

Bridged Transformer for Vision and Point Cloud 3D Object Detection

Yikai Wang¹ TengQi Ye² Lele Cao¹ Wenbing Huang³
Fuchun Sun¹✉ Fengxiang He⁴ Dacheng Tao⁴

¹Beijing National Research Center for Information Science and Technology (BNRist),
State Key Lab on Intelligent Technology and Systems,

Department of Computer Science and Technology, Tsinghua University ²ByteDance Inc.

³Institute for AI Industry Research (AIR), Tsinghua University ⁴JD Explore Academy, JD.com

wangyk17@mails.tsinghua.edu.cn, yetengqi@gmail.com, caolele@gmail.com, hwenbing@126.com,
fuchuns@tsinghua.edu.cn, fengxiang.f.he@gmail.com, dacheng.tao@gmail.com

Abstract

3D object detection is a crucial research topic in computer vision, which usually uses 3D point clouds as input in conventional setups. Recently, there is a trend of leveraging multiple sources of input data, such as complementing the 3D point cloud with 2D images that often have richer color and fewer noises. However, due to the heterogeneous geometrics of the 2D and 3D representations, it prevents us from applying off-the-shelf neural networks to achieve multimodal fusion. To that end, we propose Bridged Transformer (BrT), an end-to-end architecture for 3D object detection. BrT is simple and effective, which learns to identify 3D and 2D object bounding boxes from both points and image patches. A key element of BrT lies in the utilization of object queries for bridging 3D and 2D spaces, which unifies different sources of data representations in Transformer. We adopt a form of feature aggregation realized by point-to-patch projections which further strengthen the interaction between images and points. Moreover, BrT works seamlessly for fusing the point cloud with multi-view images. We experimentally show that BrT surpasses state-of-the-art methods on SUN RGB-D and ScanNetV2 datasets.

1. Introduction

3D object detection, which aims at identifying or locating objects in 3D scenes, is drawing increasing attention and is acting as a fundamental task towards scene understanding. Many successful attempts [3, 15, 21, 25] have been made using point cloud data as input. These attempts include converting the points to regular format (e.g., 3D voxel grids [32], polygon meshes [11], multi-views [29]), or using

3D specific operators (e.g., symmetric functions [23], voting [21]) to design grouping strategies for points. In addition, since Transformers could be naturally permutation invariant and capable of capturing large-scale data interaction, they are lately applied to 3D object detection and demonstrate superior performance [15, 19]. Besides handling point cloud learning tasks, Transformers have swept across various 2D tasks, e.g., image classification [6, 14], object detection [2, 8, 39], and semantic segmentation [33, 37].

Deep multimodal learning by leveraging the advantage of multiple modalities has shown its superiority on various applications [1, 31]. Despite the success of Transformers in 2D or 3D single-modal object detection tasks, the attempt of combining advantages from both point clouds and images remains scarce. For 3D learning tasks, the point cloud provides essential geometrical cues, while the information in rich color images can complement the point cloud by fulfilling the missing color information and correcting noise errors. As a result, the performance of 3D object detection could be potentially improved by the involvement of 2D images. One intuitive method is to lift 3-dimensional RGB vectors from images to extend the point features. A CNN-based 3D detection model, imVoteNet [20], points out the difficulty in migrating 2D/3D discrepancies by this intuitive method, and instead, imVoteNet substitutes the RGB vectors with image features extracted by a pre-trained 2D detector. However, simultaneously relying on both the image voting and point cloud voting assumptions in [20] could accumulate the intrinsic grouping errors as mentioned by [15]. To avoid the learning process of point clouds being impacted by middle-level 2D/3D feature interaction, [20] combines multimodal features over the first layer, which potentially prevents the network from fully exploiting their semantic interaction or migrating multimodal discrepancies.

In this work, we propose Bridged Transformer (BrT) – a

✉ Corresponding author: Fuchun Sun.

simple and effective Transformer framework for 3D object detection. BrT bridges the learning processes of images and point clouds inside Transformer. This approach takes the sampled points and image patches as input. To protect the self-learning process of each modality, attentions between point tokens and image patch tokens are blocked but bridged by object queries throughout the Transformer layers. To strengthen the interaction between images and points, BrT is also equipped with powerful bridging designs from two perspectives. Firstly, we leverage conditional object queries for images and points that are aware of the learned proposal points. Such design together with aligned positional embeddings tells Transformer that object queries of images and points are aligned. Secondly, despite the perspective from object queries, we perform point-to-patch projections to explicitly leverage the spatial relationships of both modalities. BrT avoids the grouping errors due to its natural ability of capturing long-range dependencies and global contextual information, and instead of lifting image features to point clouds at the beginning layer in [20], BrT allows the full propagation of feature interaction in the whole network. As an additional advantage, BrT can be extended to combine point clouds with multi-view images.

We evaluate BrT on both SUN RGB-D and ScanNetV2 datasets, where respectively, BrT achieves remarkably 2.4% and 2.2% improvements over state-of-the-art methods.

To summarize, the contributions of our work are:

- We propose BrT, a novel framework for 3D object detection that bridges the learning processes of images and point clouds inside Transformer.
- We propose to strengthen the interaction between images and points from two perspectives including conditional object queries and the point-to-patch projection.
- BrT achieves the state-of-the-art on two benchmarks, which demonstrates the superiority of our design and also the potential in multi-view scenarios.

2. Related Work

3D detection with point cloud. There are unique challenges faced by the processing of point clouds using deep neural networks (DNNs) [3, 15, 19, 21, 25, 36]. A detailed discussion around this difference can be found in [9]. The targets of object detection in 3D space are locating 3D bounding boxes and recognizing the object classes. VoxelNet [38] proposes to divide a point cloud into equally spaced 3D voxels, and then transforms the points in each voxel into a unified feature representation. VoteNet [21] reformulates Hough voting in the context of deep learning to generate better points for box proposals with grouping. Transformers are also adapted to become suitable for handling 3D points. 3DTR [19] introduces an end-to-end

Transformer with non-parametric queries and Fourier positional embeddings. Group-Free [15] adopts the attention mechanism to learn the point features, which potentially retains the information of all points to avoid the errors of previous grouping strategies. Voxel Transformer [17] effectively captures the long-range relationships between voxels.

3D detection with multimodal data. There are a few works that use deep networks to combine point clouds and images. MV3D [4] proposes an element-wise fusion of representations from different domains, based on the rigid assumption that all objects are on the same spatial plane and can be pinpointed solely from a top-down view of the point cloud. PointFusion [34] concatenates point cloud features and image features at two different levels to learn their interaction, which could not guarantee the alignment of features. ImVoteNet [20] lifts crafted semantic and texture features to the 3D seed points for fusion. However, ImVoteNet is still negatively affected by the errors of grouping and combining features only at the beginning layer, leading to highly restricted feature interaction. Different from the aforementioned methods, our BrT fully exploits the feature interaction for images and points with additional bridging processes to strengthen the interaction.

Transformer for 2D detection. Recently, Transformer achieves the cutting edge performance in computer vision tasks [2, 6, 8, 14, 18, 35, 39]. For 2D object detection based on images, DETR [2] enables the Transformer to learn relations of the objects and the global image context to directly output the final set of predictions; and it also removes the need for non-maximum suppression and anchor generation. With the help of pre-training, YOLOX [8] proposes a pure sequence-to-sequence approach that achieves competitive performance for object detection; hence it also tackles the transferability of Transformer from image recognition to object detection. Deformable DETR [39] is an efficient and fast-converging model with attention modules only paying attention to a small set of tokens instead of the whole contexts. Conditional DETR [18] learns a conditional spatial query aiming to accelerate the training process.

3. Method

In this section, we propose Bridged Transformer (BrT) for 3D object detection with both vision and point cloud as input. We describe the overall structure of BrT in Sec. 3.1, followed by the design of building blocks in Sec. 3.2. We consider two aspects to bridge the learning processes of vision and point cloud in Sec. 3.3 and Sec. 3.4, respectively.

3.1. Overall architecture

An overall architecture of our BrT is depicted in Fig. 1. Suppose we are given $N \times 3$ points representing the 3D coordinates, and an $H \times W \times 3$ image. Here, N is the number

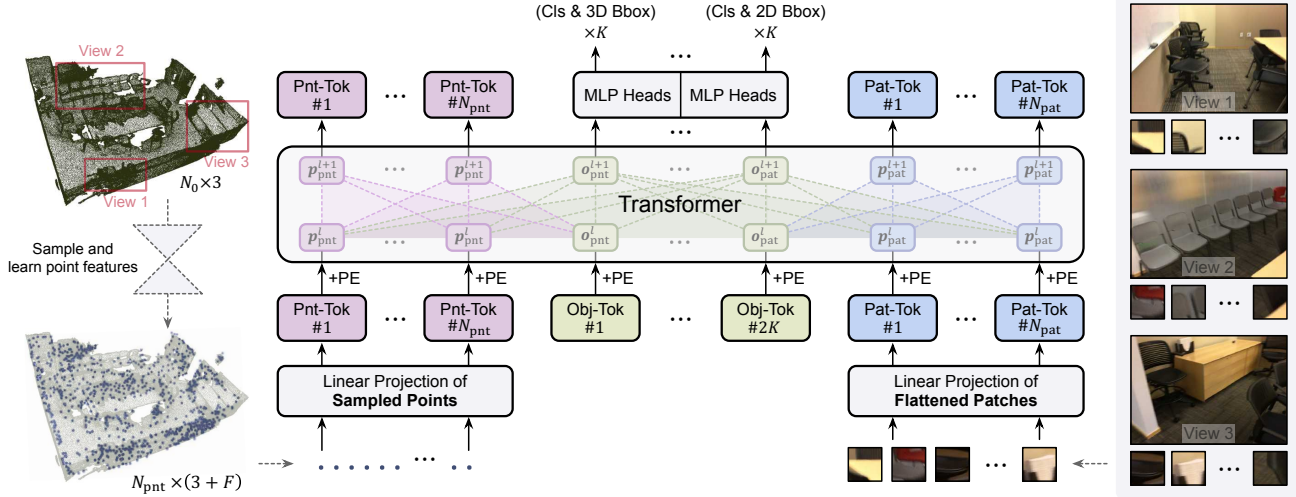


Figure 1. Overall architecture of our Bridged Transformer (BrT) for 3D object detection based on point clouds and single-view/multi-view images. For each image view, we annotate its corresponding region on the point cloud for better readability.

of points; H and W are the height and width of the image respectively. For simplicity, we first analyze one image per scene since it matches the common scenario where the camera sensors capture (depth) points and RGB at the same time. Yet our method can be extended to handle multiple images per scene with different views, at one’s disposal, as described in Sec. 3.5 and evaluated by our experiments.

Before feeding the **point cloud data** to the first Transformer stage¹, we process the data with the method adopted in [21]. Specifically, we first sample $N_{\text{pnt}} \times (3 + F)$ “seed points” from a total of $N_0 \times 3$ points, using PointNet++. Note that N_{pnt} denotes the number of sampled points; The positive integers 3 and F represent the dimension of the 3D Euclidean coordinate and point feature, respectively.

For processing the **image data**, we follow some successful practices from vision Transformers. Concretely, each image is evenly partitioned into N_{pat} patches before embedded by a multi-layer perceptron (MLP). Together with the embedded images patches, the learned object queries are sent to the model, generating output embeddings that are used to predict box coordinates and class labels.

Moreover, we adopt $2K$ learnable object queries, among which K queries for points and K for image patches. In summary, we have $N_{\text{pnt}} + N_{\text{pat}}$ basic tokens and $2K$ object queries tokens. Suppose the hidden dimension is D , The token features fed to the l -th ($l = 1, \dots, L$) Transformer stage contains point tokens $\mathbf{p}_{\text{pnt}}^l \in \mathbb{R}^{N_{\text{pnt}} \times D}$, patch tokens $\mathbf{p}_{\text{pat}}^l \in \mathbb{R}^{N_{\text{pat}} \times D}$, object queries for points $\mathbf{o}_{\text{pnt}}^l \in \mathbb{R}^{K \times D}$, and object queries for patches $\mathbf{o}_{\text{pat}}^l \in \mathbb{R}^{K \times D}$.

When given camera intrinsic and extrinsic parameters,

¹Here, each stage contains a multi-head self-attention, an MLP, and two layer normalizations.

each 3D point could be projected to the camera plane, which builds the relation between 3D coordinates and 2D image pixels. We define the projection operator $\text{proj} : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ indicating the projection process from a 3D point coordinate $\mathbf{k} = [x, y, z]^T$ to a 2D pixel coordinate $\mathbf{k}' = [u, v]^T$ on the corresponding image, and there is

$$\mathbf{k}' = \text{proj}(\mathbf{k}) = \mathbf{\Pi} \begin{bmatrix} \frac{1}{4} & 0 & 0 \\ 0 & \frac{1}{4} & 0 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{K} \mathbf{R}_t \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}, \quad (1)$$

where \mathbf{K} and \mathbf{R}_t are the intrinsic and extrinsic matrices, and $\mathbf{\Pi}$ is a perspective mapping.

BrT has $2K$ outputs which correspond to the $2K$ input object queries. An MLP is applied to the first K outputs for predicting the coordinates of 3D boxes and their class labels. For the rest of the K outputs, we use a different MLP to predict the 2D coordinates of the bounding boxes and their associated classes. It is worth mentioning that we do not need extra labels for 2D box coordinates, since they are obtained by first projecting the labels of 3D box coordinates to the 2D camera plane following Eq. (1), and then taking the axis-aligned 2D bounding boxes of projected shapes.

The optimization of BrT concerns minimizing a compound loss function that contains two parts: a regression loss for locating bounding boxes, and a classification loss for predicting the class of the associated box. The regression loss contains two components: $\mathcal{L}_{\text{obj}}^{3D}$ and $\mathcal{L}_{\text{obj}}^{2D}$ for 3D and 2D cases respectively. Likewise, for classification loss, there are also a 3D component $\mathcal{L}_{\text{cls}}^{3D}$ and a 2D component $\mathcal{L}_{\text{cls}}^{2D}$. As such, the overall loss function is formulated as

$$\mathcal{L} = \mathcal{L}_{\text{obj}}^{3D} + \alpha_1 \mathcal{L}_{\text{cls}}^{3D} + \alpha_2 \mathcal{L}_{\text{obj}}^{2D} + \alpha_3 \mathcal{L}_{\text{cls}}^{2D}, \quad (2)$$

where α_1, α_2 and α_3 are three parameters weighting relative importance between these losses. In practice, $\mathcal{L}_{\text{obj}}^{3D}, \mathcal{L}_{\text{obj}}^{2D}$ or $\mathcal{L}_{\text{cls}}^{3D}$ further consists of two sub-terms. These details will be provided in Sec. 4.1.

3.2. Transformer building block of BrT

The multi-head attention (MSA) acts as the fundamental building block of the Transformer architecture, which has three sets of input: query set, key set, and value set. Given a query set $\{q_i\}$ and a common element set $\{p_k\}$ of key set and value set, the output feature of MSA of each query element is the aggregation of the values that weighted by linear projections, formulated as:

$$\text{Att}(q_i, \{p_k\}) = \sum_{h=1}^H W_h \left(\sum_{k=1}^K A_{i,k}^h \cdot V_h p_k \right), \quad (3)$$

where h indexes over H attention heads; W_h and V_h are weights for the output projection and the value projection, respectively. $A_{i,k}^h$ is the attention weight which is computed as

$$A_{i,k}^h = \frac{\exp[(Q_h q_i)^\top (U_h p_k)]}{\sum_{k=1}^K \exp[(Q_h q_i)^\top (U_h p_k)]}, \quad (4)$$

where Q_h and U_h indicate the query projection weight and the key projection weight, respectively.

As shown in Fig. 1, during training, the BrT Transformer module prohibits the attentions between p_{pnt}^l and p_{pat}^l . Therefore, these tokens are only directly attentive to

$$\begin{aligned} p_{\text{pnt},i}^{l+1} &= \text{Att}(p_{\text{pnt},i}^l, \{o_{\text{pnt},j}^l, p_{\text{pnt},j}^l\}), \\ p_{\text{pat},i}^{l+1} &= \text{Att}(p_{\text{pat},i}^l, \{o_{\text{pat},j}^l, p_{\text{pat},j}^l\}), \end{aligned} \quad (5)$$

where subscript i and j denote token indexes. Here, calculating the attention scores based on tokens in the given scope yet leaving other tokens unattended can be achieved by applying zero masks on the attention.

Although we block the attentions between p_{pnt}^l and p_{pat}^l in the Transformer modules, we choose to explicitly connect them again by adopting the point-to-patch projection, which will be established in Sec. 3.4. Instead of allowing global attention in Transformer, we believe that the large discrepancy of coordinates essentially encourages a form of explicit 3D-2D interaction. We verify this assumption experimentally in Sec. 4.4.

In addition, o_{pnt}^l and o_{pat}^l have additional functions to further bridge the gap between 3D and 2D coordinates, which will be further detailed in Sec. 3.3. To this end, object queries o_{pnt}^l and o_{pat}^l are attentive to all tokens,

$$\begin{aligned} o_{\text{pnt},i}^{l+1} &= \text{Att}(o_{\text{pnt},i}^l, \{o_{\text{pnt},j}^l, o_{\text{pat},j}^l, p_{\text{pnt},j}^l, p_{\text{pat},j}^l\}), \\ o_{\text{pat},i}^{l+1} &= \text{Att}(o_{\text{pat},i}^l, \{o_{\text{pnt},j}^l, o_{\text{pat},j}^l, p_{\text{pnt},j}^l, p_{\text{pat},j}^l\}). \end{aligned} \quad (6)$$

3.3. Bridge by conditional object queries

The Euclidean coordinates of the 3D point cloud may vary dramatically from the camera plane coordinates of 2D image pixels, since they belong to different spaces. As a result, it may be hard for Transformer-based model to learn their relationships even after numerical normalization. In this part, we propose to leverage object queries to bridge the 3D and 2D spaces. Specifically, we adopt conditional object queries which are aware of both the 3D and 2D coordinates.

For Transformer-based object detection models, object queries are observed to probably specialize on certain areas and box sizes during the training process, even they are generated by random initialization [2]. Inspired by this, we assume that hidden features of object queries w.r.t points and images could be potentially aligned inside the Transformer. Hence instead of using randomly generated object queries, we adopt conditional object queries to boost the prediction learning process based on the object query alignment of points and images.

To align object queries of points and image patches, we first sample K points as proposals from the N_{pnt} points with kNN search, and denote the 3D coordinates and features of these K points as $k_{\text{pnt}} \in \mathbb{R}^{K \times 3}$ and $f_{\text{pnt}} \in \mathbb{R}^{K \times F}$, respectively. We then learn the 3D coordinates of proposals, denoted as $k'_{\text{pnt}} \in \mathbb{R}^{K \times 3}$, by adding k_{pnt} with additional learned biases based on f_{pnt} . The object queries for points o_{pnt}^1 are conditioned on the k'_{pnt} . Formally, there is

$$k'_{\text{pnt}} = k_{\text{pnt}} + \text{MLP}(f_{\text{pnt}}), \quad (7)$$

$$o_{\text{pnt}}^1 = \text{MLP}(k'_{\text{pnt}}) + \text{PE}, \quad (8)$$

where $\text{PE} \in \mathbb{R}^{K \times D}$ is the randomly initialized positional embeddings.

Regarding the object queries of image patches o_{pat}^1 , we project k'_{pnt} to the corresponding image and obtain the 2D coordinates of projected pixels, denoted as $\text{proj}(k'_{\text{pnt}}) \in \mathbb{R}^{K \times 2}$ where proj indicates projecting points to the image according to the per-point projection in Eq. (1). The object queries for image patches are conditionally obtained by

$$o_{\text{pat}}^1 = \text{MLP}(\text{proj}(k'_{\text{pnt}})) + \text{PE}, \quad (9)$$

where PE is the same positional embeddings as in Eq. (8). Sharing positional embeddings intuitively tells the Transformer that both object queries o_{pnt}^1 and o_{pat}^1 are aligned.

Up till here, we improve the designs of object queries to bridge the representation spaces of 3D point clouds and 2D images. Such design is described by Eq. (6), Eq. (8), and Eq. (9), which actually embody two perspectives: the attentive connections and the alignment of 3D and 2D object queries using a shared PE . Ablation studies in Sec. 4.4 verify the effectiveness of both components we propose.

3.4. Bridge by point-to-patch projection

Apart from bridging point tokens p_{pnt}^l and patch tokens p_{pat}^l with object queries, we further strengthen their relations by adding intrinsic point-to-patch projection. Denoting the 3D coordinates of N_{pnt} sampled points as $\mathbf{n}_{\text{pnt}} \in \mathbb{R}^{N_{\text{pnt}} \times 3}$, we project \mathbf{n}_{pnt} to the corresponding camera plane and obtain N_{pnt} 2D pixel coordinates that are denoted with $\text{proj}(\mathbf{n}_{\text{pnt}}) \in \mathbb{R}^{N_{\text{pnt}} \times 2}$. With proj already defined in Eq. (9), we can conveniently let \mathbf{u}_n and \mathbf{v}_n be the x -axis value and the y -axis value respectively of the n -th element of $\text{proj}(\mathbf{n}_{\text{pnt}})$, where $n = 1, 2, \dots, N_{\text{pnt}}$. If \mathbf{u}_n and \mathbf{v}_n respectively satisfy restrictions $1 \leq \mathbf{u}_n \leq H$ and $1 \leq \mathbf{v}_n \leq W$, then such a 2D coordinate could reside in the input image of size $H \times W \times 3$. Rounding \mathbf{u}_n and \mathbf{v}_n to the nearest integers obtains the valid coordinates indicating certain image pixels. It is now easy to obtain the corresponding image patch index with

$$p_n = \lfloor \lfloor \mathbf{v}_n \rfloor / S \rfloor \times \lfloor W / S \rfloor + \lfloor \lfloor \mathbf{u}_n \rfloor / S \rfloor, \quad (10)$$

where $\lfloor \cdot \rfloor$ is the rounding operator; $p_n \in \{1, 2, \dots, N_{\text{pat}}\}$ is the corresponding patch index for the n -th point; S denotes the image patch size. Our point-to-patch projection aggregates features for both points and image patches by

$$p_{\text{pnt},n}^l = p_{\text{pnt},n}^l + \text{MLP}(p_{\text{pat},p_n}^l), \quad (11)$$

where subscripts n and p_n indicate indexes of the token features p_{pnt}^l and p_{pat}^l , respectively.

3.5. Extend to multiple-view scenarios

It is challenging to directly extend current point-only or point-image methods to detect from the point cloud and **multi-view** images, which is nevertheless a common real-life situation of data organization. For example, [20] avoided using the ScanNetV2 dataset which contains rich multi-view images, probably due to the difficulty in combining the point cloud and each view.

Fortunately, with few bells and whistles, our proposed BrT has a natural advantage in combining the point cloud with multi-view images, where both point-image interaction and the interaction of multi-view images can be utilized to further improve performance. As shown in Fig. 1, when there are different views of input images for one single scene, we concatenate these images along the width-side and obtain a wide image. The following processes are the same with the single-view condition. Since for the multi-view images, each view usually contains fewer objects, we expect the number of object queries K can still handle all the objects. Our current design mainly aims to bridge between each view and the point cloud, yet it does not exploit the relations among different views with explicit projections, which is left to be our future work.

4. Experiments

Our experiments are conducted on the challenging SUN RGB-D [27] and ScanNetV2 [5] datasets. We first detail the settings for datasets and implementations in Sec. 4.1. Then, we quantitatively compare our BrT with the state-of-the-art methods in Sec. 4.2; and we present and discuss the qualitative results in Sec. 4.3. Finally, we perform analytical experiments in Sec. 4.4 to verify the advantage of each component in BrT. More details of network architectures and visualizations are provided in our Appendix.

4.1. Datasets and implementation details

Datasets. SUN RGB-D [27] is a single-view RGB-D dataset for 3D scene understanding. It consists of 10,335 RGB-D images annotated with amodal oriented 3D bounding boxes for 37 object categories, alongside corresponding camera poses. The training and validation splits are composed of 5,285 and 5,050 frames respectively. We convert depth images to point clouds using the provided camera parameters, and adopt a standard evaluation protocol to report performance on the 10 most common categories [15,20,21].

ScanNetV2 [5] is a richly annotated dataset of 3D reconstructed meshes of indoor scenes. It contains 1,513 scans covering more than 700 unique indoor scenes, out of which 1,201 scans belong to the training split, and the rest 312 scans comprise the validation subset. ScanNetV2 contains over 2.5 million images with camera poses, and their corresponding reconstructed point clouds with 3D semantic annotation for 18 object categories. Compared to single-view scans in SUN RGB-D, scenes in ScanNetV2 are more complete and cover larger areas with multiple views. In our experiments, we adopt the sample dataset from ScanNetV2 containing 25,000 frames (scannet_frames_25k) which are sampled with 100 interval frames from the whole dataset.

Implementation details. Following the common successful practice in [15, 21], we adopt PointNet++ [24] as the point cloud backbone. The backbone has four set abstraction layers where the input point cloud is sub-sampled to 2,048, 1,024, 512, and 256 points with the increasing receptive radius of 0.2, 0.4, 0.8, and 1.2, respectively. There are two feature propagation layers which successively up-sample the points to 512 and 1,024, i.e. $N_{\text{pnt}}=1,024$.

The point cloud is augmented following [21] that employs random flipping, random rotation between -5° and 5° , and random scaling with a factor from $0.9 \times$ to $1.1 \times$. We use 20k and 50k points as input for each point cloud on SUN RGB-D and ScanNetV2 datasets, respectively. For ScanNetV2, we use depths to filter out the projected 3D points which should be occluded, but visible due to the sparsity of the point cloud. Since ScanNetV2 does not provide the oriented bounding box annotation, we predict axis-aligned bounding boxes without the rotation angle, as in [15, 21].

Method	RGB	bathub	bed	bookshf	chair	desk	dresser	nightstd	sofa	table	toilet	mAP@0.25	mAP@0.5
VoteNet [21]	✗	75.5	85.6	31.9	77.4	24.8	27.9	58.6	67.4	51.1	90.5	59.1	35.8
MLCVNet [25]	✗	79.2	85.8	31.9	75.8	26.5	31.3	61.5	66.3	50.4	89.1	59.8	-
HGNet [3]	✗	78.0	84.5	35.7	75.2	34.3	37.6	61.7	65.7	51.6	91.1	61.6	-
H3DNet [36]	✗	73.8	85.6	31.0	76.7	29.6	33.4	65.5	66.5	50.8	88.2	60.1	39.0
Group-Free [15]	✗	80.0	87.8	32.5	79.4	32.6	36.0	66.7	70.0	53.8	91.1	63.0	45.2
+3-dim RGB	✓	77.1	87.2	31.2	76.5	30.8	36.2	66.3	68.1	53.0	90.7	61.7	42.0
+Faster R-CNN	✓	78.0	87.4	34.3	77.2	32.8	36.5	67.0	68.8	53.2	91.6	62.7	44.2
+YOLOS	✓	80.6	87.5	35.0	78.5	32.2	37.3	66.7	69.3	54.4	92.1	63.4	45.7
DSS [28]	✓	44.2	78.8	11.9	61.2	20.5	6.4	15.4	53.5	50.3	78.9	42.1	-
2D-driven [12]	✓	43.5	64.5	31.4	48.3	27.9	25.9	41.9	50.4	37.0	80.4	45.1	-
PointFusion [34]	✓	37.3	68.6	37.7	55.1	17.2	23.9	32.3	53.8	31.0	83.8	45.4	-
F-PointNet [22]	✓	43.3	81.1	33.3	64.2	24.7	32.0	58.1	61.1	51.1	90.9	54.0	-
imVoteNet [20]	✓	75.9	87.6	41.3	76.7	28.7	41.4	69.9	70.7	51.1	90.5	63.4	-
Our BrT	✓	82.8	88.0	40.5	79.7	33.4	40.6	67.4	71.1	55.7	93.5	65.4	48.1

Table 1. **3D object detection results on the SUN RGB-D validation set.** Evaluation metrics include the average precisions with 3D IoU threshold 0.25 (mAP@0.25) and threshold 0.5 (mAP@0.5), respectively, as proposed by [27]. Single-class metric adopts mAP@0.25 for evaluation. All listed methods adopt the geometric information (depth or point cloud), and a part of them additionally utilize RGB as input.

For default experiments, model parameters are initialized to weights of ViT-S/16, pre-trained on ImageNet-1k. Yet, the MLP heads for classification and bounding box regression and object queries are generated by Eq. (8) and Eq. (9). We set the hidden size to 384, the layer depth to 12, the patch size to 16, and the number of attention heads to 6. Images are resized to 530×730 , with $N_{\text{pat}} = \lfloor 530/16 \rfloor \times \lfloor 730/16 \rfloor = 1,485$. Besides, there is $K=256$.

For the SUN RGB-D dataset, we include an additional orientation prediction branch to predict the orientation of the 3D box, which additionally includes a classification task and an offset regression task with loss weights of 0.1 and 0.04, respectively following [15].

We train BrT with the AdamW [16] optimizer ($\beta_1=0.9$, $\beta_2=0.999$) with 600 epochs. The learning rate is initialized to 0.004 and decayed by $0.1 \times$ at the 420-th epoch, the 480-th epoch, and the 540-th epoch. We set the loss weights in Eq. (2) to $\alpha_1 = 0.2$, $\alpha_2 = 0.5$ and $\alpha_3 = 0.1$. Following [15, 21], $\mathcal{L}_{\text{obj}}^{3D}$ consists of a center offset sub-loss and a size offset sub-loss with weights 1.0 and 0.1, respectively; $\mathcal{L}_{\text{cls}}^{3D}$ consists of an object classification sub-loss and a size classification sub-loss with equal weights 1.0. In addition, following [8], $\mathcal{L}_{\text{obj}}^{2D}$ consists of a center offset sub-loss and a GIoU sub-loss with weights 1.0 and 2.0, respectively.

4.2. Comparison with state-of-the-art methods

In this section, we compare our BrT with state-of-the-art methods (including using RGB-depth, RGB-point cloud, or point cloud only) for 3D object detection on both SUN RGB-D and ScanNetV2 datasets.

Results on SUN RGB-D. In Table 1, we provide detailed per-class 3D object detection results on SUN RGB-D. We observe that our BrT achieves new records which are remarkably superior to previous methods in terms of mAP@0.25 and mAP@0.5. Specifically, BrT surpasses Group-Free [15], which is a Transformer-based model with only point clouds as input, by 2.3% (mAP@0.25) and 2.9%

(mAP@0.5); and surpasses imVoteNet [20], which is the current best CNN-based model also using RGB, by 1.9% (mAP@0.25). Note that Group-Free adopts a multi-stage ensemble over all Transformer stages to boost the performance, while our BrT uses one output for evaluation.

Since Group-Free [15] achieves the best performance among the methods that only use the geometric information, we experiment with three additional Group-Free variants: (1) “+3-dim RGB” directly appends the the 3-dimensional RGB values to the point cloud features (of the seed points sampled by PointNet++); (2) “+Faster R-CNN” adopts a pre-trained Faster R-CNN [26] (same model used in [20]), a CNN-based 2D detector, to extract region features and concatenate them to the seed points inside that 2D box frustum; (3) “+YOLOS” adopts a pre-trained YOLOS [8] (with DeiT-S [30] model), a Transformer-based 2D detector, to extract image patch features and project them by our bridging method in Sec. 3.4.

Results of these three variants (with RGB) are also provided in Table 1, where only “+YOLOS” achieved marginal gain. This result indicates that intuitive integration of RGB information is difficult to boost the performance. Besides, simply appending 3-dimensional RGB to the point features even impacts the performance, which we conjecture is owing to the discrepancy of 2D/3D representations. By comparison, our BrT is notably better than these three variants.

Results on ScanNetV2. Table 2 provides performance comparison on ScanNetv2. Similarly, we also conduct three additional experiments including “+3-dim RGB”, “+Faster R-CNN”, and “+YOLOS” as competitive baselines based on the state-of-the-art method Group-Free (24-L) [15]. We observe that these three experiments also fail to bring noticeable improvements. We experiment with two architectures for our BrT which adopt the structural designs of ViT-S/16 and ViT-B/16, respectively. The light architecture BrT (ViT-S/16) already surpasses all compared methods. BrT (ViT-B/16) further obtains an additional gain of 1.6 mAP.

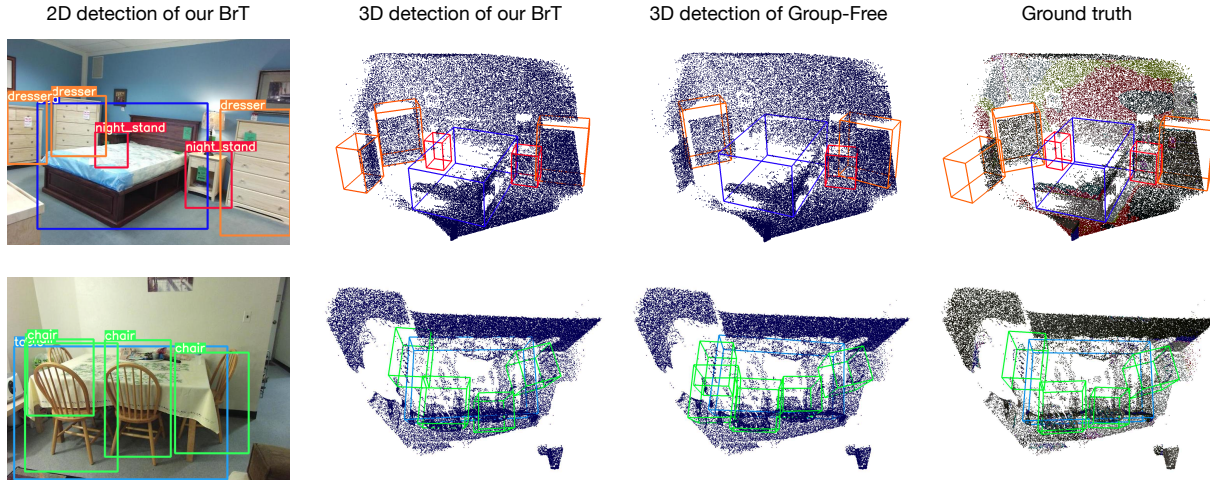


Figure 2. Visualization comparisons to show the advantage of using image information. We compare our BrT with Group-Free [15] which adopts the point cloud as input. First row: Our BrT detects nightstand even the points are very sparse with the help of image detection. Second row: Group-Free is affected by the noises and detects a false positive chair, while our BrT is not affected.

Method	Backbone	RGB	mAP@0.25	mAP@0.5
VoteNet [21]	PointNet++	×	62.9	39.9
MLCVNet [25]	PointNet++	×	64.5	41.4
H3DNet [36]	PointNet++	×	64.4	43.4
H3DNet [36]	4×PointNet++	×	67.2	48.1
HGNet [3]	GU-net	×	61.3	34.4
GSDN [10]	MinkNet	×	62.8	34.8
3D-MPA [7]	MinkNet	×	64.2	49.2
Group-Free [15] (12-L)	PointNet++	×	67.3	48.9
Group-Free [15] (24-L)	PointNet++w2×	×	69.1	52.8
+3-dim RGB	PointNet++w2×	✓	67.8	51.0
+Faster R-CNN	PointNet++w2×	✓	68.7	52.2
+YOLOS	PointNet++w2×	✓	69.2	52.6
Our BrT (ViT-S/16)	PointNet++	✓	69.7	53.0
Our BrT (ViT-B/16)	PointNet++w2×	✓	71.3	55.2

Table 2. **3D object detection results on the ScanNetV2 validation set.** All listed methods adopt the geometric information from point clouds. L denotes the number of self/cross-attention layers. PointNet++w2× expands the backbone width by 2 times.

4.3. Qualitative results and discussion

In Fig. 2, we compare the state-of-the-art method Group-Free [15] (with only point cloud input) with our BrT (with additionally image input). We observe that with the help of image clues, our BrT identifies the nightstand which is partially behind the bed acquiring few points, while Group-Free fails to detect it. In addition, in the second case, Group-Free detects a false positive chair due to the point noises, and our BrT seems to be robust to point noises thanks to the image information. We provide more visualizations on ScanNetv2 with multi-view images in our Appendix.

To highlight the advantage of our bridging techniques proposed in Sec. 3.3 and Sec. 3.4. In Fig. 3, we illustrate visualization comparisons without or with 3D&2D brid-

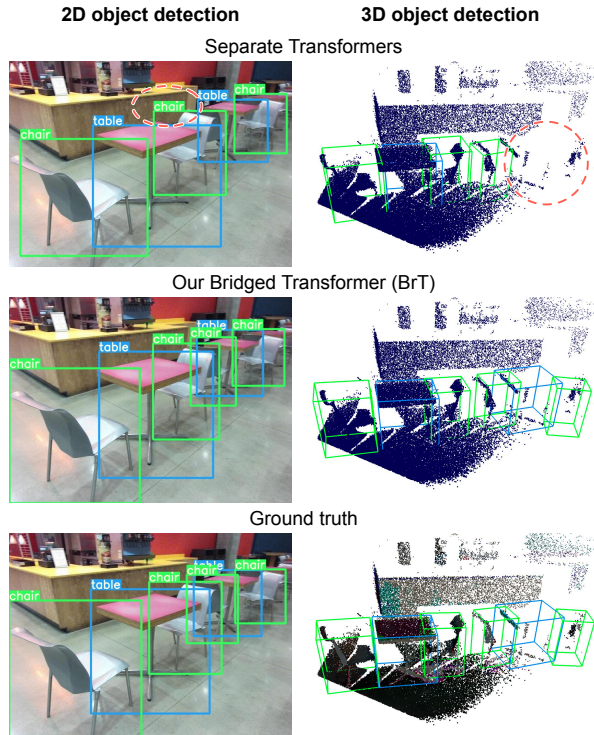


Figure 3. Visualization comparisons of using separate Transformers and our BrT. In the first row, some less obvious objects are missing (one of the occluded 2D chairs, and also the table/chair with sparse 3D points) which are highlighted with red circles. Our BrT successfully captures these difficult objects.

ing. Specifically, in the first row, we adopt two separate Transformers to learn 2D objects and 3D objects respectively, where the learning processes of both Transformers are isolated without any bridging techniques. Highlighted

Conditional object queries	Shared PEs of object queries	point-to-patch projection	mAP@0.25	mAP@0.5
✗	✗	✗	62.6	44.3
✓	✗	✗	62.9	44.7
✗	✓	✗	62.4	44.0
✓	✓	✗	64.0	46.3
✗	✗	✓	64.2	46.6
✓	✓	✓	65.4	48.1

Table 3. Ablation studies of our bridging techniques, including results when using/not using conditional object queries, shared PEs, and the point-to-patch projection proposed in Sec. 3.3 and Sec. 3.4. Experiments are performed on SUN RGB-D.

Connects. between $\mathbf{o}_{\text{pnt}}^l$ and $\mathbf{p}_{\text{pat}}^l$	Connects. between $\mathbf{o}_{\text{pat}}^l$ and $\mathbf{p}_{\text{pnt}}^l$	Connects. between $\mathbf{p}_{\text{pnt}}^l$ and $\mathbf{p}_{\text{pat}}^l$	mAP@0.25	mAP@0.5
✗	✗	✗	62.4	44.3
✓	✗	✗	64.5	47.0
✗	✓	✗	62.8	45.0
✓	✓	✓	63.9	46.4
✓	✓	✗	65.4	48.1

Table 4. Comparisons of different attention connections (connects.) in BrT with our default version (last row) proposed in Sec. 3.2. Experiments are performed on SUN RGB-D.

Pre-trained models	No pre-train	ViT-S	YOLOS-S
mAP results	62.8 (44.5)	65.4 (48.1)	63.4 (45.3)

Table 5. Results of mAP@0.25 and mAP@0.5 (in brackets) when using different pre-training strategies. Experiments are performed on SUN RGB-D.

with red circles, the result of 2D detection fails to identify the two overlapped chairs, and instead, predicts them as one single chair. In addition, its 3D detection result ignores a table and a chair due to the sparsity of points. This is a common bottleneck of 3D object detection, as points are usually unevenly distributed resulting in many sparse regions. Conversely, our BrT in the second row successfully predicts both overlapped chairs correctly in 2D detection, which is probably a consequence of the disambiguation learned by BrT from the 3D space. In addition, with the help of the image, BrT also captures all tables and chairs in 3D detection, even in the regions where points are extremely sparse.

4.4. Ablation analysis

In this section, we perform ablation studies to isolate the benefit of each important component in BrT.

Bridging technique. In Sec. 3.3 and Sec. 3.4, we propose to use conditional object queries, shared PEs of object queries, and the point-to-patch projection. We provide detailed comparison results in Table 3 to verify the necessity of each component. We observe that the combination of conditional object queries and shared PEs brings a significant performance uplift. In addition, the point-to-patch projection itself further boosts the performance. These results indicate that all of the three components are necessary.

Attention connection. As proposed in Sec. 3.2, for each layer l , we block the direct attention connections between the point tokens $\mathbf{p}_{\text{pnt}}^l$ and image patch tokens $\mathbf{p}_{\text{pat}}^l$, but we adopt globally attentive object queries $\mathbf{o}_{\text{pnt}}^l$ and $\mathbf{o}_{\text{pat}}^l$ to build relations of points and image patches. To demonstrate the advantage of this design, in Table 4, we compare different kinds of connections. We observe that using globally attentive object queries indeed help improve the performance, and our default connection setting (last row) achieves the best performance. Connections between $\mathbf{p}_{\text{pnt}}^l$ and $\mathbf{p}_{\text{pat}}^l$ lead to performance drops while increase computation costs, which is the reason that we block these connections.

Pre-training strategy. As described in Sec. 4.1, our BrT is initialized with ViT-S parameters pre-trained on ImageNet-1k. To verify the transferability of using such pre-training strategy, in Table 5, we first compare our results with training from scratch. By comparison, using ViT-S pre-training achieves 2.6 higher mAP@0.25 than training from scratch, which demonstrates the good transfer learning capability of our model, encouraging applying our architecture on larger datasets. In addition, we also try to initialize BrT with off-the-shelf parameters of YOLOS, which is first pre-trained on ImageNet-1k and then finetuned on the COCO dataset [13]. Yet we do not observe further improvements when using COCO finetuned parameters, and we speculate that parameters finetuned on COCO may be easy to overfit the SUN RGB-D dataset.

5. Conclusion

This work presents BrT for 3D object detection that is able to exploit the point clouds and images effectively. Due to the large discrepancy between the representations of 2D and 3D spaces, point tokens and image patch tokens in BrT are bridged with globally attentive object queries, instead of directly connected. In addition, we propose to bridge the 2D and 3D learning processes from two perspectives including leveraging conditional object queries and the point-to-patch projection. Our BrT surpasses state-of-the-art methods on SUN RGB-D and ScanNetV2 datasets, which also demonstrate the scalability of BrT in multi-view scenarios.

Acknowledgement

This work is funded by Major Project of the New Generation of Artificial Intelligence (No. 2018AAA0102900), the Sino-German Collaborative Research Project Crossmodal Learning (NSFC 62061136001/DFG TRR169), Major Science and Technology Innovation 2030 “New Generation Artificial Intelligence” key project (No. 2021ZD0111700), and National Natural Science Foundation of China (No. 62006137). We acknowledge the support of MindSpore, CANN and Ascend AI Processor.

References

- [1] Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. In *IEEE Trans. PAMI*, 2019. 1
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 1, 2, 4
- [3] Jintai Chen, Biwen Lei, Qingyu Song, Haochao Ying, Danny Z Chen, and Jian Wu. A hierarchical graph network for 3d object detection on point clouds. In *CVPR*, 2020. 1, 2, 6, 7
- [4] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *CVPR*, 2017. 2
- [5] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas A. Funkhouser, and Matthias Nießner. ScanNet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017. 5
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 1, 2
- [7] Francis Engelmann, Martin Bokeloh, Alireza Fathi, Bastian Leibe, and Matthias Nießner. 3d-mpa: Multi-proposal aggregation for 3d semantic instance segmentation. In *CVPR*, 2020. 7
- [8] Yuxin Fang, Bencheng Liao, Xinggang Wang, Jiemin Fang, Jiyang Qi, Rui Wu, Jianwei Niu, and Wenyu Liu. You only look at one sequence: Rethinking transformer in vision through object detection. *arXiv preprint arXiv:2106.00666*, 2021. 1, 2, 6
- [9] Yulan Guo, Hanyun Wang, Qingyong Hu, Hao Liu, Li Liu, and Mohammed Bennamoun. Deep learning for 3d point clouds: A survey. *T-PAMI*, 2020. 2
- [10] JunYoung Gwak, Christopher Choy, and Silvio Savarese. Generative sparse detection networks for 3d single-shot object detection. In *arXiv preprint arXiv:2006.12356*, 2020. 7
- [11] Iasonas Kokkinos, Michael M Bronstein, Roei Litman, and Alex M Bronstein. Intrinsic shape context descriptors for deformable shapes. In *CVPR*, 2012. 1
- [12] Jean Lahoud and Bernard Ghanem. 2d-driven 3d object detection in rgb-d images. In *CVPR*, 2017. 6
- [13] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 8
- [14] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030*, 2021. 1, 2
- [15] Ze Liu, Zheng Zhang, Yue Cao, Han Hu, and Xin Tong. Group-free 3d object detection via transformers. *arXiv preprint arXiv:2104.00678*, 2021. 1, 2, 5, 6, 7
- [16] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [17] Jiageng Mao, Yujing Xue, Minzhe Niu, Haoyue Bai, Jiashi Feng, Xiaodan Liang, Hang Xu, and Chunjing Xu. Voxel transformer for 3d object detection. In *ICCV*, 2021. 2
- [18] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *ICCV*, 2021. 2
- [19] Ishan Misra, Rohit Girdhar, and Armand Joulin. An end-to-end transformer model for 3d object detection. In *ICCV*, 2021. 1, 2
- [20] Charles R Qi, Xinlei Chen, Or Litany, and Leonidas J Guibas. Imvotenet: Boosting 3d object detection in point clouds with image votes. In *CVPR*, 2020. 1, 2, 5, 6
- [21] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. In *ICCV*, 2019. 1, 2, 3, 5, 6, 7
- [22] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *CVPR*, 2018. 6
- [23] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, 2017. 1
- [24] Charles R Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NIPS*, 2017. 5
- [25] Xie Qian, Lai Yu-kun, Wu Jing, Wang Zhoutao, Zhang Yiming, Xu Kai, and Wang Jun. Mlcvnet: Multi-level context votenet for 3d object detection. In *CVPR*, 2020. 1, 2, 6, 7
- [26] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *NIPS*, 2015. 6
- [27] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *CVPR*, 2015. 5, 6
- [28] Shuran Song and Jianxiong Xiao. Deep sliding shapes for amodal 3d object detection in rgb-d images. In *CVPR*, 2016. 6
- [29] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *ICCV*, 2015. 1
- [30] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers and distillation through attention. In *ICML*, 2021. 6
- [31] Yikai Wang, Wenbing Huang, Fuchun Sun, Tingyang Xu, Yu Rong, and Junzhou Huang. Deep multimodal fusion by channel exchanging. In *NeurIPS*, 2020. 1
- [32] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *CVPR*, 2015. 1
- [33] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *NeurIPS*, 2021. 1

- [34] Danfei Xu, Dragomir Anguelov, and Ashesh Jain. Pointfusion: Deep sensor fusion for 3d bounding box estimation. In *CVPR*, 2018. 2, 6
- [35] Yuhui Yuan, Rao Fu, Lang Huang, Weihong Lin, Chao Zhang, Xilin Chen, and Jingdong Wang. Hrformer: High-resolution transformer for dense prediction. *arXiv preprint arXiv:2110.09408*, 2021. 2
- [36] Zaiwei Zhang, Bo Sun, Haitao Yang, and Qixing Huang. H3dnet: 3d object detection using hybrid geometric primitives. *arXiv preprint arXiv:2006.05682*, 2020. 2, 6, 7
- [37] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip H.S. Torr, and Li Zhang. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *CVPR*, 2021. 1
- [38] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *CVPR*, 2018. 2
- [39] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 1, 2