

CAT-Det: Contrastively Augmented Transformer for Multi-modal 3D Object Detection

Yanan Zhang^{1,2}, Jiaxin Chen², Di Huang^{1,2*}

¹State Key Laboratory of Software Development Environment, Beihang University, Beijing, China

²School of Computer Science and Engineering, Beihang University, Beijing, China

{zhangyanan, jiaxinchen, dhuang}@buaa.edu.cn

Abstract

In autonomous driving, LiDAR point-clouds and RGB images are two major data modalities with complementary cues for 3D object detection. However, it is quite difficult to sufficiently use them, due to large inter-modal discrepancies. To address this issue, we propose a novel framework, namely Contrastively Augmented Transformer for multi-modal 3D object Detection (CAT-Det). Specifically, CAT-Det adopts a two-stream structure consisting of a Pointformer (PT) branch, an Imageformer (IT) branch along with a Cross-Modal Transformer (CMT) module. PT, IT and CMT jointly encode intra-modal and inter-modal long-range contexts for representing an object, thus fully exploring multi-modal information for detection. Furthermore, we propose an effective One-way Multi-modal Data Augmentation (OMDA) approach via hierarchical contrastive learning at both the point and object levels, significantly improving the accuracy only by augmenting point-clouds, which is free from complex generation of paired samples of the two modalities. Extensive experiments on the KITTI benchmark show that CAT-Det achieves a new state-of-the-art, highlighting its effectiveness.

1. Introduction

3D object detection is a fundamental step in autonomous driving perception systems. It mainly operates 3D point-clouds acquired by LiDAR sensors and provides important spatial clues including location, direction, and object size. Despite true and accurate geometry information recorded, the distribution of point-clouds is disordered, irregular, and sparse, making 3D object detection a challenging task.

The past few years have witnessed the fast development in 3D object detection. A large number of methods are introduced in the literature, and according to the input form in detection feature learning, the methods are roughly catego-

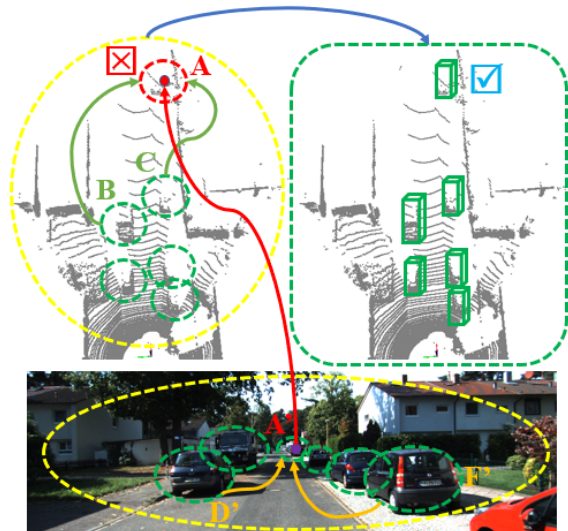


Figure 1. Illustration of the fusion process in CAT-Det. (A): A failure case due to few points at a long distance in the point-cloud modality. (A’): Its corresponding case in the image modality. Although the feature of (A) is enhanced by those of (B) and (C) in the PT branch, this is often insufficient. With the CMT module, the feature of (A) is further enhanced by that of (A’) which also integrates the contributions from (D’) and (F’) in the IT branch, and accurate detection is finally achieved.

rized into grid-based and point-based. The former initially converts point-clouds into regular grids by projecting them to images of specific views [11, 22] or subdividing them to voxels in the space [19, 43, 50, 60] and further conducts 2D or 3D Convolutional Neural Networks (CNN) to encode geometric cues. The latter directly takes raw points and applies point-cloud deep learning networks *e.g.* PointNet [34]/PointNet++ [35] or graph neural networks *e.g.* DGCNN [45] to capture shape structures [38, 40, 51, 56]. More recently, several attempts [15, 37, 52] deliver stronger models by integrating point-based and grid-based networks as hybrid representations, reporting better results.

To boost the performance of 3D object detection, another

*indicates the corresponding author.

strategy targets on multi-modal solutions, which makes use of 3D point-clouds along with 2D images. Although images have not yet proved so competent as an independent modality for this issue evidenced by inferior baselines [2,5,21,30], the combination of geometric and textural clues conveyed in point-clouds and images does lead to accuracy gains for their natural complementarity [7, 18, 33, 41, 53]. F-PointNet [33] and F-ConvNet [46] perform the fusion in series, where 3D frustum proposals are firstly cropped based on prepared 2D regions through a standard 2D CNN detector and each point within the proposal is then segmented and screened using a PointNet-like block for regression. By contrast, more studies fulfill this task in parallel. For example, [41, 48] conduct data-level fusion by enhancing 3D coordinates with point-wise 2D segmentation features; [7, 17, 18, 22, 23] achieve feature-level fusion of 2D and 3D representations from individual networks by simple concatenation or specific modules; and [32] implements box-level fusion which merges the individual candidate sets of a couple of 2D and 3D detectors in a learning manner. Different from the LiDAR only methods that continuously update for more sophisticatedly designed models and more suitable training schemes in the single point-cloud modality, the multi-modal alternatives endeavour to leverage more diverse information and suggest great potential. However, as the KITTI [12] leaderboard displays, there still exists a certain gap between the multi-modal methods and the top LiDAR only ones [59].

Such a gap is due to three aspects. (1) In multi-modal 3D object detection, PointNet++ [35]/3D sparse convolutions [50] and 2D CNNs are principal building blocks to extract point-cloud and image features respectively. Limited by their local receptive fields, contexts cannot be comprehensively acquired from both the modalities, triggering information loss. (2) The widely adopted fusion schemes, particularly the ones at the feature-level, such as direct concatenation [7, 18], additional convolution [22, 23], and simple attention [17, 53], assign no weights or coarse weights learned within limited receptive fields to different features, where crucial clues are not well highlighted. (3) Ground-truth data augmentation [50] is a common practice to facilitate LiDAR only methods; unfortunately, it is not so straightforward to apply this mechanism to multi-modal methods as augmentation in the single modality tends to cause semantic misalignment. [42] indeed presents a cross-modal augmentation technique for paired data, but the procedure on images is cumbersome and easy to incur noise.

To address the issues mentioned above, this paper proposes a novel framework for multi-modal 3D object detection, namely Contrastively Augmented Transformer Detector (CAT-Det). It adopts a two-stream structure, consisting of a Pointformer (PT) branch, an Imageformer (IT) branch together with a Cross-Modal Transformer (CMT)

module. Unlike PointNet++ and CNNs, both the PT and IT branches possess large receptive fields, which are able to respectively capture rich global context information in point-clouds and images to strengthen features of hard samples. Subsequently, the CMT module conducts cross-modal feature interaction and multi-modal feature combination, where essential cues extracted in the two modalities are sufficiently emphasized with holistically learned fine-grained weights. The integration of PT, IT, and CMT fully encodes intra-modal and inter-modal long-range dependencies as a powerful representation, thus benefiting detection performance. In addition, we propose a one-way multi-modal data augmentation (OMDA) approach through hierarchical contrastive learning, which accomplishes effective augmentation by solely performing on the point-cloud modality.

In summary, the major contributions of this paper are:

- (1) We propose a novel CAT-Det framework for multi-modal 3D object detection, with a pointformer branch, an imageformer branch and a cross-modal transformer module. To the best of our knowledge, it is the first attempt that applies the transformer structure to the given task.
- (2) We propose a one-way data augmentation approach for multi-modal 3D object detection via hierarchical contrastive learning, significantly improving the accuracy only by augmenting point-clouds, thus free from complex generation of paired samples of the two modalities.
- (3) We achieve a newly state-of-the-art mAP of all the three classes on the KITTI test set in comparison to the published counterparts, and demonstrate its advantage in detecting hard objects.

2. Related Work

Image based 3D Object Detector. Some approaches [2, 5, 6] perform 2D/3D matching via exhaustively sampling and scoring 3D proposals as representative templates. Numerous methods [8, 20, 21, 30] directly start with accurate 2D bounding boxes to roughly estimate 3D pose from geometric properties obtained by empirical observation. Another way is to first conduct depth estimation and then resort to existing point-cloud based methods [1, 44, 54]. Although 2D object detection has made remarkable advancements, images are not regarded as a good individual modality to predict 3D objects. For the absence of depth information, monocular image based methods suffer from low precision. Stereo image based methods are able to recover depth information, but it is usually coarse with additional noise.

Point-cloud based 3D Object Detector. Some methods convert point-clouds to regular grids by projecting to planes [11, 22] or subdividing to voxels [19, 43, 50, 60] so that they can be processed by 2D or 3D CNNs for feature learning. More methods take raw unordered and irregular data as input and apply point-cloud deep learning networks, such as PointNet [34] and PointNet++ [35], to encode structure features [38, 51] and a few methods [40, 56]

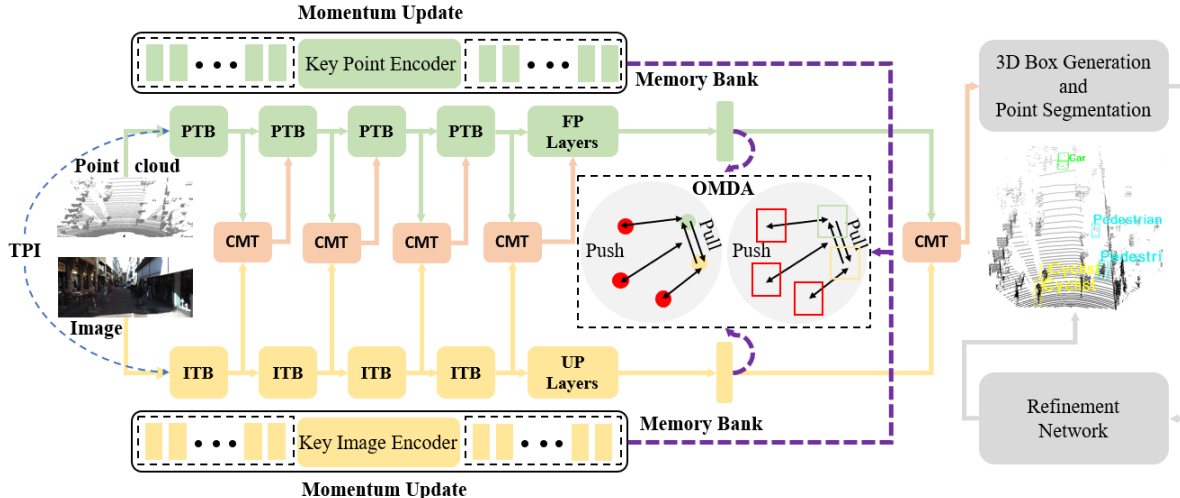


Figure 2. Framework overview. The whole framework consists of three main modules: (1) Two-stream Pointformer and Imageformer (TPI), (2) Cross-Modal Transformer (CMT), and (3) One-way Multi-modal Data Augmentation (OMDA). TPI builds intra-modal long-range context feature representation from the two modalities, and CMT performs inter-modal feature interaction and aggregation at multiple levels. In addition, OMDA conducts hierarchical contrastive learning to achieve concise yet effective data augmentation.

attempt graph neural networks in this step. Recent methods [15, 37, 52] also utilize both point-based and voxel-based networks to extract features from different representations of point-clouds. More recently, a series of methods [13, 26, 28, 29, 31, 36] have emerged based on point transformer for the property to capture global contexts.

Multi-modal 3D Object Detector. MV3D [7] and AVOD [18] take LiDAR projections and RGB images as inputs and fuse region-based features to make prediction. F-PointNet [33] and F-ConvNet [46] first leverage a 2D CNN object detector to extract 2D regions from images, then transform 2D regional coordinates to the 3D space to crop frustum proposals, and finally localize interest points within the frustum by a PointNet-like block for regression. PointPainting [41] and PI-RCNN [48] turn to semantics network for per pixel classification, and the relevant segmentation score, which serves as compact features of the image, is appended to the LiDAR points via projecting them into the segmentation mask. CLOCs [32] directly uses pre-trained 2D and 3D detectors by late fusion, making the proposals in different modalities connected without integrating the features. Recent studies [17, 22, 23, 53] combine the modalities in the feature space to obtain a multi-modal representation before feeding them into a supervised learner. Despite many efforts, to the best of our knowledge, we are the first to investigate the multi-modal transformer network for this task.

3. The Proposed Approach

3.1. Framework Overview

As Fig. 2 illustrates, CAT-Det basically adopts a two-stream structure consisting of a Pointformer (PT) branch

and an Imageformer (IT) branch, separately learning representations of LiDAR point-clouds and RGB images by exploring long-range intra-modal contexts. To complement the learning in each single modality, the Cross-Modal Transformer (CMT) module is employed to perform cross-modal feature interaction, followed by multi-modal feature aggregation with holistically learned fine-grained weights. The combination of PT, IT and CMT constitutes a novel transformer backbone. Meanwhile, a One-way Multi-modal Data Augmentation (OMDA) approach is developed to accomplish efficient data augmentation by a hierarchical contrastive learning at both the point-level and object-level, further facilitating the training of a strong deep transformer network for multi-modal 3D object detection.

3.2. Two-stream Multi-modal Transformer

Existing multi-modal 3D object detectors [17, 53] often adopt PointNet++/Sparse 3D CNN for point-clouds and 2D CNNs for images in representation learning. They mostly suffer from limited receptive fields, thus unable to fully explore global contextual information, which is important to detecting hard examples (*e.g.* tiny objects). Recent work [9, 10] has proved the effectiveness of Transformer in modeling long-range dependencies. Despite its increasing prevalence, the transformer structure has not yet been investigated in multi-modal 3D object detection. This motivates us to make the first attempt in developing a deep multi-modal transformer backbone to capture richer global contexts for 3D object detection.

To this end, we propose a novel multi-modal transformer network, consisting of the two-stream PT and IT branches connected by several CMTs as shown in Fig. 2. Given a

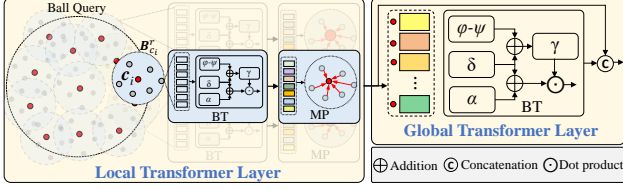


Figure 3. Point Transformer Block. By local and global combination, it captures both the dependencies from the adjacent regions and the whole scene, thus facilitating feature learning for 3D object detection. BT: basic transformer; MP: max-pooling.

paired multi-modal input $\{\mathbf{P}, I\}$, $PT(\cdot)$ and $IT(\cdot)$ learn representations for point-clouds $\mathbf{P} = \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N\} \in \mathbb{R}^{N \times 3}$ and images I , respectively. CMT performs cross-modal interaction and multi-modal aggregation at different levels.

Pointformer. Despite that a few recent attempts have investigated transformers for point-clouds [14, 57], most of them are specifically designed for classification, which only adopt a local transformer structure. However, global context information is vital to 3D detection, which cannot be fully recorded by local transformers. To address this issue, we propose a novel pointformer, composed of multiple stacked Point Transformer Blocks (PTB). As displayed in Fig. 3, PTB consists of a local transformer layer and a global transformer layer. The local layer explores geometry structures of points within a neighborhood, and the global one encodes the holistic context at the scene level. By combining them, PTB captures the context information from the points of the nearby local regions as well as the full scene.

To be specific, the local transformer layer first applies the furthest point sampling on input point-clouds \mathbf{P} to choose a subset $\mathbf{C} = \{c_1, c_2, \dots, c_{N'}\} \subset \mathbf{P}$. Then, we conduct the ball query operation by taking each point c_i as the centroid, where K points $\mathbf{B}_{c_i}^r$ are selected within a ball centered at c_i with a radius r . The subset $\mathbf{B}_{c_i}^r$ is further grouped and fed into a basic transformer block $BT(\cdot)$ for local information aggregation, which adopts the structure based on self-attention inspired by [57]. Given the input $\mathbf{B}_{c_i}^r$, the output $\mathbf{y}_i = BT(\mathbf{B}_{c_i}^r)$ is formulated as below:

$$\mathbf{y}_i = \sum_{\mathbf{p}_j \in \mathbf{B}_{c_i}^r} \rho(\gamma(\varphi(\mathbf{p}_i) - \psi(\mathbf{p}_j) + \delta)) \odot (\alpha(\mathbf{p}_j) + \delta), \quad (1)$$

where the outputs of δ , γ and α are the encoded position, the self-attention and the transformed value, respectively. φ , ψ , and α are pointwise feature transformations, such as linear projections or MLPs. δ is a position encoding function defined as $\delta = \theta(\mathbf{p}_i - \mathbf{p}_j)$. Both the mapping function γ and the encoding function θ are MLPs with two linear layers and one ReLU nonlinearity. ρ is the Softmax function.

Albeit the exploration of long-range dependency in $\mathbf{B}_{c_i}^r$ by $BT(\cdot)$, the local transformer layer processes the point-

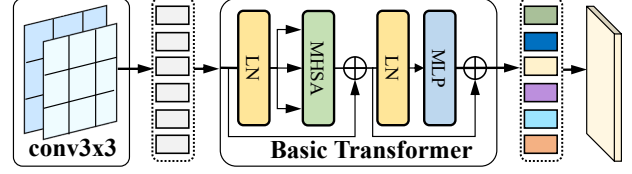


Figure 4. Image Transformer Block, which is a combination of a few convolutional layers and transformers.

cloud locally, unable to present the holistic context information. Therefore, we additionally employ a global transformer layer, which has a similar transformer structure as the local one, but takes all points \mathbf{C} as input, instead of a local subset $\mathbf{B}_{c_i}^r$. The features generated by the local and global transformer layer are concatenated to integrate both local and global contexts.

Similar to PointNet++, we adopt a Feature Propagation (FP) layer after stacked PTBs for up-sampling.

Imageformer. Vision Transformer (ViT) [10] is the first work that adopts the transformer network in the visual domain, which employs the self-attention mechanism to build holistic dependencies among visual tokens. Since raw image patches are taken as tokens, it fails to encode local visual spatial information. Some recent studies [47, 55] handle this problem by adding a few convolutional layers before the transformer layer, which is used as the basic transformer in our work. To align with Pointformer, we adopt similar structures by stacking several Image Transformer Blocks (ITB) as shown in Fig. 4. Each ITB consists of two convolutional layers for local visual context encoding, and a successive basic multi-head transformer encoder [10] for global context information exploration. Finally, ITB reshapes the transformed vector sequence into a 2D feature map for further processing. Following the stacked ITBs, an up-sampling (UP) layer is employed to recover the image resolution, generating feature maps with the same size as the original image.

Cross-Modal Transformer. PTB and ITB extensively explore contexts in the point-cloud \mathbf{P} and the image I , respectively. However, as in Fig. 1, the context in a single modality is probably incomplete due to noise, which can be complemented by that conveyed in the other modality. This motivates us to propose a module between PTB and ITB, to perform cross-modal information interaction and multi-modal feature aggregation.

Suppose the features from PTB and ITB are $\mathbf{F}_{\mathbf{P}}$ and \mathbf{F}_I respectively, where $\mathbf{F}_{\mathbf{P}}$ are representations of a set of down-sampled points $\hat{\mathbf{P}} \subset \mathbf{P}$. For each point $\mathbf{p} \in \mathbf{P}$, we project it to the corresponding pixel coordinate \mathbf{p}'_I in I by a function $f_{proj}(\cdot)$. For instance, in KITTI, $f_{proj}(\cdot)$ is formulated as:

$$\mathbf{p}'_I = f_{proj}(\mathbf{p}) = C_{rect} \cdot R_{rect} \cdot T_{cam \leftarrow LiDAR} \cdot \mathbf{p}, \quad (2)$$

where $T_{cam \leftarrow LiDAR}$ is the transformation matrix from the

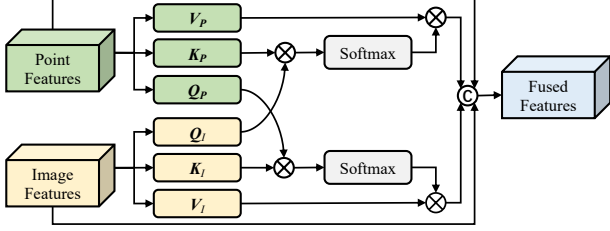


Figure 5. Cross-Modal Transformer. It performs complementary feature enhancement by adaptively learning weights for different modalities through cross transformer.

coordinate of LiDAR to camera, R_{rect} and C_{rect} are the rectifying rotation and the calibration matrices of the camera, respectively. Based on $f_{proj}(\cdot)$, we convert the 3D coordinates \hat{P} to 2D pixels $\hat{P}'_I = f_{proj}(\hat{P})$, based on which we select the features from F_I at positions \hat{P}'_I and fetch a subset F'_I spatially aligned with F_P . In other words, F_P and F'_I are point and imagery features for \hat{P} , respectively.

Subsequently, CMT projects F_P into the query $Q_P = F_P \cdot W_Q$, the key $K_P = F_P \cdot W_K$ and the value $V_P = F_P \cdot W_V$, where W_Q , W_K and W_V are learnable linear mappings. Similarly, the imagery feature F'_I is projected to Q_I , K_I and V_I . The context from the image modality is thus explored by attentive weights $A_{P \leftarrow I} = \text{Softmax}(Q_I K_P^T)$ and encoded into point features by $F_P^{cont} = A_{P \leftarrow I} \odot V_P$. Similarly, the context from the point modality can be explored and encoded into imagery features by $F_I^{cont} = \text{Softmax}(Q_P K_I^T) \odot V_I$.

The original multi-modal features F_P/F_I and the features F_P^{cont}/F_I^{cont} with cross-modal interactions are aggregated by $F_P := F_P \oplus F_P^{cont} \oplus F_I^{cont}$ as new point features, where \oplus stands for concatenation.

3.3. One-way Multi-modal Data Augmentation

Data augmentation has proved effective for object detection, which however is mostly applied within a single modality and rarely considered in the multi-modal scenario. Due to the heterogeneity between point-clouds and images, it is generally difficult to synchronize the augmenting operations across modalities, leading to severe cross-modal misalignment. Recently, [42] presents a complex approach to generate paired data, but the pipeline on images is cumbersome and easy to incur noise. Instead, we propose a novel One-way Multi-modal Data Augmentation (OMDA) approach, which performs augmentation on point-clouds only, and efficiently extends it to multiple modalities by contrastive learning.

The basic idea behind OMDA is two-fold: (1) High-quality image augmentation is generally much more complex and difficult than that on point-clouds and it is thus expected to only augment LiDAR data and then make a lightweight multi-modal extension. (2) One-way augmentation

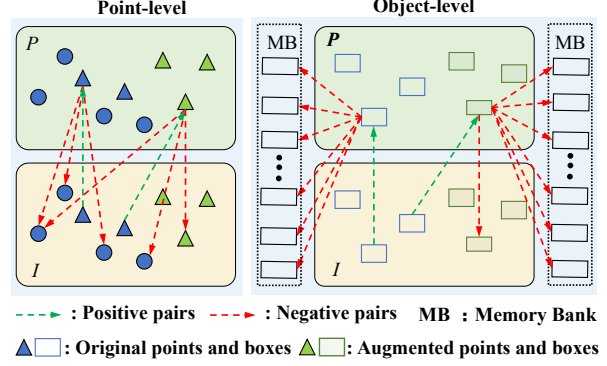


Figure 6. Positive/negative pair selection for contrastive learning.

as in (1) may bring in severe cross-modal misalignment. Inspired by the recent success of contrastive learning in self-supervised models [4, 16] and cross-modal semantic alignment [24, 25], we elaborately design a contrastive learning scheme to address such misalignment across modalities.

Specifically, OMDA adopts GT-Paste [50], widely used for LiDAR-only methods, to augment a given point-cloud by pasting extra 3D objects from other LiDAR frames without spatial collision. As there lacks the images corresponding to the augmented point-clouds, cross-modal data misalignment occurs, probably deteriorating multi-modal interaction (e.g. CMT), which implicitly assumes that point-cloud/image pairs are well-aligned. Thus, we perform contrastive learning among the raw point-clouds P , the corresponding images I , and the augmented point-clouds P_{aug} at both point- and object-level in a hierarchical manner.

Point-level Contrastive Augmentation. To preserve the supervision of raw data pair, we firstly construct cross-modal positive/negative point pairs from (P, I) for contrastive learning. Given a point $p \in P$, we fetch its corresponding 2D pixel coordinate p'_I by $p'_I = f_{proj}(p)$ as in Eq. (2). Since P and I are well aligned, (p, p'_I) indicates the 3D/2D position of the same object, thus naturally forming the positive pair. To construct negative pairs, we select the 3D points $\mathcal{N}_p \subset P$ that belong to different object classes from p , e.g. points that with confidence scores (predicted by the segmentation head) less than a threshold t . The negative pairs are chosen as $\{(p, q'_I) | q'_I \in \mathcal{N}'_{p,I}\}$, where $\mathcal{N}'_{p,I}$ is the set of 2D coordinates of \mathcal{N}_p , i.e. $\mathcal{N}'_{p,I} = f_{proj}(\mathcal{N}_p)$.

Afterwards, we construct positive/negative point pairs from the augmented point-cloud P_{aug} and the unpaired image I as in Fig. 6. Suppose $p_{aug} \in P_{aug}$ is a point from the pasted virtual object O_{vir} . Since GP-paste avoids spatial overlap with existing objects when pasting O_{vir} on P , no objects lie at the position $f_{proj}(p_{aug})$ in I . It means that $(p_{aug}, f_{proj}(p_{aug}))$ is definitely unpaired, thus forming a negative pair. To collect positive pairs w.r.t. p_{aug} , we select the 3D point $\hat{p} \in P$ that most likely belongs to the same class as p_{aug} , e.g. the point with the highest con-

confidence score predicted by the segmentation head. Thus, $(\mathbf{p}_{aug}, f_{proj}(\hat{\mathbf{p}}))$ is chosen as the positive pair.

Based on the selected positive/negative pairs $\mathcal{S}_+/\mathcal{S}_-$, we formulate the following point-level contrastive loss as:

$$\mathcal{L}_{cl-p} = - \sum_{(i,j) \in \mathcal{S}_+} \log \frac{\exp(\mathbf{f}_i^T \cdot \mathbf{f}_j / \tau)}{\sum_{(i,k) \in \mathcal{S}_-} \exp(\mathbf{f}_i^T \cdot \mathbf{f}_k / \tau)}, \quad (3)$$

where \mathbf{f}_i is the feature of the i -th point in \mathbf{P} , \mathbf{f}_j is the feature at location j of image I , and τ is the scaling factor. It can be observed that the correlation between the paired cross-modal features is increased, while being decreased for unpaired features by minimizing \mathcal{L}_{cl-p} , thus alleviating cross-modal misalignment of the augmented data.

Object-level Contrastive Augmentation. Contrastive learning at the point-level offers fine-grained point-wise semantic alignment, while detectors focus on regions. To alleviate regional semantic alignment, we additionally perform object-level contrastive learning on augmented data.

Similar to point-level contrastive learning, we first construct positive/negative pairs of objects across modalities. As in Fig. 6, the object O from the raw point-cloud \mathbf{P} and its paired object O'_I in the aligned image I naturally constitute a positive pair. The pasted virtual object O_{vir} , belonging to the same class as O , also forms a positive pair with O'_I . Due to the class imbalance frequently occurring in labeled data, \mathbf{P} and I probably only contain objects from one class. Although we can simply select background areas to form the negative pair, foreground areas belonging to objects from different classes are more desirable, as they provide stronger supervision. Inspired by [16], we employ a memory bank to generate more precise and discriminative representations for negative pair selection. As in Fig. 2, the memory bank adopts an encoder $E(\cdot)$ with the same structure as the two-stream multi-modal transformer, and maintains two queues of features denoted by \mathcal{Q}_P and \mathcal{Q}_I . \mathcal{Q}_P contains point features of objects from all classes, and \mathcal{Q}_I contains imagery features. For O or O'_{vir} , we select the elements in \mathcal{Q}_P and \mathcal{Q}_I from distinct classes to collect negative pairs.

With positive and negative object pairs $\mathcal{O}_+/\mathcal{O}_-$, object-level contrastive learning is applied by minimizing the loss:

$$\mathcal{L}_{cl-o} = - \sum_{(i,j) \in \mathcal{O}_+} \log \frac{\exp(\mathbf{g}_i^T \cdot \mathbf{g}_j / \tau)}{\sum_{(i,k) \in \mathcal{O}_-} \exp(\mathbf{g}_i^T \cdot \mathbf{g}_k / \tau)}, \quad (4)$$

where \mathbf{g} is the object-level representation by aggregating feature vectors in the object bounding box via max-pooling or is directly fetched from the memory banks $\mathcal{Q}_P/\mathcal{Q}_I$.

As in [16], we employ the momentum update mechanism to optimize the encoder $E(\cdot)$ instead of gradient update, in order to strengthen the stability of features in the memory

bank. Refer to [16] for more details about the optimization on $E(\cdot)$ and memory banks $\mathcal{Q}_P/\mathcal{Q}_I$.

Overall optimization. Besides the contrastive learning losses \mathcal{L}_{cl-p} and \mathcal{L}_{cl-o} , we also utilize the conventional detection losses \mathcal{L}_{rpn} and \mathcal{L}_{rcnn} as in [38]. The total loss for optimizing the overall transformer network is formulated as $\mathcal{L}_{tot} = \mathcal{L}_{rpn} + \mathcal{L}_{rcnn} + \lambda \cdot (\mathcal{L}_{cl-p} + \mathcal{L}_{cl-o})$, where λ is the trade-off parameter, empirically set as 0.15 by default.

4. Experiments

We evaluate CAT-Det on the widely used KITTI benchmark [12], and for fair comparison, we adopt the same protocol as in [7, 38], separating original training data into a training set and a validation set. The Average Precision (AP) is used as the metric, and the IoU thresholds are set to 0.7, 0.5, and 0.5 for car, pedestrian and cyclist, as officially specified. APs are computed by recalling 11 and 40 positions on the val and test splits respectively.

4.1. Implementation Details

Both the point-clouds and images are used in training and testing. The range of point-clouds is constrained to (0, 70.4), (-40, 40) and (-3, 1) along the X, Y and Z axis, respectively, which is further down-sampled to 16,384 points as input, and the resolution of images is 1280×384 . In the PT branch, there are four stacked PTBs, with the numbers of sampled points set to 4,096, 1,024, 256, and 64, respectively, and four FP layers which up-sample the point-cloud back to the original size with a stride of 4. Similarly, in the IT branch, there are four cascaded ITBs followed by four UP layers for parallel transposed convolutions with strides 2, 4, 8, and 16. In our memory bank, the sample size of each category is 1,024. We adopt the ADAM optimizer and the cosine annealing learning rate schedule with an initial value at 0.002. The batch size and the maximal number of learning epochs are set to 16 and 80, respectively. All the experiments are conducted on 8 GTX 1080Ti GPUs.

4.2. Comparison with State-of-the-Arts

We compare CAT-Det to the following categories of methods, including (1) LiDAR-only with non-transformer structures [3, 15, 19, 27, 37–39, 50–52, 59, 60]; (2) LiDAR-only with transformer structures [13, 31]; (3) Multi-modal (LiDAR+RGB) [7, 17, 18, 22, 23, 32, 33, 41, 46, 48, 49, 53, 58].

Table 1 summarizes the official results on the test set. As Table 1 displays, the LiDAR-only methods outperform the existing multi-modal counterparts at most cases, indicating that modeling multi-modal data indeed remains a challenging task, despite more information available. PointTransformer and M3DETR both adopt the transformer structure; however, their performance is not as good as that of the non-transformer counterparts such as [51] and [3]. By virtue of

Method	Modality	Car (%)				Pedestrian (%)				Cyclist (%)				mAP (%)
		Easy	Mod.	Hard	mAP	Easy	Mod.	Hard	mAP	Easy	Mod.	Hard	mAP	
VoxelNet [60]	L	77.47	65.11	57.73	66.77	39.48	33.69	31.51	34.89	61.22	48.36	44.37	51.32	50.99
PointRCNN [38]	L	86.96	75.64	70.70	77.77	47.98	39.37	36.01	41.12	74.96	58.82	52.53	62.10	60.33
PointPillars [19]	L	82.58	74.31	68.99	75.29	51.45	41.92	38.89	44.09	77.10	58.65	51.92	62.56	60.65
TANet [27]	L	84.39	75.94	68.82	76.38	53.72	44.34	40.49	46.18	75.70	59.44	52.53	62.56	61.71
STD [52]	L	87.95	79.71	75.09	80.92	53.29	42.47	38.35	44.70	78.69	61.59	55.30	65.19	63.60
Part-A ² [39]	L	87.81	78.49	73.51	79.94	53.10	43.35	40.06	45.50	79.17	63.52	56.93	66.54	63.99
PV-RCNN [37]	L	90.25	81.43	76.82	82.83	52.17	43.29	40.29	45.25	78.60	63.71	57.65	66.65	64.91
3DSSD [51]	L	88.36	79.57	74.55	80.83	54.64	44.27	40.23	46.38	82.48	64.10	56.90	67.82	65.01
HotSpotNet [3]	L	87.60	78.31	73.34	79.75	53.10	45.37	41.47	46.65	82.59	65.95	59.00	69.18	65.19
SA-SSD [15]	L	88.75	79.79	74.16	80.90	-	-	-	-	-	-	-	-	-
SE-SSD [59]	L	91.49	82.54	77.15	83.73	-	-	-	-	-	-	-	-	-
PointTransformer [31]	L	87.13	77.06	69.25	77.81	50.67	42.43	39.60	44.23	75.01	59.80	53.99	62.93	61.66
M3DETR [13]	L	90.28	81.73	76.96	82.99	45.70	39.94	37.66	41.10	83.83	66.74	59.03	69.87	64.65
MV3D [7]	L+R	74.97	63.63	54.00	64.20	-	-	-	-	-	-	-	-	-
ContFuse [23]	L+R	83.68	68.78	61.67	71.38	-	-	-	-	-	-	-	-	-
MMF [22]	L+R	88.40	77.43	70.22	78.68	-	-	-	-	-	-	-	-	-
PI-RCNN [48]	L+R	84.37	74.82	70.03	76.41	-	-	-	-	-	-	-	-	-
EPNet [17]	L+R	89.81	79.28	74.59	81.23	-	-	-	-	-	-	-	-	-
3D-CVF [53]	L+R	89.20	80.05	73.11	80.79	-	-	-	-	-	-	-	-	-
CLOCs [32]	L+R	88.94	80.67	77.15	82.25	-	-	-	-	-	-	-	-	-
AVOD-FPN [18]	L+R	83.07	71.76	65.73	73.52	50.46	42.27	39.04	43.92	63.76	50.55	44.93	53.08	56.84
F-PointNet [33]	L+R	82.19	69.79	60.59	70.86	50.53	42.15	38.08	43.59	72.27	56.12	49.01	59.13	57.86
PointPainting [41]	L+R	82.11	71.70	67.08	73.63	50.32	40.97	37.84	43.05	77.63	63.78	55.89	65.77	60.82
F-ConvNet [46]	L+R	87.36	76.39	66.69	76.81	52.16	43.38	38.80	44.78	81.98	65.07	56.54	67.86	63.15
CAT-Det (Ours)	L+R	89.87	81.32	76.68	82.62	54.26	45.44	41.94	47.21	83.68	68.81	61.45	71.31	67.05

Table 1. Comparison with state-of-the-art approaches on the KITTI test split. ‘-’ indicates that either the result is not reported or the source code is not publicly available. ‘L’ and ‘R’ stand for the LiDAR and RGB modalities, respectively. Best in bold.

Method	Modality	Car (%)				Pedestrian (%)				Cyclist (%)				mAP (%)
		Easy	Mod.	Hard	mAP	Easy	Mod.	Hard	mAP	Easy	Mod.	Hard	mAP	
PointPillars [19]	L	86.46	77.28	74.65	79.46	57.75	52.29	47.90	52.65	80.05	62.68	59.70	67.48	66.53
SECOND [50]	L	88.61	78.62	77.22	81.48	56.55	52.98	47.73	52.42	80.58	67.15	63.10	70.28	68.06
3DSSD [51]	L	88.55	78.45	77.30	81.43	58.18	54.31	49.56	54.02	86.25	70.48	65.32	74.02	69.82
PointRCNN [38]	L	88.72	78.61	77.82	81.72	62.72	53.85	50.24	55.60	86.84	71.62	65.59	74.68	70.67
PV-RCNN [37]	L	89.03	83.24	78.59	83.62	63.71	57.37	52.84	57.97	86.06	69.48	64.50	73.35	71.64
Part-A ² [39]	L	89.55	79.40	78.84	82.60	65.68	60.05	55.44	60.39	85.50	69.90	65.48	73.63	72.20
SE-SSD [59]	L	90.21	86.25	79.22	85.23	-	-	-	-	-	-	-	-	-
MV3D [7]	L+R	71.29	62.68	56.56	63.51	-	-	-	-	-	-	-	-	-
3D-CVF [53]	L+R	89.67	79.88	78.47	82.67	-	-	-	-	-	-	-	-	-
AVOD-FPN [18]	L+R	84.41	74.44	68.65	75.83	-	58.80	-	-	-	49.70	-	-	-
PointFusion [49]	L+R	77.92	63.00	53.27	64.73	33.36	28.04	23.38	28.26	49.34	29.42	26.98	35.25	42.75
F-PointNet [33]	L+R	83.76	70.92	63.65	72.78	70.00	61.32	53.59	61.64	77.15	56.49	53.37	62.34	65.58
SIFRNet [58]	L+R	85.62	72.05	64.19	73.95	69.35	60.85	52.95	61.05	80.87	60.34	56.69	65.97	66.99
CLOCs [32]	L+R	89.49	79.31	77.36	82.05	62.88	56.20	50.10	56.39	87.57	67.92	63.67	73.05	70.50
EPNet [17]	L+R	88.76	78.65	78.32	81.91	66.74	59.29	54.82	60.28	83.88	65.50	62.70	70.69	70.96
CAT-Det (Ours)	L+R	90.12	81.46	79.15	83.58	74.08	66.35	58.92	66.45	87.64	72.82	68.20	76.22	75.42

Table 2. Comparison with state-of-the-art approaches on the KITTI validation split. Best in bold.

the specially designed transformer structure and the effective one-way multi-modal data augmentation scheme, our method reaches a newly the-state-of-art score on KITTI. It is worth noting that CAT-Det is the first multi-modal solution that surpasses the LiDAR-only ones with a large margin, *i.e.* a gain of 1.86% over the second best HotSpotNet, suggesting the potential of multi-modal data for further improvement. Besides, CAT-Det ranks the first in many cases, particularly for the challenging classes of pedestrian and cyclist. In these two classes, there are many small or partial instances, and more context information with large receptive fields is thus required, which is fully explored by the

transformer structures of CAT-Det.

We also evaluate the methods on KITTI val. As Table 2 shows, CAT-Det again achieves the best mAP, 3.22% higher than the second best Part-A². The decent results on hard examples visualized in Fig. 7 highlight the advantage of CAT-Det for this issue.

4.3. Ablation Study

On TPI/CMT/OMDA. We individually evaluate the contributions of TPI, CMT and OMDA in CAT-Det. PointNet++ and 2D CNN with concatenation based fusion is selected as the baseline, whose mAP is 68.65%. After replacing Point-

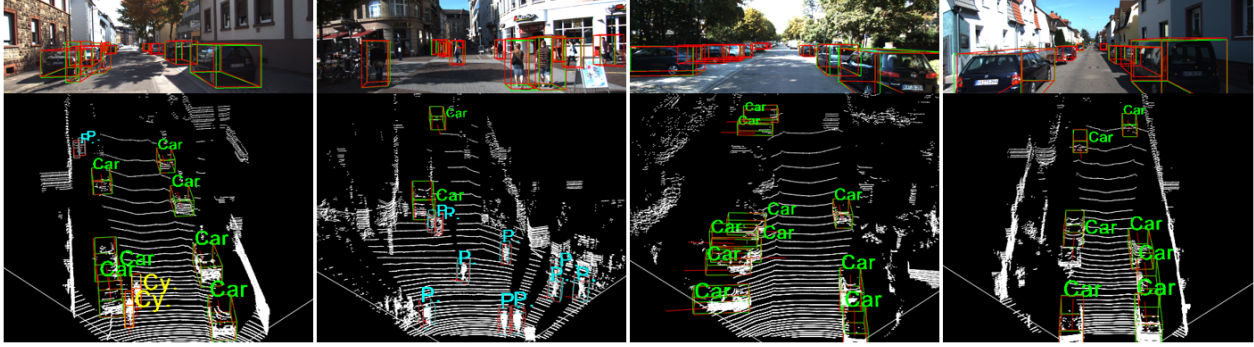


Figure 7. Visualized results by CAT-Det on KITTI val split. Red/green rectangles indicate predicted/GT bounding boxes.

Component			3D Object Detection (%)			
TPI	CMT	OMDA	Car	Ped.	Cyc.	mAP
			80.41	59.27	66.28	68.65
✓			81.12	61.26	68.94	70.44
✓	✓		81.73	63.30	70.63	71.89
✓	✓	✓	83.58	66.45	76.22	75.42

Table 3. Contributions of TPI/CMT/OMDA to CAT-Det.

Module				3D Object Detection (%)			
Concat.	AF	CMT _S	CMT	Car	Ped.	Cyc.	mAP
✓				81.12	61.26	68.94	70.44
	✓			81.31	61.74	69.45	70.83
		✓		81.57	62.68	69.99	71.41
			✓	81.73	63.30	70.63	71.89

Table 4. Results of different fusion schemes in the CMT module.

Net++ and 2D CNN with TPI while maintaining the concatenation operation, the performance is boosted by 1.79%, and it is further improved by 1.45% through incorporating CMT. We attribute this improvement to the specifically designed network structures, which extensively capture intra-modal and cross-modal global contexts. Finally, we apply OMDA and obtain the full model of CAT-Det, which gains 3.53% in mAP, validating its effectiveness.

On Fusion in CMT. To analyze CMT in integrating multi-modal features, the widely used alternatives are applied for comparison, including concatenation (Concat.) [7, 18] and attention-based fusion (AF) [17, 53]. A variant of CMT (denoted as CMT_S) is also considered, where only the cross-modal features are concatenated as $F_P^{cont} \oplus F_I^{cont}$ as the output, without single-modal features F_P and F_I . As in Table 4, CMT clearly outperforms Concat. and AF only with cross-modal features (CMT_S), and further boosts the accuracy by preserving single-modal features. It proves the advantage of CMT in encoding cross-modal global contexts.

On Components in OMDA. To validate point-level and object-level contrastive learning, three versions of OMDA are adopted: (1) point-level augmentation only (CA-P); (2) object-level augmentation with background to select negative pairs (CA-O-BG); (3) object-level augmentation with

Method			3D Object Detection (%)			
CA-P	CA-O-BG	CA-O-MB	Car	Ped.	Cyc.	mAP
			81.73	63.30	70.63	71.89
✓			82.95	66.13	74.55	74.54
✓	✓		83.37	66.27	75.36	75.00
✓		✓	83.58	66.45	76.22	75.42

Table 5. Results of different components in OMDA.

memory bank to select negative pairs (CA-O-MB). As Table 5 displays, a gain of 2.65% is achieved when only applying point-level contrastive augmentation, clearly revealing the superiority of OMDA. The performance is further boosted using a simple version of object-level contrastive augmentation, *i.e.* CA-O-BG, validating the necessity of augmentation at the object-level. By expanding object-level negative sample pairs via memory bank, the score is the best.

5. Conclusion

This paper proposes a novel framework, namely Contrastively Augmented Transformer for multi-modal 3D object Detection (CAT-Det). It aims to solve the problems of insufficient multi-modal fusion and lack of effective multi-modal data augmentation. To this end, we propose a multi-modal transformer to extensively encode the intra-modal and cross-modal long-range context information. Furthermore, we propose a one-way data augmentation via hierarchical contrastive learning, remarkably improving the accuracy only by augmenting point-clouds, thus free from complex generation of paired samples of the two modalities. Extensive experiments are conducted on KITTI, and CAT-Det reaches a newly state-of-the-art.

Acknowledgment

This work is partly supported by the National Natural Science Foundation of China (No. 62022011), the Research Program of State Key Laboratory of Software Development Environment (SKLSDE-2021ZX-04), and the Fundamental Research Funds for the Central Universities.

References

- [1] Wentao Bao, Bin Xu, and Zhenzhong Chen. Monofenet: Monocular 3d object detection with feature enhancement networks. *IEEE TIP*, 29:2753–2765, 2019. 2
- [2] Florian Chabot, Mohamed Chaouch, Jaonary Rabarisoa, Céline Teuliere, and Thierry Chateau. Deep manta: A coarse-to-fine many-task network for joint 2d and 3d vehicle analysis from monocular image. In *CVPR*, pages 2040–2049, 2017. 2
- [3] Qi Chen, Lin Sun, Zhixin Wang, Kui Jia, and Alan Yuille. Object as hotspots: An anchor-free 3d object detection approach via firing of hotspots. In *ECCV*, pages 68–84. Springer, 2020. 6, 7
- [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020. 5
- [5] Xiaozhi Chen, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. Monocular 3d object detection for autonomous driving. In *CVPR*, pages 2147–2156, 2016. 2
- [6] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Andrew G Berneshawi, Huimin Ma, Sanja Fidler, and Raquel Urtasun. 3d object proposals for accurate object class detection. In *NeurIPS*, pages 424–432. Citeseer, 2015. 2
- [7] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *CVPR*, pages 1907–1915, 2017. 2, 3, 6, 7, 8
- [8] Yilun Chen, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Dsgn: Deep stereo geometry network for 3d object detection. In *CVPR*, pages 12536–12545, 2020. 2
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3
- [10] 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. In *ICLR*, 2021. 3, 4
- [11] Martin Engelcke, Dushyant Rao, Dominic Zeng Wang, Chi Hay Tong, and Ingmar Posner. Vote3deep: Fast object detection in 3d point clouds using efficient convolutional neural networks. In *ICRA*, pages 1355–1361. IEEE, 2017. 1, 2
- [12] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, pages 3354–3361. IEEE, 2012. 2, 6
- [13] Tianrui Guan, Jun Wang, Shiyi Lan, Rohan Chandra, Zuxuan Wu, Larry Davis, and Dinesh Manocha. M3det: Multi-representation, multi-scale, mutual-relation 3d object detection with transformers. *arXiv preprint arXiv:2104.11896*, 2021. 3, 6, 7
- [14] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187–199, 2021. 4
- [15] Chenhang He, Hui Zeng, Jianqiang Huang, Xian-Sheng Hua, and Lei Zhang. Structure aware single-stage 3d object detection from point cloud. In *CVPR*, pages 11873–11882, 2020. 1, 3, 6, 7
- [16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9729–9738, 2020. 5, 6
- [17] Tengpeng Huang, Zhe Liu, Xiwu Chen, and Xiang Bai. Ep-net: Enhancing point features with image semantics for 3d object detection. In *ECCV*, pages 35–52. Springer, 2020. 2, 3, 6, 7, 8
- [18] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L Waslander. Joint 3d proposal generation and object detection from view aggregation. In *IROS*, pages 1–8. IEEE, 2018. 2, 3, 6, 7, 8
- [19] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, pages 12697–12705, 2019. 1, 2, 6, 7
- [20] Buyu Li, Wanli Ouyang, Lu Sheng, Xingyu Zeng, and Xiaogang Wang. Gs3d: An efficient 3d object detection framework for autonomous driving. In *CVPR*, pages 1019–1028, 2019. 2
- [21] Peiliang Li, Xiaozhi Chen, and Shaojie Shen. Stereo r-cnn based 3d object detection for autonomous driving. In *CVPR*, pages 7644–7652, 2019. 2
- [22] Ming Liang, Bin Yang, Yun Chen, Rui Hu, and Raquel Urtasun. Multi-task multi-sensor fusion for 3d object detection. In *CVPR*, pages 7345–7353, 2019. 1, 2, 3, 6, 7
- [23] Ming Liang, Bin Yang, Shenlong Wang, and Raquel Urtasun. Deep continuous fusion for multi-sensor 3d object detection. In *ECCV*, pages 641–656, 2018. 2, 3, 6, 7
- [24] Song Liu, Haoqi Fan, Shengsheng Qian, Yiru Chen, Wenkui Ding, and Zhongyuan Wang. Hit: Hierarchical transformer with momentum contrast for video-text retrieval. *arXiv preprint arXiv:2103.15049*, 2021. 5
- [25] Yunze Liu, Li Yi, Shanghang Zhang, Qingnan Fan, Thomas Funkhouser, and Hao Dong. P4contrast: Contrastive learning with pairs of point-pixel pairs for rgb-d scene understanding. *arXiv preprint arXiv:2012.13089*, 2020. 5
- [26] Ze Liu, Zheng Zhang, Yue Cao, Han Hu, and Xin Tong. Group-free 3d object detection via transformers. *arXiv preprint arXiv:2104.00678*, 2021. 3
- [27] Zhe Liu, Xin Zhao, Tengpeng Huang, Ruolan Hu, Yu Zhou, and Xiang Bai. Tanet: Robust 3d object detection from point clouds with triple attention. In *AAAI*, pages 11677–11684, 2020. 6, 7
- [28] 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*, pages 3164–3173, 2021. 3
- [29] Ishan Misra, Rohit Girdhar, and Armand Joulin. An end-to-end transformer model for 3d object detection. In *ICCV*, pages 2906–2917, 2021. 3

- [30] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. In *CVPR*, pages 7074–7082, 2017. 2
- [31] Xuran Pan, Zhuofan Xia, Shiji Song, Li Erran Li, and Hao Huang. 3d object detection with pointformer. In *CVPR*, pages 7463–7472, 2021. 3, 6, 7
- [32] Su Pang, Daniel Morris, and Hayder Radha. Clocs: Camera-lidar object candidates fusion for 3d object detection. In *IROS*, pages 10386–10393. IEEE, 2020. 2, 3, 6, 7
- [33] 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*, pages 918–927, 2018. 2, 3, 6, 7
- [34] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, pages 652–660, 2017. 1, 2
- [35] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NeurIPS*, pages 5099–5108, 2017. 1, 2
- [36] Hualian Sheng, Sijia Cai, Yuan Liu, Bing Deng, Jianqiang Huang, Xian-Sheng Hua, and Min-Jian Zhao. Improving 3d object detection with channel-wise transformer. In *ICCV*, pages 2743–2752, 2021. 3
- [37] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rnn: Point-voxel feature set abstraction for 3d object detection. In *CVPR*, pages 10529–10538, 2020. 1, 3, 6, 7
- [38] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointcnn: 3d object proposal generation and detection from point cloud. In *CVPR*, pages 770–779, 2019. 1, 2, 6, 7
- [39] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. From points to parts: 3d object detection from point cloud with part-aware and part-aggregation network. *IEEE TPAMI*, 2020. 6, 7
- [40] Weijing Shi and Raj Rajkumar. Point-gnn: Graph neural network for 3d object detection in a point cloud. In *CVPR*, pages 1711–1719, 2020. 1, 2
- [41] Sourabh Vora, Alex H Lang, Bassam Helou, and Oscar Beijbom. Pointpainting: Sequential fusion for 3d object detection. In *CVPR*, pages 4604–4612, 2020. 2, 3, 6, 7
- [42] Chunwei Wang, Chao Ma, Ming Zhu, and Xiaokang Yang. Pointaugmenting: Cross-modal augmentation for 3d object detection. In *CVPR*, pages 11794–11803, 2021. 2, 5
- [43] Dominic Zeng Wang and Ingmar Posner. Voting for voting in online point cloud object detection. In *Robotics: Science and Systems*, volume 1, pages 10–15. Rome, Italy, 2015. 1, 2
- [44] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *CVPR*, pages 8445–8453, 2019. 2
- [45] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *TOG*, 38(5):1–12, 2019. 1
- [46] Zhixin Wang and Kui Jia. Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal 3d object detection. In *IROS*, pages 1742–1749. IEEE, 2019. 2, 3, 6, 7
- [47] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. *arXiv preprint arXiv:2103.15808*, 2021. 4
- [48] Liang Xie, Chao Xiang, Zhengxu Yu, Guodong Xu, Zheng Yang, Deng Cai, and Xiao-fei He. Pi-rnn: An efficient multi-sensor 3d object detector with point-based attentive conv fusion module. In *AAAI*, pages 12460–12467, 2020. 2, 3, 6, 7
- [49] Danfei Xu, Dragomir Anguelov, and Ashesh Jain. Pointfusion: Deep sensor fusion for 3d bounding box estimation. In *CVPR*, pages 244–253, 2018. 6, 7
- [50] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018. 1, 2, 5, 6, 7
- [51] Zetong Yang, Yanan Sun, Shu Liu, and Jiaya Jia. 3dssd: Point-based 3d single stage object detector. In *CVPR*, pages 11040–11048, 2020. 1, 2, 6, 7
- [52] Zetong Yang, Yanan Sun, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Std: Sparse-to-dense 3d object detector for point cloud. In *ICCV*, pages 1951–1960, 2019. 1, 3, 6, 7
- [53] Jin Hyeok Yoo, Yecheol Kim, Jisong Kim, and Jun Won Choi. 3d-cvf: Generating joint camera and lidar features using cross-view spatial feature fusion for 3d object detection. In *ECCV*, pages 720–736. Springer, 2020. 2, 3, 6, 7, 8
- [54] Yurong You, Yan Wang, Wei-Lun Chao, Divyansh Garg, Geoff Pleiss, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar++: Accurate depth for 3d object detection in autonomous driving. *arXiv preprint arXiv:1906.06310*, 2019. 2
- [55] Kun Yuan, Shaopeng Guo, Ziwei Liu, Aojun Zhou, Fengwei Yu, and Wei Wu. Incorporating convolution designs into visual transformers. *arXiv preprint arXiv:2103.11816*, 2021. 4
- [56] Yanan Zhang, Di Huang, and Yunhong Wang. Pc-rnn: Point cloud completion and graph neural network for 3d object detection. In *AAAI*, pages 3430–3437, 2021. 1, 2
- [57] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *ICCV*, pages 16259–16268, 2021. 4
- [58] Xin Zhao, Zhe Liu, Ruolan Hu, and Kaiqi Huang. 3d object detection using scale invariant and feature reweighting networks. In *AAAI*, pages 9267–9274, 2019. 6, 7
- [59] Wu Zheng, Weiliang Tang, Li Jiang, and Chi-Wing Fu. Sssd: Self-ensembling single-stage object detector from point cloud. In *CVPR*, pages 14494–14503, 2021. 2, 6, 7
- [60] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *CVPR*, pages 4490–4499, 2018. 1, 2, 6, 7