

DFA3D: 3D Deformable Attention For 2D-to-3D Feature Lifting

Hongyang Li^{1,3,*†} Hao Zhang^{2,3,*†} Zhaoyang Zeng³ Shilong Liu^{4,3}
 Feng Li^{2,3} Tianhe Ren³ Lei Zhang^{1,3,‡}

¹South China University of Technology.

²The Hong Kong University of Science and Technology.

³International Digital Economy Academy (IDEA).

⁴Dept. of CST., BNRist Center, Institute for AI, Tsinghua University.

Abstract

In this paper, we propose a new operator, called 3D DeFormable Attention (DFA3D), for 2D-to-3D feature lifting, which transforms multi-view 2D image features into a unified 3D space for 3D object detection. Existing feature lifting approaches, such as Lift-Splat-based and 2D attention-based, either use estimated depth to get pseudo LiDAR features and then splat them to a 3D space, which is a one-pass operation without feature refinement, or ignore depth and lift features by 2D attention mechanisms, which achieve finer semantics while suffering from a depth ambiguity problem. In contrast, our DFA3D-based method first leverages the estimated depth to expand each view’s 2D feature map to 3D and then utilizes DFA3D to aggregate features from the expanded 3D feature maps. With the help of DFA3D, the depth ambiguity problem can be effectively alleviated from the root, and the lifted features can be progressively refined layer by layer, thanks to the Transformer-like architecture. In addition, we propose a mathematically equivalent implementation of DFA3D which can significantly improve its memory efficiency and computational speed. We integrate DFA3D into several methods that use 2D attention-based feature lifting with only a few modifications in code and evaluate on the nuScenes dataset. The experiment results show a consistent improvement of +1.41% mAP on average, and up to +15.1% mAP improvement when high-quality depth information is available, demonstrating the superiority, applicability, and huge potential of DFA3D. The code is available at <https://github.com/IDEA-Research/3D-deformable-attention.git>.

1. Introduction

3D object detection is a fundamental task in many real-world applications such as robotics and autonomous driving. Although LiDAR-based methods [37, 40, 23, 33] have achieved impressive results with the help of accurate 3D perception from LiDAR, multi-view camera-based methods have recently received extensive attention because of their low cost for deployment and distinctive capability of long-range and color perception. Classical multi-view camera-based 3D object detection approaches mostly follow monocular frameworks which first perform 2D/3D object detection in each individual view and then conduct cross-view post-processing to obtain the final result. While remarkable progress has been made [30, 31, 24, 26, 1], such a framework cannot fully utilize cross-view information and usually leads to low performance.

To eliminate the ineffective cross-view post-processing, several end-to-end approaches [25, 13, 10, 14, 32, 20, 21] have been developed. These approaches normally contain three important modules: a backbone for 2D image feature extraction, a feature lifting module for transforming multi-view 2D image features into a unified 3D space (e.g. BEV space in an ego coordinate system) to obtain lifted features, and a detection head for performing object detection by taking as input the lifted features. Among these modules, the feature lifting module serves as an important component to bridge the 2D backbone and the 3D detection head, whose quality will greatly affect the final detection performance.

To perform feature lifting, recent methods usually pre-define a set of 3D anchors in the ego coordinate system sparsely or uniformly, with randomly initialized content features. After that, they lift 2D image features into the 3D anchors to obtain the lifted features. Some methods [13, 10, 9, 34, 12] utilize a straightforward lift and splat mechanism [25] by first lifting 2D image features into pseudo LiDAR features in an ego coordinate system using estimated depth and then assigning the pseudo LiDAR fea-

*This work was done during the internship at IDEA.

†Equal contribution.

‡Corresponding author.

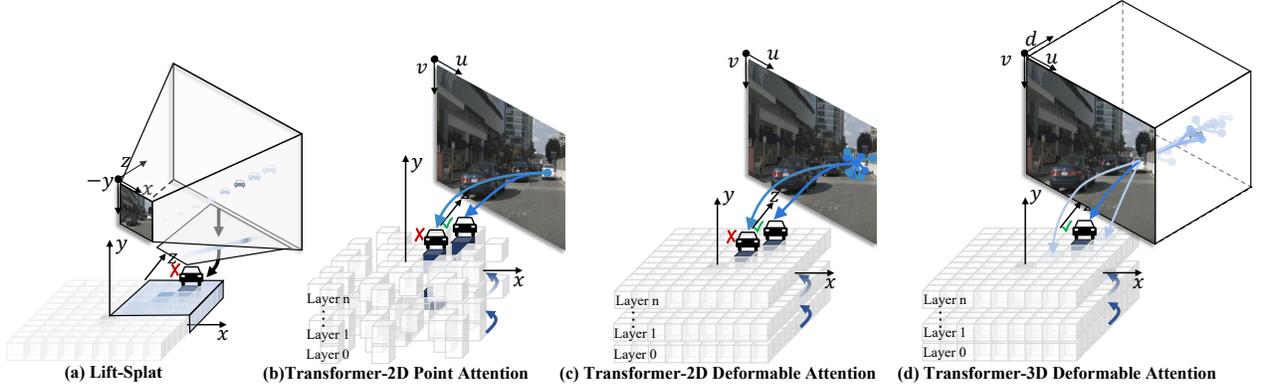


Figure 1. Comparisons of feature lifting methods. (a) Lift-Splat-based methods cannot handle depth errors after feature lifting, which can result in wrong predictions. Point attention-based (b) and 2D deformable attention-based (c) can refine the lifted features layer-by-layer. However, they suffer from depth ambiguity when multiple 3D objects are projected to the same 2D point, which may result in duplicate predictions along a ray connecting the ego car and a target object. (d) Our 3D deformable attention can effectively alleviate the depth ambiguity problem as we sample features in 3D spaces. Moreover, the multi-layer design of Transformers helps refine features and sampling points layer by layer.

tures to their closest 3D anchors to obtain the final lifted features, as shown in Fig. 1(a). Although such methods have achieved remarkable performance, due to their excessive resource consumption when lifting 2D image features, they normally cannot utilize multi-scale feature maps, which are crucial for detecting small (faraway) objects.

Instead of lift and splat, some other works utilize 2D attention to perform feature lifting. Such methods treat predefined 3D anchors as 3D queries and propose several ways to aggregate 2D image features (keys) to 3D anchors. Typical 2D attention mechanisms that have been explored include dense attention [20], point attention (a degenerated deformable attention) [3, 4], and deformable attention [14, 36, 5]. 2D attention-based methods can naturally work with Transformer-like architectures to progressively refine the lifted features layer by layer. Moreover, deformable attention makes it possible to use multi-scale 2D features, thanks to its introduced sparse attention computation. However, the major weakness of 2D attention-based approaches is that they suffer from a crucial depth ambiguity problem, as shown in Fig. 1(b,c). That is, due to the ignorance of depth information, when projecting 3D queries to a 2D view, many 3D queries will end up with the same 2D position with similar sampling points in the 2D view. This will result in highly entangled aggregated features and lead to wrong predictions along a ray connecting the ego car and a target object (as visualized in Fig. 7). Although some concurrent works [15, 39] have tried to alleviate this issue, they can not fix the problem from the root. The root cause lies in the lack of depth information when applying 2D attention to sample features in a 2D pixel coordinate system¹.

To address the above problems, in this paper, we propose a basic operator called 3D deformable attention (DFA3D),

built upon which we develop a novel feature lifting approach, as shown in Fig. 1(d). We follow [25] to leverage a depth estimation module to estimate a depth distribution for each 2D image feature. We expand the dimension of each single-view 2D image feature map by computing the outer product of them and their estimated depth distributions to obtain the expanded 3D feature maps in a 3D pixel coordinate system². Each 3D query in the BEV space is projected to the 3D pixel space with a set of predicted 3D sampling offsets to specify its 3D receptive field and pool features from the expanded 3D feature maps. In this way, a 3D query close to the ego car only pools image features with smaller depth values, whereas a 3D query far away from the ego car mainly pools image features with larger depth values, and thus the depth ambiguity problem is effectively alleviated.

In our implementation, as each of the expanded 3D feature maps is an outer product between a 2D image feature map and its estimated depth distributions, we do not need to maintain a 3D tensor in memory to avoid excessive memory consumption. Instead, we seek assistance from math to simplify DFA3D into a mathematically equivalent depth-weighted 2D deformable attention and implement it through CUDA, making it both memory-efficient and fast. The Pytorch interface of DFA3D is very similar to 2D deformable attention and only requires a few modifications to replace 2D deformable attention, making our feature lifting approach easily portable.

In summary, our main contributions are:

1. We propose a basic operator called 3D deformable attention (DFA3D) for feature lifting. Leveraging the property of outer product between 2D features and their estimated depth, we develop a memory-efficient and fast implementation.

¹The pixel coordinate system that only considers the u, v axes.

²The pixel coordinate system that considers all of the u, v, d axes.

- Based on DFA3D, we develop a novel feature lifting approach, which not only alleviates the depth ambiguity problem from the root, but also benefits from multi-layer feature refinement of a Transformer-like architecture. Thanks to the simple interface of DFA3D, our DFA3D-based feature lifting can be implemented in any method that utilizes 2D deformable attention (also its degeneration)-based feature lifting with only a few modifications in code.
- The consistent performance improvement (+1.41 mAP on average) in comparative experiments on the nuScenes [2] dataset demonstrate the superiority and generalization ability of DFA3D-based feature lifting.

2. Related Work

Multi-view 3D Detectors with Lift-Splat-based Feature Lifting:

To perform 3D tasks in the multi-view context in an end-to-end pipeline, Lift-Splat [25] proposes to first lift multi-view image features by performing outer product between 2D image feature maps and their estimated depth, and then transform them to a unified 3D space through camera parameters to obtain the lifted pseudo LiDAR features. Following the LiDAR-based methods, the pseudo LiDAR features are then splatted to predefined 3D anchors through a voxel-pooling operation to obtain the lifted features. The lifted features are then utilized to facilitate the following downstream tasks. For example, BEVDet [10] proposes to adopt such a pipeline for 3D detection and verifies its feasibility. Although being effective, such a method suffers from the problem of huge memory consumption, which prevents them from using multi-scale feature maps. Moreover, due to the fixed assignment rule between the pseudo LiDAR features and the 3D anchors, each 3D anchor can only interact with image features once without any further adjustment later. This problem makes Lift-Splat-based methods reliant on the quality of depth. The recently proposed BEVDepth [13] makes a full analysis and validates that the learned depth quality has a great impact on performance. As the implicitly learned depth in previous works does not satisfy the accuracy requirement, BEVDepth proposes to resort to monocular depth estimation for help. By explicitly supervising monocular depth estimation with ground truth depth obtained by projecting LiDAR points onto multi-view images, BEVDepth achieves a remarkable performance. Inspired by BEVDepth, instead of focusing on the above-discussed problems of refinement or huge memory consumption, many works [12, 34] try to improve their performance through improving depth quality by resorting to multi-view stereo or temporal information. However, due to small overlap regions between multi-view images and also the movements of surrounding objects, they still rely on the monocular depth estimation, which is an ill-posed problem and can not be accurately addressed.

Multi-view 3D Detectors with 2D Attention-based Feature Lifting:

Motivated by the progress in 2D detection [27, 28, 18, 19, 38, 11, 17], many works propose to introduce attention into camera-based 3D detection. They treat multi-view image features as keys and values and the predefined 3D anchors as 3D queries in the unified 3D space. The 3D queries are projected onto multi-view 2D image feature maps according to the extrinsic and intrinsic camera parameters and lift features from the image feature maps through cross attention like the decoder part of 2D detection methods [3, 41]. PETR [20] and PETRv2 [21] are based on the classical (dense) attention mechanism, where the interactions between each 3D query and all image features (keys) make it inefficient. DETR3D-like methods [32, 4] proposes to utilize the point attention which lets a query interact with only one key obtained by bilinear interpolation according to its projected location (also called point deformable attention [14]). Although being efficient, the point attention has a small receptive field in multi-view 2D image feature maps and normally results in a relatively weak performance. By contrast, BEVFormer [14, 36] proposes to apply 2D deformable attention to let each 3D query interact with a local region around the location where the 3D query is projected to and obtains a remarkable result. OA-BEV [5] tries to introduce 2D detection into the pipeline to guide the network to further focus on the target objects. Benefiting from the Transformer-like architecture, the lifted features obtained by 2D attention can be progressively refined layer-by-layer. However, the methods with 2D attention-based feature lifting suffer from the problem of depth ambiguity: when projecting 3D queries onto a camera view, those with the same projected coordinates (u, v) but different depth values d end up with the same reference point and similar sampling points in the view. Consequently, these 3D queries will aggregate features that are strongly entangled and thus result in duplicate predictions along a ray connecting the ego car and a target object in BEV. Some concurrent works [39, 15] have noticed this problem, but none of them tackle the problem from the root.

Instead of conducting feature lifting using existing operators, such as voxel-pooling and 2D deformable attention, we develop a basic operator DFA3D to perform feature lifting. As summarized in Table 1, our proposed DFA3D-based feature lifting can not only benefit from multi-layer refinement but also address the depth ambiguity problem. Moreover, DFA3D-based feature lifting can take multi-scale feature maps into consideration efficiently.

Table 1. Comparison with other feature lifting methods. † indicates the direct utilization of multi-scale feature maps without sampling.

Methods	depth differentiation	multi-layer refinement	multi-scale†
Lift-Splat-based	✓		
2D Attention-based		✓	✓
DFA3D-based (ours)	✓	✓	✓

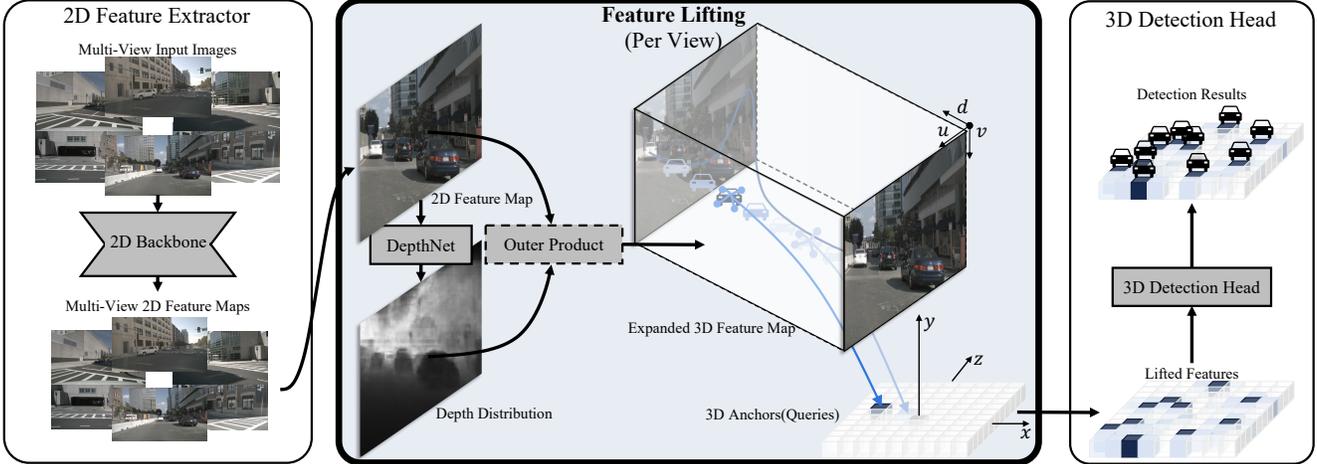


Figure 2. The overview of an end-to-end multi-view 3D detection pipeline with DFA3D. The 2D feature extractor takes multi-view images as inputs and extracts multi-scale 2D image feature maps for each view independently. After that, the 2D image feature maps are fed into DFA3D-based feature lifting process and lifted to predefined 3D anchors, which are considered as 3D queries during the process. Finally, a 3D detection head takes the lifted features as input and predicts 3D bounding boxes as the final detection results. The “Outer Product” enclosed by a dashed border indicates that we do not actually conduct it. For better visualization, we assume the feature maps have one scale. We present 2D image feature maps by RGB images and depth distributions by a gray-scale depth map. Best viewed in color.

3. Approach

The overview of an end-to-end multi-view 3D detection approach with DFA3D-based feature lifting is shown in Fig. 2. In this section, we introduce how to construct the expanded 3D feature maps for DFA3D in Sec. 3.1 and how DFA3D works to realize feature lifting in Sec. 3.2. In addition, we explain our efficient implementation of DFA3D in Sec. 3.3. Finally, we provide a detailed analysis and comparison between our DFA3D-based feature lifting and the current mainstream feature lifting methods in Sec. 3.4.

3.1. Feature Expanding

In the context of multi-view 3D detection, the depth information of input images is unavailable. We follow [13, 34, 12, 10] to introduce a monocular depth estimation to supplement depth information for 2D image features. Specifically, we follow [13] to adopt a DepthNet [7] module, which takes as input the 2D image feature maps from the backbone and generates a discrete depth distribution for each 2D image feature. The DepthNet module can be trained by using LiDAR information [13] or ground truth 3D bounding box information [5] as supervision.

We denote the multi-view 2D image feature maps as $\mathbf{X} \in \mathbb{R}^{V \times H \times W \times C}$, where V , H , W , and C indicates the number of views, spatial height, spatial width, and the number of channels, respectively. We feed the multi-view 2D image feature maps into DepthNet to obtain the discrete depth distributions. The distributions have the same spatial shape as \mathbf{X} and can be denoted as $\mathbf{D} \in \mathbb{R}^{V \times H \times W \times D}$, where D indicates the number of pre-defined discretized

depth bins. Here, for simplicity, we assume that \mathbf{X} only has one feature scale. In the context of multi-scale, in order to maintain the consistency of depths across multiple scales, we select the feature maps of one scale to participate in the generation of depth distributions. The generated depth distributions are then interpolated to obtain the depth distributions for multi-scale feature maps.

After obtaining the discrete depth distributions \mathbf{D} , we expand the dimension of \mathbf{X} into 3D by conducting the outer product between \mathbf{D} and \mathbf{X} , which is formulated as $\mathbf{F} = \mathbf{D} \otimes_{-1} \mathbf{X}$, where \otimes_{-1} indicates the outer product conducted at the last dimension, and $\mathbf{F} \in \mathbb{R}^{V \times H \times W \times D \times C}$ denotes the multi-view expanded 3D feature maps. For more details about the expanded 3D feature maps, please refer to the further visualization in appendix.

Directly expanding 2D image feature maps through outer products will lead to high memory consumption, especially when taking multi-scale into consideration. To address this issue, we develop a memory-efficient algorithm to solve this issue, which will be explained in Sec. 3.3.

3.2. 3D Deformable Attention and Feature Lifting

After obtaining the multi-view expanded 3D feature maps, DFA3D is utilized as the backend to transform them into a unified 3D space to obtain the lifted features. Specifically, in the context of feature lifting, DFA3D treats the predefined 3D anchors as 3D queries, the expanded 3D feature maps as 3D keys and values, and performs deformable attention in the 3D pixel coordinate system.

To be more specific, for a 3D query located at (x, y, z) in the ego coordinate system, its view reference point $\mathbf{R}_n =$

(u_n, v_n, d_n) in the 3D pixel coordinate system for the n -th view is calculated by³

$$d_n \begin{bmatrix} u_n \\ v_n \\ 1 \end{bmatrix} = \begin{bmatrix} f_n^x & 0 & u_n^0 \\ 0 & f_n^y & v_n^0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix}, \quad (1)$$

where $f_n^x, f_n^y, u_n^0, v_n^0$ are the corresponding camera intrinsic parameters. By solving Eq. 1, we can obtain

$$u_n = f_n^x \frac{x}{d_n} + u_n^0, v_n = f_n^y \frac{y}{d_n} + v_n^0, d_n = z. \quad (2)$$

We denote the content feature of a 3D query as $\mathbf{q} \in \mathbb{R}^{C_q}$, where C_q indicates the dimension of \mathbf{q} . The 3D deformable attention mechanism in the n -th view can be formulated by⁴

$$\begin{aligned} \Delta \mathbf{S} &= \mathbf{W}^S \mathbf{q}, \Delta \mathbf{S}_k = \Delta \mathbf{S}_{[3k:3(k+1)]}, \mathbf{A} = \mathbf{W}^A \mathbf{q}, \\ \mathbf{q}_n &= \sum_{k=1}^K \mathbf{A}_k \text{Trili}(\mathbf{F}_n, \mathbf{R}_n + \Delta \mathbf{S}_k), \end{aligned} \quad (3)$$

where \mathbf{q}_n is the querying result in the n -th view, K denotes the number of sampling points, $\mathbf{W}^S \in \mathbb{R}^{3K \times C_q}$ and $\mathbf{W}^A \in \mathbb{R}^{K \times C_q}$ are learnable parameters, $\Delta \mathbf{S}$ and \mathbf{A} represent the view sampling offsets in the 3D pixel coordinate system and their corresponding attention weights respectively. $\text{Trili}(\cdot)$ indicates the trilinear interpolation used to sample features in the expanded 3D feature maps. The detailed implementation of $\text{Trili}(\cdot)$ and our simplification will be explained in Sec. 3.3.

After obtaining the querying results through DFA3D in multiple views, the final lifted feature \mathbf{q}' for the 3D query is obtained by query result aggregation,

$$\mathbf{q}' = \sum_n \mathbf{V}_n \mathbf{q}_n, \quad (4)$$

where \mathbf{V}_n indicates the visibility of the 3D query in the n -th view. The lifted feature is used to update the content feature of the 3D query for feature refinement.

3.3. Efficient 3D Deformable Attention

Although the 3D deformable attention has been formulated in Eq. 3, maintaining the multi-view expanded 3D feature maps $\mathbf{F} \in \mathbb{R}^{V \times H \times W \times D \times C}$ will lead to tremendous memory consumption, especially when taking multi-scale features into consideration. We prove that in the context of feature expanding through the outer product between 2D image feature maps and their discrete depth distributions,

³Since there is only a rigid transformation between the camera coordinate system of n -th view and the ego coordinate system, we assume the extrinsic matrix as identity matrix and drop it for notation simplicity.

⁴Here we set the number of attention heads as 1 and the number of feature scales as 1, and drop their indices for notation simplicity.

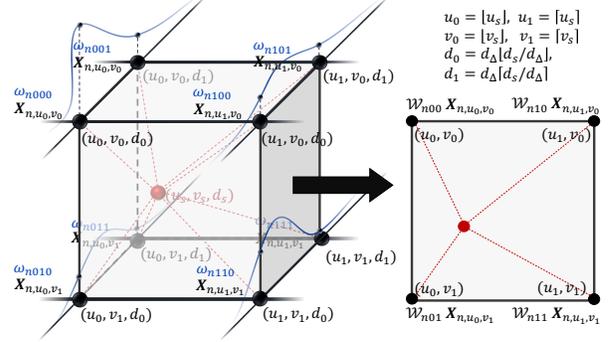


Figure 3. The original trilinear interpolation shown on the left can be simplified to a weighted bilinear interpolation as shown right.

the trilinear interpolation can be transformed into a depth-weighted bilinear interpolation. Hence the 3D deformable attention can also be transformed into a depth-weighted 2D deformable attention accordingly.

As shown in Fig. 3, without loss of generality, for a specific view (the n -th view for example) with an expanded 3D image feature map $\mathbf{F}_n = \mathbf{D}_n \otimes_{-1} \mathbf{X}_n$, $\mathbf{F}_n \in \mathbb{R}^{H \times W \times D \times C}$, assume a sampling point $\mathbf{s} = (u_s, v_s, d_s)$ in the 3D pixel coordinate system falls in a cuboid formed by eight points: $\{(u_i, v_j, d_k)\}$, $i, j, k \in \{0, 1\}$, which are the sampling point's eight nearest points in \mathbf{F}_n . The sampling point performs trilinear interpolation on the features of these eight points to obtain the sampling result. Since the eight points with the same u, v coordinates share the same 2D image feature, their corresponding contributions in the interpolation process can be presented as $\{\omega_{nij k} \mathbf{X}_{n, u_i, v_j}\}$, where $\omega_{nij k}$ denotes $\mathbf{D}_{n, u_i, v_j, d_k}$ for notation simplicity. $\omega_{nij k}$ indicates the depth confidence scores of point (u_i, v_j, d_k) selected from the discrete depth distributions \mathbf{D}_{n, u_i, v_j} that it belongs to, according to its depth values d_k . Therefore, the trilinear interpolation can be formulated as:

$$\begin{aligned} t_u &= \frac{u_s - \lfloor u_s \rfloor}{\lceil u_s \rceil - \lfloor u_s \rfloor}, t_v = \frac{v_s - \lfloor v_s \rfloor}{\lceil v_s \rceil - \lfloor v_s \rfloor}, \\ t_d &= \frac{d_s - \lfloor d_s/d_\Delta \rfloor d_\Delta}{\lceil d_s/d_\Delta \rceil d_\Delta - \lfloor d_s/d_\Delta \rfloor d_\Delta}, \\ \text{Trili}(\mathbf{F}_n, \mathbf{s}) &= \\ &\omega_{n000} (1 - t_d) (1 - t_u) (1 - t_v) \mathbf{X}_{n, u_0, v_0} + \\ &\omega_{n001} t_d (1 - t_u) (1 - t_v) \mathbf{X}_{n, u_0, v_0} + \\ &\omega_{n100} (1 - t_d) t_u (1 - t_v) \mathbf{X}_{n, u_1, v_0} + \\ &\omega_{n101} t_d t_u (1 - t_v) \mathbf{X}_{n, u_1, v_0} + \\ &\omega_{n010} (1 - t_d) (1 - t_u) t_v \mathbf{X}_{n, u_0, v_1} + \\ &\omega_{n001} t_d (1 - t_u) t_v \mathbf{X}_{n, u_0, v_1} + \\ &\omega_{n110} (1 - t_d) t_u t_v \mathbf{X}_{n, u_1, v_1} + \\ &\omega_{n111} t_d t_u t_v \mathbf{X}_{n, u_1, v_1}, \end{aligned} \quad (5)$$

where d_Δ indicates the depth interval between two adjacent

depth bins. Eq. 5 can be further re-formulated into

$$\begin{aligned}
 \mathcal{W}_{nij} &= \omega_{nij0}(1 - t_d) + \omega_{nij1}t_d, \\
 \text{Trili}(\mathbf{F}_n, \mathbf{s}) &= \\
 \mathcal{W}_{n00} (1 - t_u) (1 - t_v) \mathbf{X}_{n,u_0,v_0} &+ \\
 \mathcal{W}_{n10} t_u (1 - t_v) \mathbf{X}_{n,u_1,v_0} &+ \\
 \mathcal{W}_{n01} (1 - t_u) t_v \mathbf{X}_{n,u_0,v_1} &+ \\
 \mathcal{W}_{n11} t_u t_v \mathbf{X}_{n,u_1,v_1}. &
 \end{aligned} \tag{6}$$

Comparing Eq. 5 and Eq. 6, we can find that the trilinear interpolation here can be actually split into two parts: 1) a simple linear interpolation of the estimated discrete depth distribution along the depth axis to obtain the depth scores $\mathcal{W}_n = \{\mathcal{W}_{nij}\}$, and 2) a depth score-weighted bilinear interpolation in the 2D image feature maps. Compared with bilinear interpolation, the only extra computation we need is to sample depth scores \mathcal{W}_n through linear interpolation. With such an optimization, we can calculate the 3D deformable attention on the fly rather than maintaining all the expanded features. Thus the entire attention mechanism can be implemented in a more efficient way. Moreover, such an optimization can reduce half of the multiplication operations and thus speedup the overall computation as well.

We conduct a simple experiment to compare the efficiency. The results are shown in Table 2. For a fair comparison, we show the resource consumption of all steps in detail. The results show that our efficient implementation only takes about 3% time cost and 1% memory consumption compared with the vanilla one.

Table 2. The comparisons of resource consumption (Memory / Time) between the vanilla one and our efficient one. ‘‘Expand’’, ‘‘Aggregate’’ and ‘‘Full’’ indicate feature expanding, feature aggregation, and overall process respectively. Our efficient one does not need ‘‘Expand’’.

Method	Expand	Aggregate	Full
Vanilla	25303MB / 76.5ms	674MB / 85.4ms	3204MB / 161.9ms
Efficient	- / -	29MB / 5.2ms	29MB / 5.2ms

3.4. Analysis

Comparison with Lift-Splat-based Feature Lifting. Lift-Splat-based feature lifting follows the pipeline of LiDAR-based methods. To construct pseudo LiDAR features as if they are obtained from LiDAR data, it also utilizes depth estimation to supplement depth information. It obtains the pseudo LiDAR features by constructing single-scale 3D expanded feature maps explicitly and transforming them from a 3D pixel coordinate system into the ego coordinate system according to the camera parameters. The pseudo LiDAR features are further assigned to their nearest 3D anchors to generate lifted features for downstream tasks. Since the locations of both 3D anchors and pseudo LiDAR features are

constant, the assignment rule between them is fixed based on their geometrical relationship.

In DFA3D, the relationship between 3D queries and the expanded 3D features that are computed based on the estimated sampling locations can also be considered as an assignment rule. Instead of being fixed, DFA3D can progressively refine the assignment rule by updating the sampling locations layer by layer in a Transformer-like architecture. Besides, the efficient implementation of DFA3D enables the utilization of multi-scale 3D expanded features, which are more crucial to detecting small (faraway) objects.

Comparison with 2D Attention-based Feature Lifting. Similar to DFA3D, 2D attention-based feature lifting also transforms each 3D query into a pixel coordinate system. However, to satisfy the input of off-the-shelf 2D attention operators, a 3D query needs to be projected to a 2D one, where its depth information is discarded. It is a compromise to implementation and can cause multiple 3D queries to collapse into one 2D query. Hence multiple 3D queries can be entangled greatly. This can be proved easily according to the pin-hole camera model. In such a scenario, according to Eq. 2, as long as a 3D query’s coordinate (x, y, z) in the ego coordinate system satisfies

$$\frac{x}{z} = \frac{u_n - u_n^0}{f_n^x}, \frac{y}{z} = \frac{v_n - v_n^0}{f_n^y}, \tag{7}$$

its reference point in the target 2D pixel coordinate system will be projected to the shared (u_n, v_n) coordinate, resulting in poorly differentiated lifted features that will be used to update the 3D queries’ content and downstream tasks.

4. Experiments

We conduct extensive experiments to evaluate the effectiveness of DFA3D-based feature lifting compared with other feature lifting methods. We integrate DFA3D into several open-source methods to verify its generalization ability and portability.

4.1. Experiment Setup

Dataset and Metrics We follow previous works [14, 13, 32] to conduct experiments on the nuScenes dataset [2]. The nuScenes dataset contains 1,000 sequences, which are captured by a variety of sensors (e.g. cameras, RADAR, LiDAR). Each sequence is about 20 seconds, and is annotated in 2 frames/second. Since we focus on the camera-based 3D detection task, we use only the image data taken from 6 cameras. There are 40k samples with 1.4M annotated 3D bounding boxes in total. We use the classical splits of training, validation, and testing set with 28k, 6k, and 6k samples, respectively. We follow previous works [14, 13, 32] to take the official evaluation metrics, including mean average precision (mAP), ATE, ASE, AOE, AVE, and AAE

for our evaluations. These six metrics evaluate results from center distance, translation, scale, orientation, velocity, and attribute, respectively. In addition, a comprehensive metric called nuScenes detection score (NDS) is also provided.

Implementation Details For fair comparisons, when implementing DFA3D on the open-source methods, we keep all configurations the same without bells and whistles. For all experiments, without specification, we take BEVFormer-base [14] as our baseline, which uses a ResNet101-DCN [8, 6] backbone initialized from FCOS3D [31] checkpoint and a Transformer with 6 encoder layers and 6 decoder layers. All experiments are trained with 24 epochs using AdamW optimizer [22] with a base learning rate of 2×10^{-4} on 8 NVIDIA Tesla A100 GPUs.

4.2. Comparisons of Feature Lifting Methods

To compare different feature lifting methods fairly, we set feature lifting as the only variable and keep the others the same. More specifically, we block the temporal information and utilize the same 2D backbone and 3D detection head. The results are shown in Table 3.

We first compare four different feature lifting methods in the first four rows in Table 3. All models use one layer for a fair comparison. The results show that DFA3D outperforms all previous works, which demonstrates the effectiveness. Although the larger receptive field makes 2D deformable attention (DFA2D)-based method obtains much better result than the point attention-based method, its performance is not satisfactory enough compared with the Lift-Splat-based one. One key difference between DFA2D and Lift-Splat is their depth utilization. The extra depth information enables the Lift-Splat-based method to achieve a better performance, although the depth information is predicted and not accurate enough. Different from the DFA2D-based feature lifting, DFA3D enables the usage of depth information in the deformable attention mechanism and helps our DFA3D-based feature lifting achieve the best performance.

As discussed in Sec. 3.4, the adjustable assignment rule enables attention-based methods to conduct multi-layer refinement. With the help of one more layer refinement, DFA2D-based feature lifting surpasses Lift-Splat. However, benefitting from the better depth utilization, our DFA3D-based feature lifting still maintains the superiority.

4.3. Generalizations for Different Models

To verify the generalization and portability of DFA3D-based feature lifting, we evaluate it on various open-sourced methods that rely on 2D deformable attention or point attention-based feature lifting. Thanks to the mathematical simplification of DFA3D, our DFA3D-based feature lifting can be easily integrated into these methods with only a few code modifications (please refer to the appendix for a more intuitive comparison).

Table 3. Comparisons of feature lifting methods. DFA2D denotes the 2D deformable attention. We use the same architecture and supervision for depth estimation as our DFA3D for the Lift-Splat-based method. We also leverage SECOND FPN [35] to enable multi-scale feature maps for the Lift-Splat-based method.

Method	#layers	mAP \uparrow	mATE \downarrow	mASE \downarrow	mAOE \downarrow	mAVE \downarrow	mAAE \downarrow	NDS \uparrow
PointAttn	1	35.0	76.3	27.8	42.9	84.9	20.7	42.2
DFA2D	1	35.9	74.6	27.8	42.5	85.8	20.9	42.8
Lift-Splat	1	36.9	73.1	28.0	44.7	85.0	23.6	43.0
DFA3D (Ours)	1	37.3	73.4	27.5	41.7	85.2	22.5	43.6
DFA2D	2	37.1	72.7	27.7	43.5	81.4	20.6	44.0
DFA3D (Ours)	2	37.9	72.4	27.7	38.1	78.1	22.3	45.1

Table 4. Comparisons of our method and baselines. Without specifications, all models use ResNet101 [8] as the backbone. \dagger denotes that the model adopts ResNet101-DCN [6] as the backbone. Sparse4D* denotes the Sparse4D without Depth Reweight Module (DRM) proposed by [15].

Method	mAP \uparrow	mATE \downarrow	mASE \downarrow	mAOE \downarrow	mAVE \downarrow	mAAE \downarrow	NDS \uparrow
DETR3D \dagger [32]	34.7	76.5	26.8	39.2	87.6	21.1	42.2
DETR3D-DFA3D \dagger	35.5(+0.8)	74.4	26.8	41.6	86.4	20.7	42.8(+0.6)
BEVFormer-t [14]	25.2	90.0	29.4	65.5	65.7	21.6	35.4
BEVFormer-t-DFA3D	26.9(+1.7)	88.0	29.2	60.6	60.6	23.1	37.3(+1.9)
BEVFormer-s [14]	37.0	72.1	28.0	40.7	43.6	22.0	47.9
BEVFormer-s-DFA3D	40.1(+3.1)	72.1	27.9	41.1	39.1	19.6	50.1(+2.2)
BEVFormer-b [14]	41.6	67.3	27.4	37.2	39.4	19.8	51.7
BEVFormer-b-DFA3D \dagger	43.0(+1.4)	65.4	27.1	37.4	34.1	20.5	53.1(+1.4)
DA-BEV-S \dagger [39]	42.8	63.3	27.3	33.1	32.7	18.8	53.9
DA-BEV-S-DFA3D \dagger	43.5(+0.7)	64.5	27.0	32.8	31.6	20.2	54.2(+0.3)
DA-BEV \dagger [39]	43.3	62.3	27.1	35.1	30.9	18.8	54.5
DA-BEV-DFA3D \dagger	44.1(+0.8)	62.6	27.4	33.4	31.1	19.4	54.7(+0.2)
Sparse4D* \dagger [39]	43.2	-	-	37.9	-	-	53.7
Sparse4D*-DRM \dagger [39]	43.6(+0.4)	63.3	27.9	36.3	31.7	17.7	54.1(+0.4)
Sparse4D*-DFA3D \dagger	44.6(+1.4)	63.2	27.1	38.8	30.5	18.1	54.5(+0.8)

We compare baselines to those that integrate DFA3D-based feature lifting in Table 4. The results show that DFA3D brings consistent improvements in different methods, indicating its generalization ability across different models. Furthermore, DFA3D introduces significant gains of 1.7, 3.1, and 1.4 mAP on BEVFormer-t⁵, BEVFormer-s, and BEVFormer-b, respectively, showing the effectiveness of DFA3D across different model sizes. Note that, for a fair comparison, we use one sampling point and fixed sampling offsets when conducting experiments on DETR3D [32].

For a more comprehensive comparison, we also conduct experiments on two concurrent works DA-BEV [39] and Sparse4D [15], who design their own modules to address the depth ambiguity problem implicitly or through post-refinement. Differently, we solve the problem from the root at the feature lifting process, which is a more principled solution. The results show that, even with their efforts as a foundation, DFA3D achieves +0.7 mAP, +0.8 mAP and +1.0 mAP improvement over DA-BEV-S, DA-BEV [39] and Sparse4D [15] respectively, which verifies the necessity of our new feature lifting approach.

⁵We use the BEVFormer-t, BEVFormer-s, and BEVFormer-b for the BEVFormer-tiny, BEVFormer-small, and BEVFormer-base variants in <https://github.com/fundamentalvision/BEVFormer>.

Table 5. The effect of different numbers of DFA3D-based feature lifting used in the BEVFormer’s encoder layers. N_{2D} and N_{3D} are the numbers of layers that use DFA2D-based and DFA3D-based feature lifting respectively.

N_{2D}	N_{3D}	mAP↑	mATE↓	mASE↓	mAOE↓	mAVE↓	mAAE↓	NDS↑
4	2	42.0	66.9	27.1	35.3	34.2	19.1	52.7
2	4	42.7	66.0	27.5	38.0	32.6	20.2	52.9
0	6	43.0	65.4	27.1	37.4	34.1	20.5	53.1

Table 6. The effects of DepthNet module and depth quality. “Depth Sup.” means that only supervise the DepthNet but do not use the estimated depth in feature lifting. “Unsup.”, “Sup.”, and “GT” denote the depth used in feature lifting, that is unsupervised learned, supervised learned, and ground truth (from LiDAR).

Row	Method	mAP↑	mATE↓	mASE↓	mAOE↓	mAVE↓	mAAE↓	NDS↑
1	BEVFormer	41.6	67.3	27.4	37.2	39.4	19.8	51.7
2	BEVFormer-Depth Sup.	42.0	66.9	27.4	37.6	38.6	19.9	52.1
3	BEVFormer-DFA3D Unsup.	41.7	67.3	28.0	36.7	35.3	19.6	51.8
4	BEVFormer-DFA3D Sup.	43.0	65.4	27.1	37.4	34.1	20.5	53.1
5	BEVFormer-DFA3D GT	56.7	45.3	26.8	24.7	32.8	19.2	62.4

4.4. Ablations

Mixture of 3D and 2D Deformable Attention-based Feature Lifting. We verify the effectiveness of our DFA3D-based feature lifting by mixing DFA2D-based and DFA3D-based feature lifting in the encoder of BEVFormer. As shown in Table 5, in the first N_{2D} layers we use the DFA2D-based feature lifting, and in the following N_{3D} layers we use the DFA3D-based one. The results show a consistent improvement with the usage of DFA3D-based feature lifting increases. The progressive improvements indicate that the more DFA3D introduced, the less depth ambiguity problem is, which finally results in a better performance.

Effects of DepthNet Module and Depth Quality. As shown in Table 6, we first present the influence of the DepthNet module. Compared with the 1.4% mAP gain (Row 4 vs. Row 1) brought by DFA3D, simply equipping BEVFormer with DepthNet and supervising it brings in a 0.4% mAP gain (Row 2 vs. Row 1), which is relatively marginal. The comparison indicates that the main improvement does not come from the additional parameters and supervision of DepthNet, but from a better feature lifting.

We then evaluate the influence of depth qualities. We experiment with three different depth, 1) learned from unsupervised learning, 2) learned from supervised learning, and 3) ground truths generated from LiDAR. The corresponding results of these different depth qualities are available in Rows 3, 4, and 5 in Table 6. The results show that the performance improves with better depth quality. Remarkably, the model variant with ground truth depth achieves 56.7% mAP, outperforming baselines by 15.1% mAP. It indicates a big improvement space which deserves a further in-depth study of obtaining better depth estimation.

Table 7. Comparison of different feature querying methods. Trilinear indicates querying features on the expanded 3D feature maps only at the reference points.

Transformation	mAP↑	mATE↓	mASE↓	mAOE↓	mAVE↓	mAAE↓	NDS↑
Trilinear	40.6	68.7	27.7	38.9	35.5	18.9	51.3
DFA3D	43.0	65.4	27.1	37.4	34.1	20.5	53.1

Table 8. The evaluation of efficiency and resource consumption of DFA3D based on BEVFormer.

Method	Speed (FPS)	GPU Mem (GB)	#Param (M)
BEVFormer	2.8	6.9	62.1
BEVFormer-DFA3D	2.4	7.0	68.6

Effects of Feature Querying Methods. Deformable attention enables querying features dynamically. It generates sampling points according to the queried features from the previous layer. Inspired by previous work [32], we freeze the position of sampling points, under which the attention operation actually degenerates to the trilinear interpolation. We compare the degenerated one with the default one in Table 7. The results show the effectiveness of the deformable operation. The dynamic sampling can help focus on more important positions and enable better feature querying.

Efficiency and resource consumption. We evaluate the efficiency and resource consumption of DFA3D based on BEVFormer, as shown in Table 8. DFA3D-based feature lifting requires a little more resources and is slightly slower. We clarify that the side effects, especially the extra parameters mainly come from DepthNet, which is adjustable and can be optimized with the development of depth estimation.

5. Conclusion

In this paper, we have presented a basic operator called 3D deformable attention (DFA3D), built upon which we develop a novel feature lifting approach. Such an approach not only takes depth into consideration to tackle the problem of depth ambiguity but also benefits from the multi-layer refinement mechanism. We seek assistance from math to simplify DFA3D and develop a memory-efficient implementation. The simplified DFA3D makes querying features in 3D space through deformable attention possible and efficient. The experimental results show a consistent improvement, demonstrating the superiority and generalization ability of DFA3D-based feature lifting.

Limitations and Future Work. In this paper, we simply take monocular depth estimation to provide depth information for DFA3D, which is not accurate and stable enough. Recently proposed methods [16, 29] have proved that long temporal information can provide the network with a better depth sensing ability. How to make full use of the superior depth sensing ability to generate high-quality depth maps explicitly, and utilize them to help improve the performance of DFA3D-based feature lifting is still an issue, which we leave as future work.

References

- [1] Tom Bruls, Horia Porav, Lars Kunze, and Paul Newman. The Right (Angled) Perspective: Improving the Understanding of Road Scenes Using Boosted Inverse Perspective Mapping. In *2019 IEEE Intelligent Vehicles Symposium (IV)*, pages 302–309. IEEE, 2019. **1**
- [2] 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 multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. **3, 6**
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End Object Detection with Transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020. **2, 3**
- [4] Zehui Chen, Zhenyu Li, Shiquan Zhang, Liangji Fang, Qin-hong Jiang, and Feng Zhao. Graph-DETR3D: Rethinking Overlapping Regions for Multi-View 3D Object Detection. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 5999–6008, 2022. **2, 3**
- [5] Xiaomeng Chu, Jiajun Deng, Yuan Zhao, Jianmin Ji, Yu Zhang, Houqiang Li, and Yanyong Zhang. OA-BEV: Bringing Object Awareness to Bird’s-Eye-View Representation for Multi-Camera 3D Object Detection. *arXiv preprint arXiv:2301.05711*, 2023. **2, 3, 4**
- [6] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable Convolutional Networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017. **7**
- [7] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep Ordinal Regression Network for Monocular Depth Estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2002–2011, 2018. **4**
- [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. **7**
- [9] Junjie Huang and Guan Huang. BEVDet4D: Exploit Temporal Cues in Multi-camera 3D Object Detection. *arXiv preprint arXiv:2203.17054*, 2022. **1**
- [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, 3, 4**
- [11] 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**
- [12] Yin hao Li, Han Bao, Zheng Ge, Jinrong Yang, Jianjian Sun, and Zeming Li. BEVStereo: Enhancing Depth Estimation in Multi-View 3D Object Detection with Temporal Stereo. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 1486–1494, 2023. **1, 3, 4**
- [13] Yin hao 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. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 1477–1485, 2023. **1, 3, 4, 6**
- [14] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. BEVFormer: Learning Bird’s-Eye-View Representation from Multi-Camera Images via Spatiotemporal Transformers. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*, pages 1–18. Springer, 2022. **1, 2, 3, 6, 7**
- [15] Xuewu Lin, Tianwei Lin, Zixiang Pei, Lichao Huang, and Zhizhong Su. Sparse4D: Multi-view 3D Object Detection with Sparse Spatial-Temporal Fusion. *arXiv preprint arXiv:2211.10581*, 2022. **2, 3, 7**
- [16] Xuewu Lin, Tianwei Lin, Zixiang Pei, Lichao Huang, and Zhizhong Su. Sparse4D v2: Recurrent Temporal Fusion with Sparse Model. *arXiv preprint arXiv:2305.14018*, 2023. **8**
- [17] 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. In *International Conference on Learning Representations*, 2021. **3**
- [18] Shilong Liu, Tianhe Ren, Jiayu Chen, Zhaoyang Zeng, Hao Zhang, Feng Li, Hongyang Li, Jun Huang, Hang Su, Jun Zhu, and Lei Zhang. Detection Transformer with Stable Matching, 2023. **3**
- [19] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection. *arXiv preprint arXiv:2303.05499*, 2023. **3**
- [20] Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. PETR: Position Embedding Transformation for Multi-view 3D Object Detection. In *European Conference on Computer Vision*, pages 531–548, 2022. **1, 2, 3**
- [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, 3**
- [22] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*, 2018. **7**
- [23] Xuran Pan, Zhuofan Xia, Shiji Song, Li Erran Li, and Gao Huang. 3D Object Detection With Pointformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7463–7472, 2021. **1**
- [24] 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. **1**
- [25] 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, 3**

- [26] Lennart Reiher, Bastian Lampe, and Lutz Eckstein. A Sim2Real Deep Learning Approach for the Transformation of Images from Multiple Vehicle-Mounted Cameras to a Semantically Segmented Image in Bird’s Eye View. In *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pages 1–7. IEEE, 2020. [1](#)
- [27] Tianhe Ren, Shilong Liu, Feng Li, Hao Zhang, Ailing Zeng, Jie Yang, Xingyu Liao, Ding Jia, Hongyang Li, He Cao, Jianan Wang, Zhaoyang Zeng, Xianbiao Qi, Yuhui Yuan, Jianwei Yang, and Lei Zhang. detrex: Benchmarking Detection Transformers, 2023. [3](#)
- [28] Tianhe Ren, Jianwei Yang, Shilong Liu, Ailing Zeng, Feng Li, Hao Zhang, Hongyang Li, Zhaoyang Zeng, and Lei Zhang. A Strong and Reproducible Object Detector with Only Public Datasets, 2023. [3](#)
- [29] Shihao Wang, Yingfei Liu, Tiancai Wang, Ying Li, and Xiangyu Zhang. Exploring Object-Centric Temporal Modeling for Efficient Multi-View 3D Object Detection. *arXiv preprint arXiv:2303.11926*, 2023. [8](#)
- [30] 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. [1](#)
- [31] 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. [1](#), [7](#)
- [32] 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](#), [3](#), [6](#), [7](#), [8](#)
- [33] Yue Wang and Justin M Solomon. Object DGCNN: 3D Object Detection using Dynamic Graphs. *Advances in Neural Information Processing Systems*, 34:20745–20758, 2021. [1](#)
- [34] Zengran Wang, Chen Min, Zheng Ge, Yinhao Li, Zeming Li, Hongyu Yang, and Di Huang. STS: Surround-view Temporal Stereo for Multi-view 3D Detection. *arXiv preprint arXiv:2208.10145*, 2022. [1](#), [3](#), [4](#)
- [35] Yan Yan, Yuxing Mao, and Bo Li. SECOND: Sparsely Embedded Convolutional Detection. *Sensors*, 18(10):3337, 2018. [7](#)
- [36] Chenyu Yang, Yuntao Chen, Hao Tian, Chenxin Tao, Xizhou Zhu, Zhaoxiang Zhang, Gao Huang, Hongyang Li, Yu Qiao, Lewei Lu, et al. BEVFormer v2: Adapting Modern Image Backbones to Bird’s-Eye-View Recognition via Perspective Supervision. *arXiv preprint arXiv:2211.10439*, 2022. [2](#), [3](#)
- [37] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-based 3D Object Detection and Tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11784–11793, 2021. [1](#)
- [38] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel Ni, and Heung-Yeung Shum. DINO: DETR with Improved DeNoising Anchor Boxes for End-to-End Object Detection. In *The Eleventh International Conference on Learning Representations*, 2022. [3](#)
- [39] Hao Zhang, Hongyang Li, Xingyu Liao, Feng Li, Shilong Liu, Lionel M. Ni, and Lei Zhang. DA-BEV: Depth Aware BEV Transformer for 3D Object Detection, 2023. [2](#), [3](#), [7](#)
- [40] Yin Zhou and Oncel Tuzel. VoxelNet: End-to-End Learning for Point Cloud Based 3D Object Detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4490–4499, 2018. [1](#)
- [41] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: Deformable Transformers for End-to-End Object Detection. In *International Conference on Learning Representations*, 2020. [3](#)

Appendix

A. Closer visualization of the expanded 3D image features.

In Sec. 3.1, for a specific view (the n^{th} view for example), we obtain the expanded 3D image feature map by conducting outer product at last dimension between the estimated depth distribution $D_n \in \mathbb{R}^{H \times W \times D}$ and 2D image feature maps $X_n \in \mathbb{R}^{H \times W \times C}$. To further explain the resulting expanded 3D image feature $F_n \in \mathbb{R}^{H \times W \times D \times C}$, we show a closer visualization in Fig. 4. For the features with the same (u_i, v_j) coordinate (zoomed in), they share the same 2D image feature X_{n,u_i,v_j} and depth distribution D_{n,u_i,v_j} . However, since they have different depth values d_k , they select different depth confidence scores D_{n,u_i,v_j,d_k} , resulting in different features $\{D_{n,u_i,v_j,d_k} \cdot X_{n,u_i,v_j}\}$.

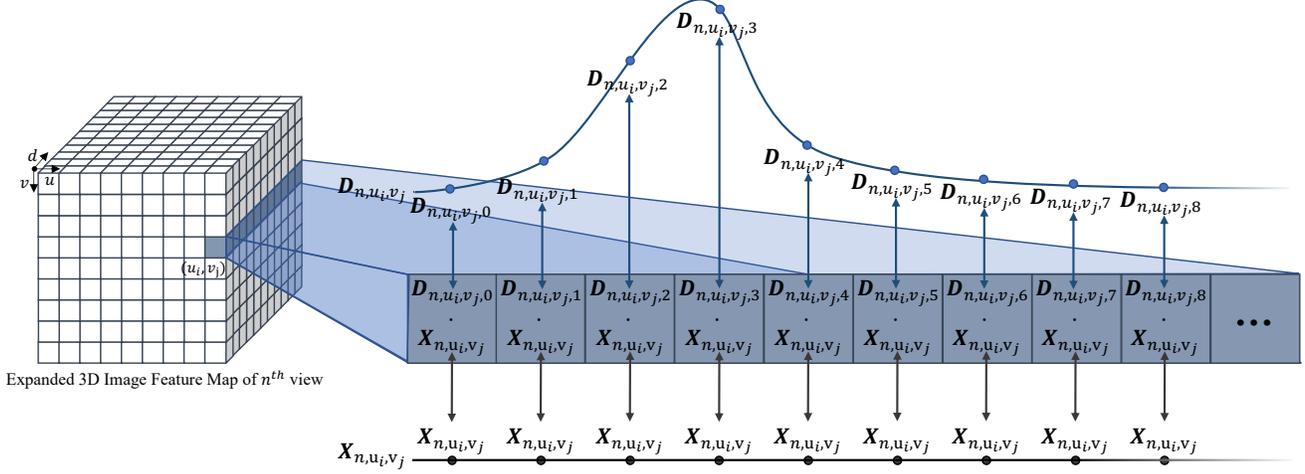


Figure 4. A closer visualization of the expanded 3D image features.

B. Can 3D deformable attention be trivially implemented by feature weighting followed with a common 2D deformable attention?

As each sampling point in 3D deformable attention has its own 3D location, it will lead to its own depth score for 4 adjacent image features when conducting weighted bilinear interpolation. Inevitably, there will be features (X_{n,u_1,v_0} and X_{n,u_1,v_1} in Fig. 5) referred by more than one sampling point when sampling points are located in neighboring grids. In such cases, feature weighting will result in conflicts. Thus, we can not simply prepare 2D features by depth-weighting and then conduct the common 2D deformable attention. The depth-weighted feature computation should be conducted on the fly.

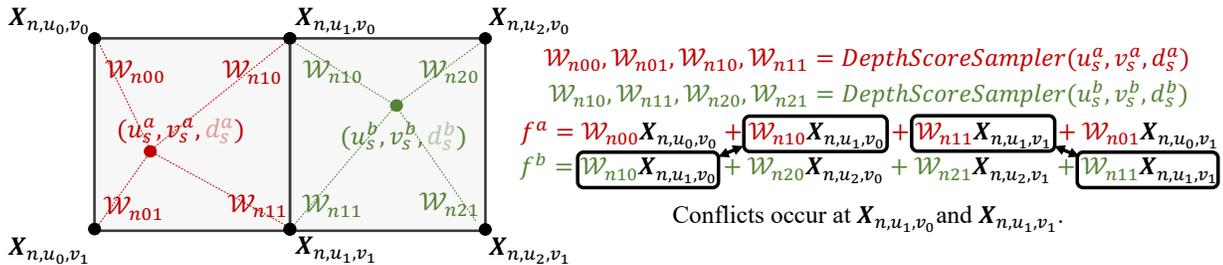


Figure 5. In 3D deformable attention, when sampling features for two sampling points (u_s^a, v_s^a, d_s^a) and (u_s^b, v_s^b, d_s^b) respectively, we first sample depth scores $(W_{n00}, W_{n10}, W_{n11}, W_{n01})$ and $(W_{n10}, W_{n20}, W_{n21}, W_{n11})$ for these two points respectively. After that, we conduct depth-weighted bilinear interpolation for these two points based on their own depth scores independently.

C. Applicability

As the comparison shown in Fig. 6, when integrating our DFA3D in any 2D deformable attention-based feature lifting only requires a few modifications in code. The main modification lies in the addition of DepthNet and replacing 2D deformable attention with our 3D deformable attention.

<pre> multiview_feature_maps = Backbone(multiview_input_RGB) multiview_reference_point = ProjEgo2Pixel(positions_BEV, extrinsic_param, intrinsic_param) </pre>	<pre> multiview_feature_maps = Backbone(multiview_input_RGB) multiview_depth_dist = DepthNet(multiview_feature_maps) multiview_reference_point = ProjEgo2Pixel(positions_BEV, multiview_extrinsic_param, multiview_intrinsic_param) </pre>
<pre> multiview_reference_point = multiview_reference_point[... :2] # [N 3] → [N 2] </pre>	
<pre> for layer_id in range(num_layers): multiview_sample_offset = MLP_2D(content_BEV) # [N C] → [N K 2] multiview_sample_point = multiview_reference_point + multiview_sample_offset content_BEV = DeformableAttention_2D(</pre>	<pre> for layer_id in range(num_layers): multiview_sample_offset = MLP_3D(content_BEV) # [N C] → [N K 3] multiview_sample_point = multiview_reference_point + multiview_sample_offset content_BEV = DeformableAttention_3D(</pre>
<pre> multiview_feature_maps_2D, multiview_sample_point) lifted_feature = content_BEV </pre>	<pre> multiview_depth_dist, multiview_feature_maps_2D, multiview_sample_point) lifted_feature = content_BEV </pre>
2D Deformable Attention-based Feature Lifting	3D Deformable Attention-based Feature Lifting

Figure 6. Integrate DFA3D into any 2D deformable attention-based feature lifting requires only a few modifications in code.

D. Visualization

We visualize the predictions of BEVFormer and BEVFormer-DFA3D in Fig. 7. The shaded triangles correspond to camera rays, in which BEVFormer makes more duplicate predictions behind or in front of ground truth objects compared with BEVFormer-DFA3D. It demonstrates the negative effects of depth ambiguity in 2D deformable attention. After integrating with DFA3D, the wrong predictions caused by the depth ambiguity problem are reduced.

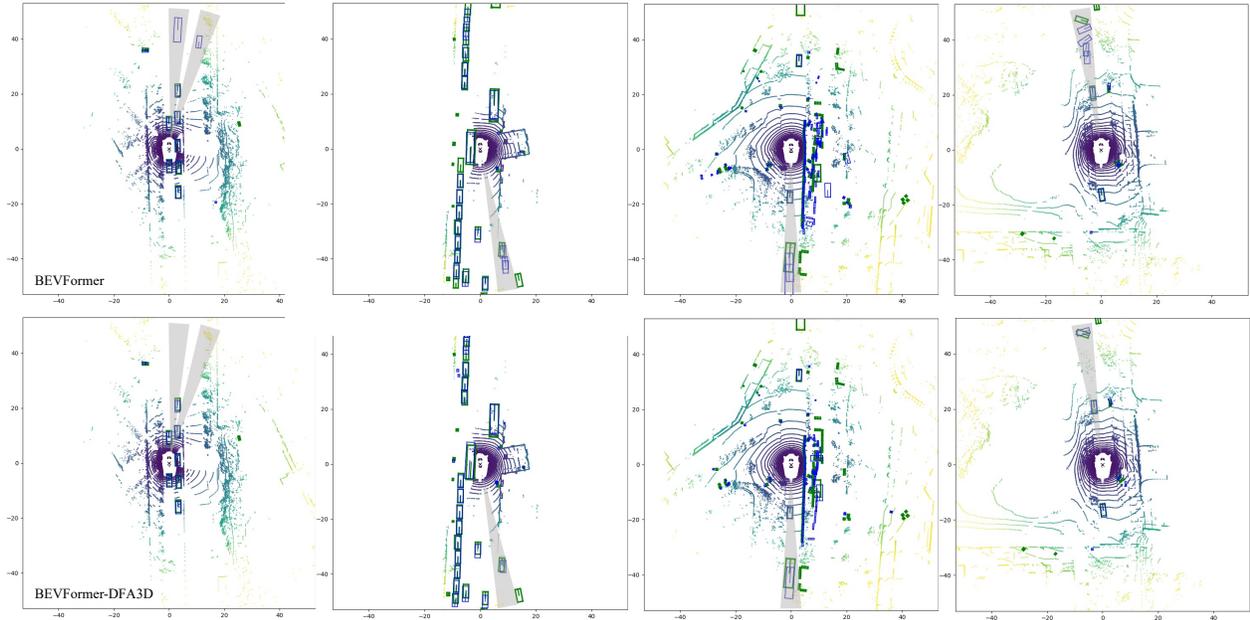


Figure 7. The visualization of predictions in BEV using green boxes to represent ground truth boxes and blue boxes to represent predicted ones. Duplicate predictions, which are caused by the depth ambiguity problem, are enclosed by shaded triangles.