*e.g.*, NMS) is adopted to filter out duplicated predictions. However, Image-based Param. raises some problems:

- Estimating depth from single image is inherently an ill-posed inverse problem. The predicted depth is of large error and thus object’s positioning accuracy is poor. For the surround-view camera system, adjacent views overlap with each other. The correlation among views can be leveraged to promote perception performance, which is neglected in Image-based Param.
- Post-processing among views is tricky and unstable. When the predictions from different views do not overlap in 3D space, NMS fails to filter out duplicated predictions.

Cartesian Parametrization. For Cartesian Param. (Fig. 1 mid), object’s position is parameterized by 3D cartesian coordinate (x, y, z) . Correspondingly, the perception range is a rectangular region, denoted as $\{(x, y), |x| < X_{max}, |y| < Y_{max}\}$. Cartesian Param. is adopted in some recent studies [7,28]. Compared with monocular methods, they consider the correlation among views, *i.e.*, taking all views together as input to jointly regress object’s 3D coordinate (x, y, z) .

But Cartesian Param. also raise problems. We take a special case for example to illustrate them. As shown in Fig. 2, assume that object A appears in different views at timestamp t_1 and t_2 . They have the same pixel coordinate (u, v) and radial distance d . Their image patterns are also the same.

- Cartesian Param. leads to ambiguity in label assignment because of its rectangular perception range. Ground-truth objects out of the perception range are ignored in label assignment. In the case of Fig. 2, though A_{t_1} and A_{t_2} are at the same distance d and have the same image patterns, A_{t_1} is ignored while A_{t_2} is kept. A_{t_1} and A_{t_2} are not treated equally in the training phase, which is contradictory and harmful to the convergence of network.
- Cartesian Param. neglects the view symmetry of surround-view cameras. From the perspective of machine learning, the detector aims at approximating a function F that maps from image patterns X to the prediction targets Y , *i.e.*, $Y = F(X)$. In the case of Fig. 2, for Image-based Param., A_{t_1} and A_{t_2} have the same image patterns X and prediction targets (u, v, d) , and can be treated as the same mapping in F . A_{t_1} and A_{t_2} are symmetrical in terms of view. Image-based Param. leverages the view symmetry as inductive bias of the detector, making it easier to approximate F . While for Cartesian Param., object A_{t_1} and A_{t_2} correspond to different prediction targets $((x_{t_1}, y_{t_1}, z_{t_1})$ for A_{t_1} and $(x_{t_2}, y_{t_2}, z_{t_2})$ for A_{t_2}). Since the view symmetry is neglected, the mappings for A_{t_1} and A_{t_2} are learned separately. Thus, F is more complicated and optimization becomes harder.

In this work, we propose a surround-view 3D DEtection TRansformer, named PolarDETR. To adapt to the view symmetry of surround-view camera system, we parameterize object’s position by polar (cylindrical)³ coordinate (r, α, z) and decomposes velocity to radial velocity v_{rad} and tangential velocity v_{tan} . Besides, we reformulate perception range, label assignment and loss function in polar coordinate system. We term it as Polar Parametrization (Polar Param. for short).

Polar Param. establishes explicit associations between image patterns and prediction targets. Radial distance r is associated with object’s size in image. Azimuth α is associated with the index of pixel. We add 3D positional encoding to each pixel which contains explicit clues about azimuth. Radial velocity v_{rad} is associated with the changing rate of object’s size. Tangential velocity v_{tan} is associated with the object’s movement in image plane (similar to optical flow). With these explicit associations, the mapping is simplified and the detector achieves better convergence and performance.

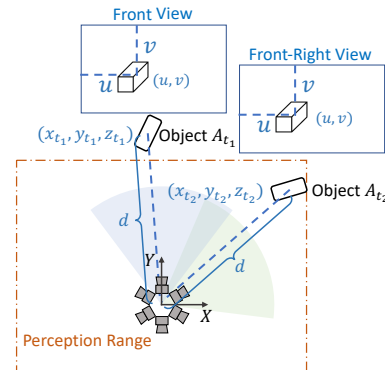


Fig. 2. A special case for illustrating the problems of Cartesian Param. Assume A_{t_1} and A_{t_2} are totally the same in image, except appearing in different views.

³ 'Cylindrical' is more proper for describing 3D space. But considering 'polar' better reflects the insight of this work and Z dimension is not critical, we adopt 'polar' for method description

Besides, PolarDETR enables center-context feature aggregation to enhance the information interaction between object queries and images, and adopts pixel ray as positional encoding to provide 3D spatial priors and help to predict azimuth α .

We evaluate PolarDETR on the challenging nuScenes [1] benchmark. PolarDETR achieves promising performance-speed trade-off on different backbone configurations. And we submit PolarDETR for official evaluation on the nuScenes *test* set. PolarDETR **ranks 1st** on the highly competitive 3D detection and tracking leaderboard at the submission time (Mar. 4th, 2022). Besides, thorough ablation studies are provided to validate the effectiveness of PolarDETR.

2 Related Work

2.1 Polar Coordinate

Polar coordinate is leveraged in some previous works for data representation. In the field of point cloud semantic segmentation, some grid-based methods [31,35] partition 3D space into polar grids and transform LiDAR point clouds into the grid representation, in order to ease the problem of long-tailed distribution. In the field of image instance segmentation, PolarMask [29] formulates the instance mask as a set of contour points represented in the polar coordinate system centered at the instance. Differently, in order to adapt to the view symmetry of surround-view cameras, PolarDETR parameterizes object’s position in polar coordinate system and reformulates label assignment and loss function accordingly.

2.2 DETR-based Object Detection

Recently, DETR [2] formulates object detection as a set prediction problem and exploits a standard transformer [21], which contains a simple encoder-decoder structure without hand-craft designs and achieves promising performance. Based on Hungarian algorithm [20], the set of object queries can be matched with targets in a one-to-one manner, which helps remove the non-maximum suppression (NMS). Deformable DETR [37] motivated by [36] adopts multi-scale deformable attention in the transformer and achieves faster convergence and better performance than DETR. [14,5] also address the convergence problem in DETR by incorporating spatial information into object queries. Several works [26,19] propose to prune the redundant tokens in transformers to minimize the computation cost of DETR. In terms of 3D object detection, DETR3D [28] extends DETR for 3D domain. Inspired by Deformable DETR, DETR3D projects 3D object queries to 2D reference points to aggregate features from all views. DETR3D builds up a simple DETR-based pipeline for 3D object detection, while it suffers from (1) ambiguity in label assignment, (2) neglecting view symmetry and (3) insufficient contextual information. To solve these problems, PolarDETR adopts Polar Parametrization to ease optimization and enables center-context feature aggregation to enhance the feature interaction.

2.3 Vision-based 3D Object Detection

Vision-based 3D detection is a basic perception task in autonomous driving. Early studies are mainly based on KITTI [6] dataset, which provides front-view object annotations. KITTI boosts the development of monocular 3D object detection methods [13,12,32,34,18,8,22,23]. Recently, with nuScenes [1] dataset available, which contains 360° annotations around the ego vehicle, new 3D detection paradigms have been proposed. Some works [24,16,25] still follows the monocular pipeline to detect objects, and then project the multi-view detection results to the same coordinate system and adopt NMS to merge results. DETR3D [28] extends DETR for 3D object detection. BEVDet [7] projects surround-view image features to Bird-Eye-View (BEV) space, and set detection head on BEV features. We presents a new paradigm specially designed for surround-view camera systems, in which the view symmetry is exploited to ease optimization and boost performance.

3 PolarDETR

3.1 Overview

Fig. 3 illustrates the proposed PolarDETR, which follows the DETR [2,37,28] paradigm. Given surround-view images $\mathcal{I} = \{I_1, \dots, I_K\}$ (K denotes the number of views), a shared CNN backbone extracts image features $\mathcal{F}_{\text{img}} = \{F_1, \dots, F_K\}$. A set of object queries $\mathcal{Q} = \{q_1, \dots, q_N\}$ (N denotes the number of queries) is used to detect objects. Specifically, each object query encodes the semantic features and positional information of the corresponding object. And a series of decoder layers aggregate features from surround-view feature maps and update queries iteratively. The feed forward network (FFN) follows the decoder layers and predicts polar box encodings B_{enc} , polar velocity components (v_{rad}, v_{tan}) , and class labels based on queries.

3.2 Polar Parametrization

In PolarDETR, object’s position is parameterized by the polar coordinate. FFN outputs a 9-tuple polar box encoding, denoted as,

$$B_{\text{enc}} = (b_r, b_{\sin \alpha}, b_{\cos \alpha}, b_z, b_l, b_w, b_h, b_{\sin \theta}, b_{\cos \theta}). \quad (1)$$

We then decode B_{enc} to get the predicted box vector B_{pred} represented in polar coordinate system, *i.e.*,

$$\begin{aligned} B_{\text{pred}} &= (r, \sin \alpha, \cos \alpha, z, l, w, h, \sin \theta, \cos \theta), \\ r &= \sigma(b_r) \cdot R_{\text{max}}, \quad z = \sigma(b_z) \cdot (Z_{\text{max}} - Z_{\text{min}}) + Z_{\text{min}}, \\ \sin \alpha &= \frac{b_{\sin \alpha}}{\sqrt{b_{\sin \alpha}^2 + b_{\cos \alpha}^2}}, \quad \cos \alpha = \frac{b_{\cos \alpha}}{\sqrt{b_{\sin \alpha}^2 + b_{\cos \alpha}^2}}, \\ l &= \exp(b_l), \quad w = \exp(b_w), \quad h = \exp(b_h), \\ \sin \theta &= \frac{b_{\sin \theta}}{\sqrt{b_{\sin \theta}^2 + b_{\cos \theta}^2}}, \quad \cos \theta = \frac{b_{\cos \theta}}{\sqrt{b_{\sin \theta}^2 + b_{\cos \theta}^2}}, \end{aligned} \quad (2)$$

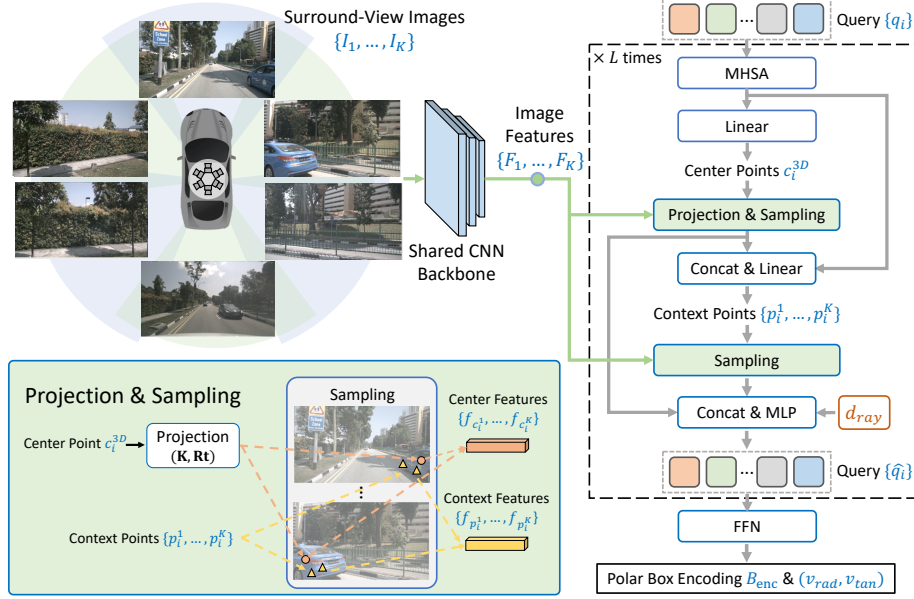


Fig. 3. The framework of PolarDETR. A shared CNN backbone extracts image features $\mathcal{F}_{img} = \{F_1, \dots, F_K\}$ from surround-view images $\mathcal{I} = \{I_1, \dots, I_K\}$. Decoder layers aggregate center and context features $\{f_{c_1^i}, \dots, f_{c_i^k}, f_{p_1^i}, \dots, f_{p_i^k}\}$ from surround-view feature maps and update queries $\{q_i\}$ iteratively. Pixel ray d_{ray} is introduced to provide spatial priors. FFN predicts polar box encodings B_{enc} , polar velocity components (v_{rad}, v_{tan}) , and class labels.

where r, α, z are the radial distance, azimuth and height, indicating the position of object’s geometry center in 3D space. (l, w, h) and θ denote the size and orientation of bounding box, respectively. Z_{max} and Z_{min} denote the perception range along the Z dimension. R_{max} denotes the maximum perception distance. σ is sigmoid function. To make sure the continuity of the regression space [33], we parameterize both α and θ by a 2-D $(\sin(\cdot), \cos(\cdot))$ pair.

For Polar Param., object’s position is decoupled into the radial distance and the azimuth. The radial distance is symmetrical in different views. It’s highly correlated with the object size in image, and can be learned from image patterns. The azimuth is relative to the indices of cameras and pixels which capture the object, and can be learned from the positional encodings. Compared with Cartesian Param., which localizes objects by regressing the cartesian coordinate (x, y) , Polar Param. makes more sense.

Polar Decomposition for Velocity Estimation. Most methods decompose velocity along the cartesian axes and get velocity components v_x and v_y . Differently, PolarDETR decomposes velocity to radial velocity v_{rad} and tangential velocity v_{tan} . Radial velocity v_{rad} is associated with the changing rate of object’s

size. Tangential velocity v_{tan} is associated with the object’s movement in image plane (similar to optical flow). PolarDETR establishes explicit associations between image patterns (input) and velocity (prediction target), resulting in more accurate velocity estimation.

3.3 Decoder Layer

Decoder layers iteratively aggregate features and update queries. Each decoder layer begins with a multi-head self-attention module (MHSA) for inter-query information interaction. Then a linear layer extracts object’s 3D position from each query, *i.e.*,

$$(b_r, b_{\sin \alpha}, b_{\cos \alpha}, b_z) = \text{Linear}(\text{MHSA}(q_i)). \quad (3)$$

We decode $(b_r, b_{\sin \alpha}, b_{\cos \alpha}, b_z)$ with Eq. (2) and get object’s 3D center point $c_i^{3D} = (r, \alpha, z)$.

Center-Context Feature Aggregation. Then we aggregate features from surround-view feature maps through a two-step center-context procedure.

Firstly, we project the 3D center point c_i^{3D} to each image and get a set of 2D center points $\{c_i^1, \dots, c_i^K\}$, which denote object’s centers in all views. *I.e.*,

$$c_i^k = \mathbf{K}^k \cdot \mathbf{Rt}^k \cdot c_i^{3D}, \quad (4)$$

where \mathbf{K}^k and \mathbf{Rt}^k are the projection matrices of view k derived from camera’s intrinsics and extrinsics respectively. Based on $\{c_i^1, \dots, c_i^K\}$, we sample center features $\{f_{c_i^1}, \dots, f_{c_i^K}\}$ from surround-view feature maps $\{F_1, \dots, F_K\}$ through bilinear interpolation. A 3D point may be invisible in some views and the projected 2D points are out of range of the images. In this case, the corresponding point features are set to zero.

Considering that center features are not informative enough for localizing objects, we further include context features to enhance the interaction between queries and surround-view features. Specifically, based on both the center features $\{f_{c_i^1}, \dots, f_{c_i^K}\}$ and the query embeddings q_i , we generate a set of context points $\{p_i^1, \dots, p_i^K\}$ by predicting the offset relative to the center points. *I.e.*,

$$\Delta u_i^k, \Delta v_i^k = \text{Linear}(\text{Concat}(f_{c_i^k}, q_i)), \quad p_i^k = c_i^k + (\Delta u_i^k, \Delta v_i^k). \quad (5)$$

For clarity, in notation we only consider one context point for each view to avoid introducing extra subscripts. In our implementation, more context points are adopted for better context modeling. Similarly, based on context points $\{p_i^1, \dots, p_i^K\}$, we sample context features $\{f_{p_i^1}, \dots, f_{p_i^K}\}$ from surround-view image features $\{F_1, \dots, F_K\}$ through bilinear interpolation.

Pixel Ray. Motivated by [15,27], we introduce pixel ray as 3D spatial priors. As shown in Fig. 4, pixel ray travels from camera’s optical center through the pixel to the corresponding 3D point. It encodes the correspondence between 2D image pixel and 3D point, and contains explicit clues about azimuth. We leverage pixel

ray as extra positional encodings. Specifically, for each center or context point, the unit direction vector of corresponding pixel ray d_{ray} is concatenated with point features in the channel dimension (refer to Eq. (6)).

Query Update. Then we aggregate center and context features with d_{ray} to update query embeddings, *i.e.*,

$$\hat{q}_i = \text{MLP}(\text{Concat}(\{f_{c_i^1}, \dots, f_{c_i^K}, f_{p_i^1}, \dots, f_{p_i^K}\}, d_{\text{ray}})) + q_i. \quad (6)$$

The updated query embeddings encode more accurate positional information of the object and contribute to better feature aggregation in the next decoder layer.

3.4 Perception Range, Label Assignment and Loss Function

Polar Param. reformulates perception range, label assignment and loss function in polar coordinate system.

Perception Range. As discussed above, for Cartesian Param. the rectangular perception range introduces ambiguity into label assignment. For Polar Param., the concerned perception range is a circular region with radius R_{max} centered at the ego vehicle. It avoids ambiguity and fits with the common sense. It’s worth noting that the evaluation region of nuScenes benchmark [1] corresponds with the circular perception range of Polar Param.

Label Assignment. Label assignment between predictions and ground-truth (GT) objects is also performed based on polar coordinate. 3D box annotations are first transformed to polar coordinate representation, *i.e.*,

$$B_{gt} = (\bar{r}, \sin \bar{\alpha}, \cos \bar{\alpha}, \bar{z}, \bar{l}, \bar{w}, \bar{h}, \sin \bar{\theta}, \cos \bar{\theta}). \quad (7)$$

Then we adopt bipartite matching to uniquely assign predictions with ground-truth boxes. Assume there exist N predictions and M GTs. The pair-wise matching cost between prediction i and GT j is defined as,

$$\begin{aligned} C(i, j) &= C_{\text{cls}}(i, j) + C_{\text{box}}(i, j), \\ C_{\text{box}}(i, j) &= |r - \bar{r}| + k_{\text{scaling}} \cdot (|\sin \alpha - \sin \bar{\alpha}| + |\cos \alpha - \cos \bar{\alpha}|), \end{aligned} \quad (8)$$

where $C_{\text{cls}}(i, j)$ is the class term inherited from DETR [2]. $C_{\text{box}}(i, j)$ is the box term based on Polar Param. k_{scaling} is the scaling factor, which numerically scales up the value of azimuth. In radial direction, $r \in [0, R_{\text{max}}]$ (R_{max} is set to 50 in nuScenes). In tangential direction, $\sin \alpha, \cos \alpha \in [-1, 1]$. Their values differ by more than an order of magnitude. Without k_{scaling} , the radial direction will dominate the assignment and there would be significant assignment error in the tangential direction. By computing the cost of each prediction-GT pair, we get the cost matrix $\mathbf{H} \in \mathbb{R}^{M \times N}$. Then Hungarian algorithm [20] is adopted to find the optimal assignment.

Loss Function. We reformulate the loss function based on polar coordinate. The bipartite matching loss consists of two parts: a focal loss [10] for class labels and a \mathcal{L}_1 loss for polar box parameters $(r, \sin \alpha, \cos \alpha, z, l, w, h, \sin \theta, \cos \theta)$ and polar velocity components (v_{rad}, v_{tan}) . The terms about azimuth, *i.e.*, $\sin \alpha$ and $\cos \alpha$, are also scaled up by $k_{scaling}$ to balance the error distribution between tangential and radial direction.

3.5 Temporal Information

Temporal information is important for velocity estimation and occlusion cases. We straightforwardly and elegantly extend PolarDETR to PolarDETR-T, which takes streaming camera frames as input. The 3D center point c_i^{3D} are projected to past frames for fetching image features. Taking frame $t - n$ for example,

$$c_i^{k(t-n)} = \mathbf{K}^k \cdot \mathbf{Rt}^k \cdot \mathbf{Pose}^{(t-n)} \cdot c_i^{3D}, \quad (9)$$

where $c_i^{k(t-n)}$ denotes the corresponding 2D point in the frame $t - n$, and $\mathbf{Pose}^{(t-n)}$ denotes the pose transformation matrix which reflects the movement of the ego-vehicle in the time interval $[t - n, t]$. We sample center and context features from past frames in the same manner of current frame as mentioned above. And all the sampled features (both current and past ones) are aggregated together to update query embeddings.

For a running system with streaming input, for efficient inference, we can cache image feature maps of the past frames. For each moment, we only need to process images of current frame t with backbone to get $\mathcal{F}_{img}^{(t)}$. Feature maps of past frames $\{\mathcal{F}_{img}^{(t-1)}, \mathcal{F}_{img}^{(t-2)}, \dots\}$ are directly fetched from cache, avoiding duplicated computation. Since most computation cost lies in backbone, PolarDETR-T can run at a similar FPS compared with PolarDETR.

4 Experiments

4.1 Dataset

We validate the effectiveness of PolarDETR on the large-scale nuScenes [1] dataset, which is currently the most popular benchmark for vision-based methods. NuScenes contains 1000 driving sequences, with 700, 150 and 150 sequences for training, validation and testing, respectively. Each sequence is approximately

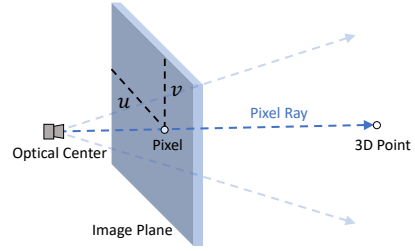


Fig. 4. Illustration of pixel ray. Pixel ray travels from camera’s optical center through the pixel to the corresponding 3D point, which encodes the correspondence between 2D image pixel and 3D point.

20-second long and provides 6 surround-view images per frame with the resolution of 1600×900 . We submit *test* set results to the online server for official evaluation to get the leaderboard results. And other experiments are evaluated on the *val* set.

4.2 Experimental Settings

We implement PolarDETR with PyTorch [17] framework and MMDetection3D [3] toolbox. The results are based on three backbone configurations. We adopt ImageNet [4] pretraining for ResNet-50, FCOS3D [24] pretraining for ResNet-101, and DD3D [16] pretraining for VoVNet. Unless specified, we adopt ResNet-50 as backbone for ablation experiments. We train PolarDETR on eight RTX3090 GPUs with the total batch size 8. Inference speeds of all models are measured on one RTX3090 GPU. We adopt mixed precision (float32 and float16) to accelerate training and disable it when measuring the inference speeds. We do not use test-time augmentation or model ensemble during inference. We adopt the implementation of [36] for the context point sampling. The number of attention heads is set to 8 as the common practice. Results of PolarDETR-T are based on only one past frame. Using more past frames are feasible and brings more gain. More details about experimental settings will be available in the code.

In addition, we simply extend PolarDETR for 3D object tracking by leveraging the tracking-by-detection algorithm [30]. Specifically, we project objects of the current frame back to the previous frame with the predicted velocity, and then match them with the tracked objects by closest distance matching. We evaluate the tracking performance on nuScenes tracking benchmark.

4.3 Metrics

Detection. For 3D object detection, we report the official metrics: NuScenes Detection Score (NDS), mean Average Precision (mAP) and true positive (TP) metrics (including Average Translation Error (ATE), Average Scale Error (ASE), Average Orientation Error (AOE), Average Velocity Error (AVE) and Average Attribute Error (AAE)). The main metric, NDS, is a weighted sum of the other metrics for comprehensively judging the detection capacity, defined as,

$$\text{NDS} = \frac{1}{10} [5\text{mAP} + \sum_{\text{mTP} \in \text{TP}} (1 - \min(1, \text{mTP}))]. \quad (10)$$

Tracking. For 3D object tracking, we report the official metrics: Average Multi Object Tracking Accuracy (AMOTA), Average Multi Object Tracking Precision (AMOTP), False Positives (FP), False Negatives (FN), Identity Switches (IDS), Track Initialization Duration (TID) and Longest Gap Duration (LGD). AMOTA serves as the main metric, which penalizes ID switches, false positive, and false negatives and is averaged among various recall thresholds.

Table 1. Performance comparison (*val* set). PolarDETR achieves promising performance-speed trade-off. FPS are measured on one RTX3090 GPU.

Method	Backbone	Input	NDS↑	FPS↑	mAP↑	mATE↓	mASE↓	mAOE↓	mAVE↓	mAAE↓
DETR3D [28]	Res-50	1600 × 900	0.373	6.3	0.302	0.811	0.282	0.493	0.979	0.212
BEVDet [7]	Res-50	704 × 256	0.372	7.8	0.286	0.724	0.278	0.590	0.873	0.247
BEVDet [7]	Res-50	1056 × 384	0.381	4.2	0.304	0.719	0.272	0.555	0.903	0.257
PolarDETR	Res-50	1600 × 900	0.409	6.0	0.338	0.768	0.284	0.443	0.883	0.221
PolarDETR-T	Res-50	1600 × 900	0.458	~6.0	0.354	0.748	0.277	0.432	0.539	0.197
FCOS3D [24]	Res-101	1600 × 900	0.372	1.7	0.295	0.806	0.268	0.511	1.315	0.170
PGD [25]	Res-101	1600 × 900	0.409	1.4	0.335	0.732	0.263	0.423	1.285	0.172
DETR3D [28]	Res-101	1600 × 900	0.425	3.7	0.346	0.773	0.268	0.383	0.842	0.216
BEVDet [7]	Res-101	704 × 256	0.373	7.1	0.288	0.722	0.269	0.538	0.911	0.270
BEVDet [7]	Res-101	1056 × 384	0.389	3.8	0.317	0.704	0.273	0.531	0.940	0.250
PolarDETR	Res-101	1600 × 900	0.444	3.5	0.365	0.742	0.269	0.350	0.829	0.197
PolarDETR-T	Res-101	1600 × 900	0.488	~3.5	0.383	0.707	0.269	0.344	0.518	0.196
DETR3D [28]	VoVNet	1600 × 900	0.509	2.7	0.445	0.687	0.261	0.271	0.727	0.191
PolarDETR	VoVNet	1600 × 900	0.532	2.5	0.462	0.628	0.262	0.263	0.658	0.180

4.4 Main Results

Performance Comparison. In Tab. 1, we compare PolarDETR with other state-of-the-art methods on three backbone configurations. On both ResNet-50 and ResNet-101, PolarDETR significantly outperforms DETR3D [28] and BEVDet [7] while achieving comparable inference speed. And on ResNet-101, PolarDETR outperforms FCOS3D [24] and PGD [25] in terms of both performance and speed. On VoVNet, PolarDETR significantly outperforms DETR3D [28]. With temporal information, PolarDETR-T achieves much higher results than PolarDETR, especially in terms of mAVE.

NuScenes Leaderboard. Tab. 2 shows the 3D detection leaderboard of nuScenes benchmark [1] for camera modality at the submission time (Mar. 4th, 2022). For fair comparison, we adopt the pretrained VoVNet [9] (refer to DD3D [16]) as backbone, the same with the other top methods on the leaderboard. PolarDETR **ranks 1st** on this highly competitive leaderboard. And it’s worth noting that our implementation is compact and elegant, without test-time augmentation, multi-model ensemble or other tricks.

Tab. 3 shows the tracking leaderboard of nuScenes benchmark [1] for camera modality (Mar. 4th, 2022). At the submission time, PolarDETR **ranks 1st** on the leaderboard and outperforms other methods by a large margin. We only adopt a simple algorithm to generate tracking results. The tracking performance could be further improved with well-designed tracking algorithm.

4.5 Ablation Study

Key Components. We provide ablation experiments to validate the key components of PolarDETR. As shown in Tab. 4, Polar Param., context point and pixel ray respectively improve NDS by 1.5%, 0.9% and 0.6% and bring negligible

Table 2. 3D detection leaderboard of nuScenes benchmark (*test set*) (Mar. 4th, 2022). PolarDETR ranks 1st on the highly competitive leaderboard at the submission time.

Method	NDS↑	mAP↑	mATE↓	mASE↓	mAOE↓	mAVE↓	mAAE↓
DD3D	0.477	0.418	0.572	0.249	0.368	1.014	0.124
DETR3D	0.479	0.412	0.641	0.255	0.394	0.845	0.133
BEVDet	0.488	0.424	0.524	0.242	0.373	0.950	0.148
TPD-e	0.488	0.440	0.534	0.248	0.391	0.998	0.146
PolarDETR	0.493	0.431	0.588	0.253	0.408	0.845	0.129

Table 3. 3D tracking leaderboard of nuScenes benchmark (*test set*) (Mar. 4th, 2022).

Method	AMOTA↑	AMOTP↓	FP↓	FN↓	IDS↓	TID↓	LGD↓
DEFT	0.177	1.564	22163	60565	6901	1.600	3.080
QD-3DT	0.217	1.550	16495	60156	6856	1.620	2.961
PolarDETR	0.273	1.185	18853	59150	2170	0.990	2.330

overhead. It’s worth noting that Polar Param. only reformulates the optimization problem without introducing any computational budget. The significant improvement brought by Polar Param. validates its effectiveness.

Table 4. Ablations about the key components of PolarDETR.

-	Polar Param.	Context Point	Pixel Ray	NDS↑	FPS↑
A)				0.373	6.3
B)	✓			0.388	6.3
C)	✓	✓		0.397	6.0
D)	✓	✓	✓	0.403	6.0

Polar Velocity Decomposition. Tab. 5 presents the ablation study about the velocity decomposition. We compare cartesian and polar decomposition based on temporal input (frame $t - 1$ and t). Polar decomposition results in much lower mAVE (mean Average Velocity Error) and achieves much accurate velocity estimation.

Scaling Factor. Tab. 6 shows the ablations about the scaling factor k_{scaling} . Without numerically scaling up azimuth, *i.e.*, $k_{\text{scaling}} = 1$, the performance is poor because of the numerical unbalance between the tangential and radial direction. We use coarse grid search to tune k_{scaling} , and find $k_{\text{scaling}} = 20$ corresponds to relatively good results. With finer tuning, higher performance can be expected.

Table 5. Ablations about polar velocity decomposition.

Velocity	mAVE ↓
Cartesian (v_x, v_y)	0.556
Polar (v_{rad}, v_{tan})	0.539

Table 6. Ablations about the scaling factor.

k_{scaling}	NDS↑	mAP↑	mATE↓	mASE↓	mAOE↓	mAVE↓	mAAE↓
1	0.294	0.235	0.935	0.292	0.718	1.073	0.290
10	0.389	0.319	0.806	0.285	0.476	0.906	0.235
20	0.403	0.330	0.771	0.277	0.459	0.873	0.237
30	0.398	0.322	0.789	0.282	0.449	0.882	0.226

Context Point. Tab. 7 shows the ablations about the number of context points. Without context points, NDS is relatively low (0.390) because of lack of contextual information. Performance improves with the number increasing. But the gain get saturated with 4 context points (NDS = 0.403). And further increasing the number has negative effects. By default 4 context points are adopted in PolarDETR.

In Tab. 8, we further explore the impact of input for generating context points. As mentioned in Fig. 3 and Eq. (5), context points are predicted based on a combination of query embeddings q_i and center features $\{f_{c_i^1}, \dots, f_{c_i^K}\}$. The results prove that both terms contribute to better context feature aggregation and higher performance.

Table 7. Ablations about the number of context points.

Point	NDS↑	mAP↑	mATE↓	mASE↓	mAOE↓	mAVE↓	mAAE↓
0	0.390	0.324	0.797	0.279	0.546	0.883	0.215
2	0.402	0.329	0.787	0.281	0.466	0.869	0.223
4	0.403	0.330	0.771	0.277	0.459	0.873	0.237
8	0.399	0.326	0.787	0.279	0.460	0.880	0.236

Decoder Layer. Ablations about the number of decoder layers are shown in Tab. 9. With only 1 decoder layer (no iteration), results are poor. In PolarDETR, we adopt 6 decoder layers, which corresponds to relatively good results.

5 Qualitative Results

Center-Context Feature Aggregation. Visualizations about center and context points are in Fig. 6. Center points focus on object’s center for localization while context points focus on a wider region to capture more information. They complement each other and contribute to better feature aggregation.

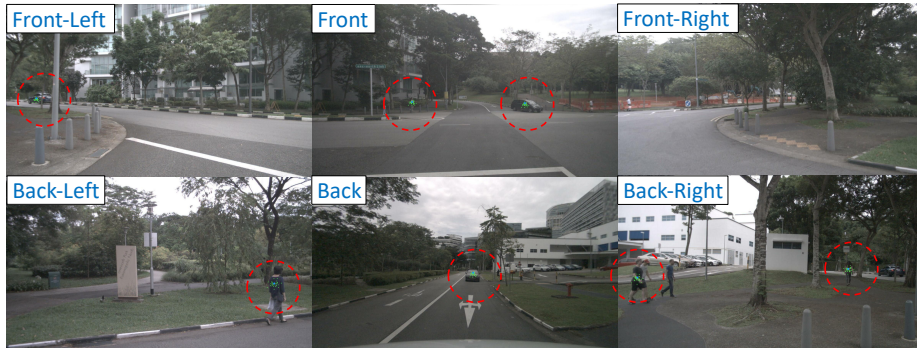
Table 8. Ablations about the input for generating context points.

Query Embeddings	Center Features	NDS \uparrow	mAP \uparrow
✓		0.393	0.324
	✓	0.403	0.325
✓	✓	0.403	0.330

**Fig. 5. Qualitative results about final predictions.** For visualization, 3D bounding box predictions are projected onto surround-view images. Blue ones denote predictions while green ones denote GTs.

Table 9. Ablations about the number of decoder layers.

Layer	NDS \uparrow	mAP \uparrow	mATE \downarrow	mASE \downarrow	mAOE \downarrow	mAVE \downarrow	mAAE \downarrow
1	0.256	0.233	0.901	0.318	1.167	1.182	0.389
2	0.358	0.316	0.792	0.286	0.662	1.044	0.262
3	0.383	0.328	0.790	0.286	0.532	0.957	0.243
4	0.395	0.323	0.786	0.277	0.461	0.890	0.248
5	0.401	0.329	0.769	0.277	0.472	0.884	0.240
6	0.403	0.330	0.771	0.277	0.459	0.873	0.237
7	0.398	0.325	0.780	0.280	0.440	0.904	0.235

**Fig. 6. Visualizations about center (blue) and context (green) points.** Several queries are selected for presentation (circled with red). Please zoom in for better view.

Detection Results. Qualitative results about final predictions are shown in Fig. 5. Blue and green boxes respectively denote predictions and GTs. Highly occluded or far-away objects may result in bad cases, which is a common problem faced by all detectors. Expect these challenging cases, PolarDETR achieves stable and satisfactory detection results.

6 Conclusion

In this paper, we present Polar Param. to exploit the view symmetry of surround-view camera system. Polar Param. establishes explicit associations between image patterns and prediction targets, superior to Image-based Param. and Cartesian Param.. Based on Polar Param., PolarDETR achieves promising performance in terms of both 3D detection and 3D tracking on the challenging nuScenes benchmark. And Polar Param. can be extended to other perception task and even planning task. We leave it as further work.

References

1. Caesar, H., Bankiti, V., Lang, A.H., Vora, S., Liong, V.E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuscenes: A multimodal dataset for autonomous driving. In: CVPR (2020) 4, 5, 8, 9, 11
2. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: ECCV (2020) 4, 5, 8
3. Contributors, M.: MMDetection3D: OpenMMLab next-generation platform for general 3D object detection. <https://github.com/open-mmlab/mmdetection3d> (2020) 10
4. Deng, J., Dong, W., Socher, R., Li, L., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR (2009) 10
5. Gao, P., Zheng, M., Wang, X., Dai, J., Li, H.: Fast convergence of DETR with spatially modulated co-attention. In: ICCV (2021) 4
6. Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: Vision meets robotics: The KITTI dataset. *Int. J. Robotics Res.* (2013) 5
7. Huang, J., Huang, G., Zhu, Z., Du, D.: Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. arXiv:2112.11790 (2021) 1, 2, 5, 11
8. Kumar, A., Brazil, G., Liu, X.: Groomed-nms: Grouped mathematically differentiable NMS for monocular 3d object detection. In: CVPR (2021) 5
9. Lee, Y., Hwang, J., Lee, S., Bae, Y., Park, J.: An energy and gpu-computation efficient backbone network for real-time object detection. In: CVPRW (2019) 11
10. Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV (2017) 9
11. Liu, Z., Wu, Z., Tóth, R.: SMOKE: single-stage monocular 3d object detection via keypoint estimation. In: CVPRW (2020) 2
12. Liu, Z., Zhou, D., Lu, F., Fang, J., Zhang, L.: Autoshape: Real-time shape-aware monocular 3d object detection. arXiv:2108.11127 (2021) 1, 5
13. Lu, Y., Ma, X., Yang, L., Zhang, T., Liu, Y., Chu, Q., Yan, J., Ouyang, W.: Geometry uncertainty projection network for monocular 3d object detection. arXiv:2107.13774 (2021) 1, 5
14. Meng, D., Chen, X., Fan, Z., Zeng, G., Li, H., Yuan, Y., Sun, L., Wang, J.: Conditional DETR for fast training convergence. ICCV (2021) 4
15. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: ECCV (2020) 7
16. Park, D., Ambrus, R., Guizilini, V., Li, J., Gaidon, A.: Is pseudo-lidar needed for monocular 3d object detection? arXiv:2108.06417 (2021) 1, 2, 5, 10, 11
17. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. In: NIPS (2019) 10
18. Reading, C., Harakeh, A., Chae, J., Waslander, S.L.: Categorical depth distribution network for monocular 3d object detection. In: CVPR (2021) 5
19. Roh, B., Shin, J., Shin, W., Kim, S.: Sparse DETR: efficient end-to-end object detection with learnable sparsity. In: ICLR (2022) 4
20. Stewart, R., Andriluka, M., Ng, A.Y.: End-to-end people detection in crowded scenes. In: CVPR (2016) 4, 8
21. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need (2017) 4

22. Wang, L., Du, L., Ye, X., Fu, Y., Guo, G., Xue, X., Feng, J., Zhang, L.: Depth-conditioned dynamic message propagation for monocular 3d object detection. In: CVPR (2021) 5
23. Wang, L., Zhang, L., Zhu, Y., Zhang, Z., He, T., Li, M., Xue, X.: Progressive coordinate transforms for monocular 3d object detection. arXiv:2108.05793 (2021) 5
24. Wang, T., Zhu, X., Pang, J., Lin, D.: FCOS3D: fully convolutional one-stage monocular 3d object detection. In: ICCVW (2021) 1, 2, 5, 10, 11
25. Wang, T., Zhu, X., Pang, J., Lin, D.: Probabilistic and geometric depth: Detecting objects in perspective. In: CoRL (2021) 1, 2, 5, 11
26. Wang, T., Yuan, L., Chen, Y., Feng, J., Yan, S.: Pnp-detr: Towards efficient visual analysis with transformers. In: ICCV (2021) 4
27. Wang, T., Zhang, J., Cai, Y., Yan, S., Feng, J.: Direct multi-view multi-person 3d pose estimation. In: NeurIPS (2021) 7
28. Wang, Y., Guizilini, V.C., Zhang, T., Wang, Y., Zhao, H., Solomon, J.: DETR3d: 3d object detection from multi-view images via 3d-to-2d queries. In: CoRL (2021) 1, 2, 4, 5, 11
29. Xie, E., Sun, P., Song, X., Wang, W., Liu, X., Liang, D., Shen, C., Luo, P.: Polarmask: Single shot instance segmentation with polar representation. In: CVPR (2020) 4
30. Yin, T., Zhou, X., Krähenbühl, P.: Center-based 3d object detection and tracking. In: CVPR (2021) 10
31. Zhang, Y., Zhou, Z., David, P., Yue, X., Xi, Z., Gong, B., Foroosh, H.: Polarnet: An improved grid representation for online lidar point clouds semantic segmentation. In: CVPR (2020) 4
32. Zhang, Y., Lu, J., Zhou, J.: Objects are different: Flexible monocular 3d object detection. In: CVPR (2021) 1, 5
33. Zhou, Y., Barnes, C., Lu, J., Yang, J., Li, H.: On the continuity of rotation representations in neural networks. In: CVPR (2019) 6
34. Zhou, Y., He, Y., Zhu, H., Wang, C., Li, H., Jiang, Q.: Monocular 3d object detection: An extrinsic parameter free approach. In: CVPR (2021) 5
35. Zhu, X., Zhou, H., Wang, T., Hong, F., Ma, Y., Li, W., Li, H., Lin, D.: Cylindrical and asymmetrical 3d convolution networks for lidar segmentation. arXiv:2011.10033 (2020) 4
36. Zhu, X., Hu, H., Lin, S., Dai, J.: Deformable convnets V2: more deformable, better results. In: CVPR (2019) 4, 10
37. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable DETR: deformable transformers for end-to-end object detection. In: ICLR (2021) 4, 5