

H3DNet: 3D Object Detection Using Hybrid Geometric Primitives

Zaiwei Zhang¹, Bo Sun^{*1}, Haitao Yang^{*1}, and Qixing Huang¹

The University of Texas at Austin, Austin, Texas, USA, 78710

Abstract. We introduce H3DNet, which takes a colorless 3D point cloud as input and outputs a collection of oriented object bounding boxes (or BB) and their semantic labels. The critical idea of H3DNet is to predict a hybrid set of geometric primitives, i.e., BB centers, BB face centers, and BB edge centers. We show how to convert the predicted geometric primitives into object proposals by defining a distance function between an object and the geometric primitives. This distance function enables continuous optimization of object proposals, and its local minimums provide high-fidelity object proposals. H3DNet then utilizes a matching and refinement module to classify object proposals into detected objects and fine-tune the geometric parameters of the detected objects. The hybrid set of geometric primitives not only provides more accurate signals for object detection than using a single type of geometric primitives, but it also provides an overcomplete set of constraints on the resulting 3D layout. Therefore, H3DNet can tolerate outliers in predicted geometric primitives. Our model achieves state-of-the-art 3D detection results on two large datasets with real 3D scans, ScanNet and SUN RGB-D. [Our code is open-sourced at here.](#)

Keywords: 3D Deep Learning, Geometric Deep Learning, 3D Point Clouds, 3D Bounding Boxes, 3D Object Detection

1 Introduction

Object detection is a fundamental problem in visual recognition. In this work, we aim to detect the 3D layout (i.e., oriented 3D bounding boxes (or BBs) and associated semantic labels) from a colorless 3D point cloud. This problem is fundamentally challenging because of the irregular input and a varying number of objects across different scenes. Choosing suitable intermediate representations to integrate low-level object cues into detected objects is key to the performance of the resulting system. While early works [39, 40] classify sliding windows for object detection, recent works [33, 4, 17, 30, 55, 29, 52, 36, 52, 28, 46] have shown the great promise of designing end-to-end neural networks to generate, classify, and refine object proposals.

This paper introduces H3DNet, an end-to-end neural network that utilizes a novel intermediate representation for 3D object detection. Specifically, H3DNet

* Equal Contribution



Fig. 1: Our approach leverages a hybrid and overcomplete set of geometric primitives to detect and refine 3D object bounding boxes (BBs). Note that red BBs are initial object proposals, green BBs are refined object proposals, and blue surfaces and lines are hybrid geometric primitives.

first predicts a hybrid and overcomplete set of geometric primitives (i.e., BB centers, BB face centers, and BB edge centers) and then detects objects to fit these primitives and their associated features. This regression methodology, which is motivated from the recent success of keypoint-based pose regression for 6D object pose estimation [20, 26, 12, 22, 27, 37], displays two appealing advantages for 3D object detection. First, each type of geometric primitives focuses on different regions of the input point cloud (e.g., points of an entire object for predicting the BB center and points of a planar boundary surface for predicting the corresponding BB face center). Combining diverse primitive types can add the strengths of their generalization behaviors. On new instances, they offer more useful constraints and features than merely using one type of primitives. Second, having an overcomplete set of primitive constraints can tolerate outliers in predicted primitives (e.g., using robust functions) and reduce the influence of individual prediction errors. The design of H3DNet fully practices these two advantages.

Specifically, H3DNet consists of three modules. The first module computes dense pointwise descriptors and uses them to predict geometric primitives and their latent features. The second module converts these geometric primitives into object proposals. A key innovation of H3DNet is to define a parametric distance function that evaluates the distance between an object BB and the predicted primitives. This distance function can easily incorporate diverse and overcomplete geometric primitives. Its local minimums naturally correspond to object proposals. This method allows us to optimize object BBs continuously and generate high-quality object proposals from imprecise initial proposals.

The last module of H3DNet classifies each object proposal as a detected object or not, and also predicts for each detected object an offset vector of its geometric parameters and a semantic label to fine-tune the detection result. The performance of this module depends on the input. As each object proposal is associated with diverse geometric primitives, H3DNet aggregates latent features associated with these primitives, which may contain complementary semantic and geometric information, as the input to this module. We also introduce a network design that can handle a varying number of geometric primitives.

We have evaluated H3DNet on two popular benchmark datasets ScanNet and SUN RGB-D. On ScanNet, H3DNet achieved 67.2% in mAP (0.25), which corresponded to a 8.5% relative improvement from state-of-the-art methods

that merely take the 3D point positions as input. On SUN RGB-D, H3DNet achieved 60.1% in mAP (0.25), which corresponded to a 2.4% relative improvement from the same set of state-of-the-art methods. Moreover, on difficult categories of both datasets (i.e., those with low mAP scores), the performance gains of H3DNet are significant (e.g., from 38.1/47.3/57.1 to 51.9/61.0/75.3/ on window/door/shower-curtain, respectively). We have also performed an ablation study on H3DNet. Experimental results justify the importance of regressing a hybrid and overcomplete set of geometric primitives for generating object proposals and aggregating features associated with matching primitives for classifying and refining detected objects. In summary, the contributions of our work are:

- Formulation of object detection as regressing and aggregating an overcomplete set of geometric primitives
- Predicting multiple types of geometric primitives that are suitable for different object types and scenes
- State-of-the-art results on SUN RGB-D and ScanNet with only point clouds

2 Related Works

3D object detection. From the methodology perspective, there are strong connections between 3D object detection approaches and their 2D counterparts. Most existing works follow the approach of classifying candidate objects that are generated using a sliding window [39, 40] or more advanced techniques [30, 55, 29, 52, 36, 52, 28, 46]. Objectness classification involves template-based approaches or deep neural networks. The key differences between 2D approaches and 3D approaches lie in feature representations. For example, [24] leverages a pairwise semantic context potential to guide the proposals’ objectness score. [33] uses clouds of oriented gradients (COG) for object detection. [10] utilizes the power of 3D convolution neural networks to identify locations and keypoints of 3D objects. Due to the computational cost in the 3D domain, many methods utilize 2D-3D projection techniques to integrate 2D object detection and 3D data processing. For example, MV3D [4] and VoxelNet [55] represent the 3D input data in a bird’s-eye view before proceeding to the rest of the pipeline. Similarly, [14, 17, 30] first process 2D inputs to identify candidate 3D object proposals.

Point clouds have emerged as a powerful representation for 3D deep learning, particularly for extracting salient geometric features and spatial locations (c.f. [31, 32]). Prior usages of point-based neural networks include classification [31, 32, 21, 45, 16, 49, 47, 48, 9, 11], segmentation [32, 41, 2, 21, 45, 43, 16, 49, 47, 48, 9, 11, 7, 46], normal estimation [2], and 3D reconstruction [42, 6, 50].

There are also growing interests in object detection from point clouds [29, 52, 36, 52, 28, 46]. H3DNet is most relevant to [29], which leverages deep neural networks to predict object bounding boxes. The key innovation of H3DNet is that it utilizes an overcomplete set of geometric primitives and a distance function to integrate them for object detection. This strategy can tolerate inaccurate primitive predictions (e.g., due to partial inputs).

Multi-task 3D understanding. Jointly predicting different types of geometric primitives is related to multi-task learning [3, 13, 35, 34, 53, 25, 23, 19, 28, 56, 54,

8], where incorporating multiple relevant tasks together boosts the performance of feature learning. In a recent work HybridPose [37], Song et al. show that predicting keypoints, edges between keypoints, and symmetric correspondences jointly lift the prediction accuracies of each type of features. In this paper, we show that predicting BB centers, BB face centers, and BB edge centers together help to improve the generalization behavior of primitive predictions.

Overcomplete constraints regression. The main idea of H3DNet is to incorporate an overcomplete set of constraints. This approach achieves considerable performance gains from [29], which uses a single type of geometric primitives. At a conceptual level, similar strategies have been used in tasks of object tracking [44], zero-shot fine-grained classification [1], 6D object pose estimation [37] and relative pose estimation between scans [51], among others. Compared to these works, the novelties of H3DNet lie in designing hybrid constraints that are suitable for object detection, continuous optimization of object proposals, aggregating hybrid features for classifying and fine-tuning object proposals, and end-to-end training of the entire network.

3 Approach

This section describes the technical details of H3DNet. Section 3.1 presents an approach overview. Section 3.2 to Section 3.5 elaborate on the network design and the training procedure of H3DNet.

3.1 Approach Overview

As illustrated in Figure 2, the input of H3DNet is a dense set of 3D points (i.e., a point cloud) $S \in \mathbb{R}^{3 \times n}$ ($n = 40000$ in this paper). Such an input typically comes from depth sensors or the result of multi-view stereo matching. The output is given by a collection of (oriented) bounding boxes (or BB) $\mathcal{O}_S \in \overline{\mathcal{O}}$, where $\overline{\mathcal{O}}$ denotes the space of all possible objects. Each object $o \in \overline{\mathcal{O}}$ is given by its class label $l_o \in \mathcal{C}$, where \mathcal{C} is pre-defined, its center $\mathbf{c}_o = (c_o^x, c_o^y, c_o^z)^T \in \mathbb{R}^3$ in a world coordinate system, its scales $\mathbf{s}_o = (s_o^x, s_o^y, s_o^z)^T \in \mathbb{R}^3$, and its orientation $\mathbf{n}_o = (\mathbf{n}_o^x, \mathbf{n}_o^y)^T \in \mathbb{R}^2$ in the xy-plane of the same world coordinate system (note that the upright direction of an object is always along the z axis).

H3DNet consists of three modules, starting from geometric primitive prediction, to proposal generation, to object refinement. The theme is to predict and integrate an overcomplete set of geometric primitives, i.e., BB centers, BB face centers, and BB edge centers. The entire network is trained end-to-end.

Geometric primitive module. The first module of H3DNet takes a point cloud S as input and outputs a set of geometric primitives \mathcal{P}_S that predicts locations of BB centers, BB face centers, and BB edge centers of the underlying objects. The network design extends that of [29]. Specifically, it combines a sub-module for extracting dense point-wise descriptors and sub-modules that take point-wise descriptors as input and output offset vectors between input points and the corresponding centers. The resulting primitives are obtained through

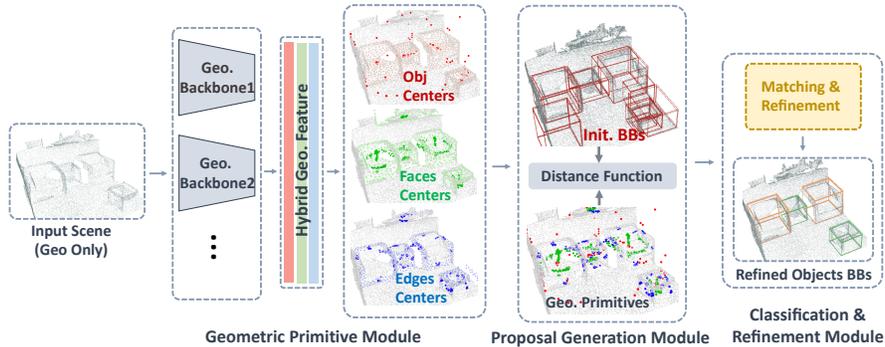


Fig. 2: H3DNet consists of three modules. The first module computes a dense descriptor and predicts three geometric primitives, namely, BB centers, BB face centers, and BB edge centers. The second module converts geometric primitives into object proposals. The third module classifies object proposals and refines the detected objects.

clustering. In addition to locations, each predicted geometric primitive also possesses a latent feature that is passed through subsequent modules of H3DNet.

In contrast to [29], H3DNet exhibits two advantages. First, since only a subset of predicted geometric primitives is sufficient for object detection, the detected objects are insensitive to erroneous predictions. Second, different types of geometric primitives show complementary strength. For example, BB centers are accurate for complete and solid objects, while BB face centers are suitable for partial objects that possess rich planar structures.

Proposal generation module. The second module takes predicted geometric primitives as input and outputs a set of object proposals. A critical innovation of H3DNet is to formulate object proposals as local minimums of a distance function. This methodology is quite flexible in several ways. First, it is easy to incorporate overcomplete geometric primitives, each of which corresponds to an objective term of the distance function. Second, it can handle outlier predictions and mispredictions using robust norms. Finally, it becomes possible to optimize object proposals continuously, and this property relaxes the burden of generating high-quality initial proposals.

Classification and refinement module. The last module of H3DNet classifies each object proposal into a detected object or not. This module also computes offset vectors to refine the BB center, BB size, and BB orientation of each detected object, and a semantic label. The key idea of this module is to aggregate features of the geometric primitives that are close to the corresponding primitives of each object proposal. Such aggregated features carry rich semantic information that is unavailable in the feature associated with each geometric primitive.

3.2 Primitive Module

The first module of H3DNet predicts a set of geometric primitives from the input point cloud. Each geometric primitive provides some constraints on the

detected objects. In contrast to most prior works that compute a minimum set of primitives, i.e., that is sufficient to determine the object bounding boxes, H3DNet leverages an overcomplete set of geometric primitives, i.e., BB centers, BB face centers, and BB edge centers. In other words, these geometric primitives can provide up-to 19 positional constraints for one BB. As we will see later, they offer great flexibilities in generating, classifying, and refining object proposals.

Similar to [29], the design of this module combines a descriptor sub-module and a prediction sub-module. The descriptor sub-module computes dense point-wise descriptors. Its output is fed into the prediction sub-module, which consists of three prediction branches. Each branch predicts one type of geometric primitives. Below we provide the technical details of the network design.

Descriptor sub-module. The output of the descriptor sub-module provides semantic information to group points for predicting geometric primitives (e.g., points of the same object for BB centers and points of the same planar boundary faces for BB face centers). Instead of using a single descriptor computation tower [29], H3DNet integrates four separate descriptor computation towers. The resulting descriptors are concatenated together for primitive prediction and subsequent modules of H3DNet. Our experiments indicate that this network design can learn distinctive features for predicting each type of primitives. However, it does not lead to a significant increase in network complexity.

BB center prediction. The same as [29], H3DNet leverages a network with three fully connected layers to predict the offset vector between each point and its corresponding object center. The resulting BB centers are obtained through clustering (c.f. [29]). Note that in addition to offset vectors, H3DNet also computes an associated feature descriptor for each BB center. These feature descriptors serve as input feature representations for subsequent modules of H3DNet.

Predictions of BB centers are accurate on complete and rectangular shaped objects. However, there are shifting errors for partial and/or occluded objects, and thin objects, such as pictures or curtains, due to imbalanced weighting for offset prediction. This motivates us to consider centers of BB faces and BB edges.

BB face center prediction. Planar surfaces are ubiquitous in man-made scenes and objects. Similar to BB center, H3DNet uses 3 fully connected layers to perform point-wise predictions. The predicted attributes include a flag that indicates whether a point is close to a BB face or not and if so, an offset vector between that point and its corresponding BB face center. For training, we generate the ground-truth labels by computing the closest BB face for each point. We say a point lies close to a BB face (i.e., a positive instance) if that distance is smaller than $0.2m$. Similar to BB centers, each BB face center prediction also possesses a latent feature descriptor that is fed into the subsequent modules.

Since face center predictions are only affected by points that are close to that face, we found that they are particularly useful for objects with rich planar patches (e.g., refrigerator and shower-curtain) and incomplete objects.

BB edge center prediction. Boundary line features form another type of geometric cues in all 3D scenes and objects. Similar to BB faces, H3DNet employs 3 fully connected layers to predict for each point a flag that indicates whether it is close to a BB edge or not and if so, an offset vector between that point and

the corresponding BB edge center. The same as BB face centers, we generate ground-truth labels by computing the closest BB edge for each point. We say a point lies close to a BB edge if the closest distance is smaller than $0.2m$. Again, each BB edge center prediction possesses a latent feature of the same dimension. Compared to BB centers and BB face centers, BB edge centers are useful for objects where point densities are irregular (e.g., with large holes) but BB edges appear to be complete (e.g., window and computer desk).

As analyzed in details in the supplemental material, error distributions of different primitives are largely uncorrelated with each other. Such uncorrelated prediction errors provide a foundation for performance boosting when integrating them together for detecting objects.

3.3 Proposal Module

After predicting geometric primitives, H3DNet proceeds to compute object proposals. Since the predicted geometric primitives are overcomplete, H3DNet converts them into a distance function and generates object proposals as local minimums of this distance function. This approach, which is the crucial contribution of H3DNet, exhibits several appealing properties. First, it automatically incorporates multiple geometric primitives to determine the parameters of each object proposal. Second, the distance function can optimize object proposals continuously. The resulting local minimums are insensitive to initial proposals, allowing us to use simple initial proposal generators. Finally, each local minimum is attached to different types of geometric primitives, which carry potentially complementary semantic information. As discussed in Section 3.4, the last module of H3DNet builds upon this property to classify and refine object proposals.

Proposal distance function. The proposal distance function $F_S(o)$ measures a cumulative proximity score between the predicted geometric primitives \mathcal{P}_S and the corresponding object primitives of an object proposal o . Recall that $l_o \in \mathcal{C}$, $\mathbf{c}_o \in \mathbb{R}^3$, $\mathbf{s}_o \in \mathbb{R}^3$, and $\mathbf{n}_o \in \mathbb{R}^2$ denote the label, center, scales, and orientation of o . With $\mathbf{o} = (\mathbf{c}_o^T, \mathbf{s}_o^T, \mathbf{n}_o^T)^T$ we collect all the geometric parameters of o . Note that each object proposal o has 19 object primitives (i.e., one BB center, six BB face centers, and twelve BB edge centers). Let $\mathbf{p}_i(o), 1 \leq i \leq 19$ be the location of the i -th primitive of o . Denote $t_i \in \mathcal{T} := \{\text{center, face, edge}\}$ as the type of the i -th primitive. Let $\mathcal{P}_{t,S} \subseteq \mathcal{P}_S$ collect all predicted primitives with type $t \in \mathcal{T}$. We define

$$F_S(o) := \sum_{t \in \mathcal{T}} \beta_t \sum_{\mathbf{c} \in \mathcal{P}_{t,S}} \min \left(\min_{1 \leq i \leq 19, t_i = t} \|\mathbf{c}_i - \mathbf{p}_i(o)\|^2 - \delta, 0 \right). \quad (1)$$

In other words, we employ the truncated L2-norm to match predicted primitives and closest object primitives. β_t describes the trade-off parameter for type t . Both β_t and the truncation threshold δ are determined via cross-validation.

Initial proposals. H3DNet detects object proposals by exploring the local minimums of the distance function from a set of initial proposals. From the perspective of optimization, we obtain the same local minimum from any initial solution that is sufficiently close to that local minimum. This means the initial proposals

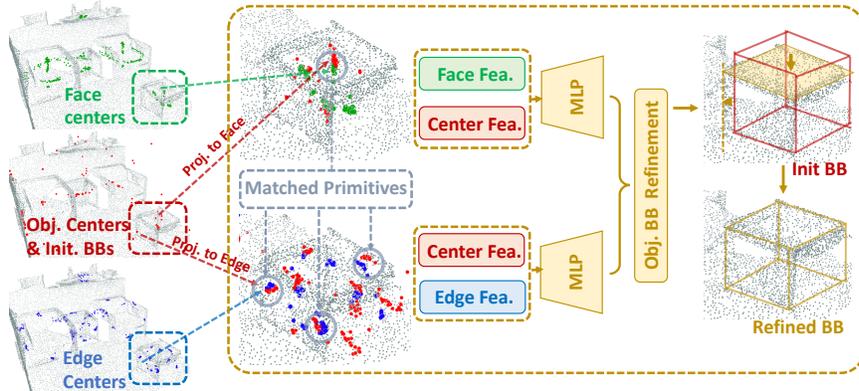


Fig. 3: Illustration of the matching, feature aggregation and refinement process.

do not need to be exact. In our experiments, we found that a simple object proposal generation approach is sufficient. Specifically, H3DNet utilizes the method of [29], which initializes an object proposal from each detected BB center.

Proposal refinement. By minimizing F_S , we refine each initial proposal. Note that different initial proposals may share the same local minimum. The final object proposals only collect distinctive local minimums.

3.4 Classification and Refinement Module

The last module of H3DNet takes the output of the proposal module as input and outputs a collection of detected objects. This module combines a classification sub-module and a refinement sub-module. The classification sub-module determines whether each object proposal is an object or not. The refinement sub-module predicts for each detected object the offsets in BB center, BB size, and BB orientation and a semantic label.

The main idea is to aggregate features associated the primitives (i.e., object centers, edge centers, and face centers) of each object proposal. Such features capture potentially complementary information, yet only at this stage (i.e., after we have detected groups of matching primitives) it becomes possible to fuse them together to determine and fine-tune the detected objects.

As illustrated in Figure 3, we implement this sub-module by combing four fully connected layers. The input layer concatenates input features of 19 object primitives of an object proposal (i.e., one BB center, six BB face centers, and twelve BB edge centers). Each input feature integrates features associated with primitives that are in the neighborhood of the corresponding object primitive. To address the issue that there is a varying number of neighborhood primitives (e.g., none or multiple), we utilize a variant of the max-pooling layer in PointNet [31, 32] to compute the input feature. Specifically, the input to each max-pooling layer consists of the feature associated with the input object proposal, which addresses the issue of no matching primitives, and 32 feature points that are randomly

sampled in the neighborhood of each object primitive. In our implementation, we determine the neighboring primitives via range query, and the radius is $0.05m$.

The output of this module combines the label that indicates objectiveness, offsets in BB center, BB size, and BB orientation, and a semantic label.

3.5 Network Training

Training H3DNet employs a loss function with five objective terms:

$$\begin{aligned} \min_{\theta_g, \theta_p, \theta_c, \theta_o} \quad & \lambda_g l_g(\theta_g) + \lambda_p l_p(\theta_g, \theta_p) + \lambda_f l_f(\theta_g, \theta_p, \theta_o) \\ & + \lambda_c l_c(\theta_g, \theta_p, \theta_c) + \lambda_o l_o(\theta_g, \theta_p, \theta_o) \end{aligned} \quad (2)$$

where l_g trains the geometric primitive module θ_g , l_p trains the proposal module θ_p , l_f trains the potential function and refinement sub-network θ_o , l_c trains the classification sub-network θ_c , and l_o trains the refinement sub-network. The trade-off parameters λ_g , λ_p , λ_f , λ_c , and λ_o are determined through 10-fold cross-validation. Intuitively, l_c , l_o and l_f provide end-to-end training of H3DNet, while l_g and l_p offer intermediate supervisions.

Formulation. Formulations of l_g , l_p , l_c , and l_o follow common strategies in the literature. Specifically, both l_g and l_p utilize L2 regression losses and a cross-entropy loss for geometric primitive location and existence flag prediction, and initial proposal generation; l_c applies a cross-entropy loss to train the object classification sub-network; l_o employs L2 regression losses for predicting the shape offset, and a cross-entropy loss for predicting the semantic label. Since these four loss terms are quite standard, we leave the details to the supplemental material.

l_f seeks to match the local minimums of the potential function and the underlying ground-truth objects. Specifically, consider a parametric potential function $f_\Theta(\mathbf{x})$ parameterized by Θ . Consider a local minimum \mathbf{x}_Θ^* which is a function of Θ . Let \mathbf{x}^{gt} be the target location of \mathbf{x}_Θ^* . We define the following alignment potential to pull \mathbf{x}_Θ^* to close to \mathbf{x}^{gt} :

$$l_m(\mathbf{x}_\Theta^*, \mathbf{x}^{gt}) := \|\mathbf{x}_\Theta^* - \mathbf{x}^{gt}\|^2. \quad (3)$$

The following proposition describes how to compute the derivatives of l_m with respect to Θ . The proof is deferred to the supp. material.

Proposition 1. *The derivatives of l_m with respect to Θ is given by*

$$\frac{\partial l_m}{\partial \Theta} := 2(\mathbf{x}_\Theta^* - \mathbf{x}^{gt})^T \cdot \frac{\partial \mathbf{x}_\Theta^*}{\partial \Theta}, \quad \frac{\partial \mathbf{x}_\Theta^*}{\partial \Theta} := -\left(\frac{\partial^2 f_\Theta(\mathbf{x}^*)}{\partial^2 \mathbf{x}}\right)^{-1} \cdot \frac{\partial^2 f_\Theta(\mathbf{x}^*)}{\partial \mathbf{x} \partial \Theta}. \quad (4)$$

We proceed to use l_m to define l_f . For each scene S , we denote the set of ground-truth objects and the set of local minimums of potential function F_S as \mathcal{O}^{gt} and \mathcal{O}^* , respectively. Note that \mathcal{O}^* depends on the network parameters and hyper-parameters. Let $\mathcal{C}_S \subset \mathcal{O}^{gt} \times \mathcal{O}^*$ collect the nearest object in \mathcal{O}^* for each object in \mathcal{O}^{gt} . Consider a training set of scenes \mathcal{S}_{train} , we define

$$l_f := \sum_{S \in \mathcal{S}_{train}} \sum_{(\mathbf{o}^*, \mathbf{o}^{gt}) \in \mathcal{C}_S} l_m(\mathbf{o}^*, \mathbf{o}^{gt}). \quad (5)$$

Table 1: 3D object detection results on ScanNet V2 val dataset. We show per-category results of average precision (AP) with 3D IoU threshold 0.25 as proposed by [38], and mean of AP across all semantic classes with 3D IoU threshold 0.25.

	RGB	cab	bed	chair	sofa	tabl	door	wind	bkskf	pic	cntr	desk	curt	fridg	showr	toil	sink	bath	ofurn	mAP
3DSIS-5[10]	✓	19.8	69.7	66.2	71.8	36.1	30.6	10.9	27.3	0.0	10.0	46.9	14.1	53.8	36.0	87.6	43.0	84.3	16.2	40.2
3DSIS[10]	✗	12.8	63.1	66.0	46.3	26.9	8.0	2.8	2.3	0.0	6.9	33.3	2.5	10.4	12.2	74.5	22.9	58.7	7.1	25.4
Votenet[29]	✗	36.3	87.9	88.7	89.6	58.8	47.3	38.1	44.6	7.8	56.1	71.7	47.2	45.4	57.1	94.9	54.7	92.1	37.2	58.7
Ours	✗	49.4	88.6	91.8	90.2	64.9	61.0	51.9	54.9	18.6	62.0	75.9	57.3	57.2	75.3	97.9	67.4	92.5	53.6	67.2
w/o refine	✗	37.2	89.3	88.4	88.5	64.4	53.0	44.2	42.2	11.1	51.2	59.8	47.0	54.3	74.3	93.1	57.0	85.6	43.5	60.2

Computing the derivatives of l_f with respect to the network parameters is a straightforward application of Prop.1.

Training. We train H3DNet end-to-end and from scratch with the Adam optimizer [15]. Please defer to the supplemental material for hyper-parameters used in training, such as learning rate etc.

4 Experimental Results

In this section, we first describe the experiment setup in Section 4.2. Then, we compare our method with current state-of-the-art 3D object detection methods quantitatively, and analyze our results in Section 4.2, where we show the importance of using geometric primitives and discuss our advantages. Finally, we show ablation results in Section 4.3 and qualitative comparison in Figures (6) and (5). More results and discussions can be found in the supplemental material.

4.1 Experimental Setup

Datasets. We employ two popular datasets ScanNet V2[5] and SUN RGB-D V1[38]. ScanNet is a dataset of richly-annotated 3D reconstructions of indoor scenes. It contains 1513 indoor scenes annotated with per-point instance and semantic labels for 40 semantic classes. SUN RGB-D is a single-view RGB-D dataset for 3D scene understanding, which contains 10335 indoor RGB and depth images with per-point semantic labels and object bounding boxes. For both datasets, we use the same training/validation split and BB semantic classes (18 classes for ScanNet and 10 classes for SUN RGB-D) as in VoteNet[29] and sub-sample 40000 points from every scene.

Evaluation protocol. We use Average Precision(AP) and the mean of AP across all semantic classes (mAP)[38] under different IoU values (the minimum IoU to consider a positive match). Average precision computes the average precision value for recall value over 0 to 1. IoU is given by the ratio of the area of intersection and area of union of the predicted bounding box and ground truth bounding box. Specifically, we use AP/mAP@0.25 and AP/mAP@0.5.

Baseline Methods We compare H3DNet with STAR approaches: VoteNet [29] is a geometric-only detector that combines deep point set networks and a voting procedure. GSPN[52] uses a generative model for instance segmentation. Both

Table 2: **Left:** 3D object detection results on ScanNetV2 val set. **Right:** results on SUN RGB-D V1 val set. We show mean of average precision (mAP) across all semantic classes with 3D IoU threshold 0.25 and 0.5.

	Input	mAP@0.25	mAP@0.5		Input	mAP@0.25	mAP@0.5
DSS[40]	Geo + RGB	15.2	6.8	DSS[40]	Geo + RGB	42.1	-
F-PointNet[30]	Geo + RGB	19.8	10.8	COG[33]	Geo + RGB	47.6	-
GSPN[52]	Geo + RGB	30.6	17.7	2D-driven[18]	Geo + RGB	45.1	-
3D-SIS [10]	Geo + 5 views	40.2	22.5	F-PointNet[30]	Geo + RGB	54.0	-
VoteNet [29]	Geo only	58.7	33.5	VoteNet [29]	Geo only	57.7	32.9
Ours	Geo only	67.2	48.1	Ours	Geo only	60.1	39.0
w/o refine	Geo only	60.2	37.3	w/o refine	Geo only	58.5	34.2

3D-SIS [10] and DSS [40] extract features from 2D images and 3D shapes to generate object proposals. F-PointNet [30] and 2D-Driven [18] first propose 2D detection regions and project them to 3D frustum for 3D detection. Cloud of gradient(COG) [33] integrates sliding windows with a 3D HoG-like feature.

4.2 Analysis of Results

As shown in Table 2, our approach leads to an average mAP score of 67.2%, with 3D IoU threshold 0.25 (mAP@0.25), on ScanNet V2, which is 8.5% better than the top-performing baseline approach [29]. In addition, our approach is 14.6% better than the baseline approach [29] with 3D IoU threshold 0.5 (mAP@0.5). For SUN RGB-D, our approach gains 2.4% and 6.1% in terms of mAP, with 3D IoU threshold 0.25 and 0.5 respectively. On both datasets, the performance gains of our approach under mAP@0.5 are larger than those under mAP@0.25, meaning our approach offers more accurate predictions than baseline approaches. Such improvements are attributed to using an overcomplete set of geometric primitives and their associated features for generating and refining object proposals. We can also understand the relative less salient improvements on SUN RGB-D than ScanNet in a similar way, i.e., labels of the former are less accurate than the latter, and the strength of H3DNet is not fully utilized on SUN RGB-D. Except for the classification and refinement module, our approach shares similar computation pipeline and complexity with VoteNet. The computation on multiple descriptor towers and proposal modules can be paralleled, which should not increase computation overhead. In our implementation, our approach requires 0.058 seconds for the last module per scan. Conceptually, our approach requires 50% more time compared to [29] but operates with a higher detection accuracy.

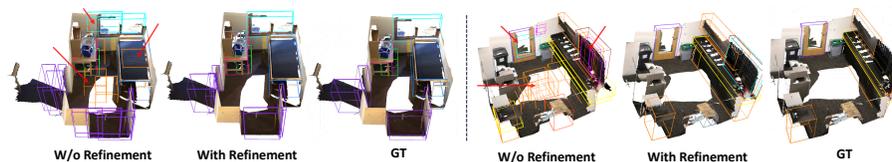


Fig. 4: Effect of geometric primitive matching and refinement.

Improvement on thin objects. One limitation of the current top-performing baseline [29] is predicting thin objects in 3D scenes, such as doors, windows and pictures. In contrast, with face and edge primitives, H3DNet is able to extract better features for those thin objects. For example, the frames of window or picture provide dense edge feature, and physical texture of curtain or shower-curtain provide dense face/surface feature. As shown in Table 1, H3DNet leads to significant performance gains on thin objects, such as door (13.7%), window (13.8%), picture (10.8%), curtain (10.1%) and shower-curtain (18.2%).

Improvement on objects with dense geometric primitives. Across the individual object classes in ScanNet in Table 1, other than those thin objects, our approach also leads to significant performance gain on cabinet (13.1%), table (6.1%), bookshelf (10.3%), refrigerator (11.8%), sink (12.7%) and other-furniture (16.4%). One explanation is that the geometric shapes of these object classes possess rich planar structures and/or distinct edge structures, which contribute greatly on geometric primitive detection and object refinement.

Effect of primitive matching and refinement. Using a distance function to refine object proposals and aggregating features of matching primitives are crucial for H3DNet. On ScanNet, merely classifying the initial proposals results in a 14.6% drop on mAP 0.5. Figure 4 shows qualitative object detection results, which again justify the importance of optimizing and refining object proposals.

4.3 Ablation Study

Effects of using different geometric primitives. H3DNet can utilize different groups of geometric primitives for generating, classifying, and refining

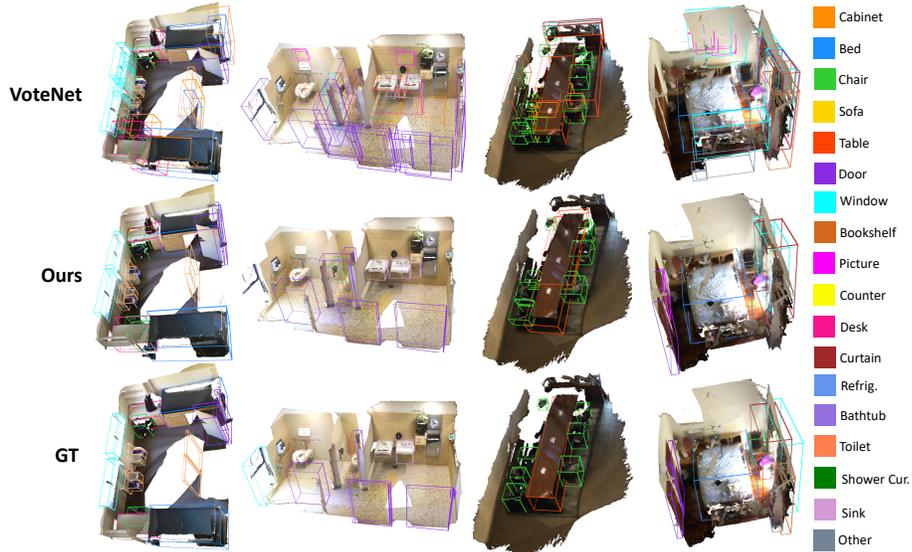


Fig. 5: Qualitative baseline comparisons on ScanNet V2.

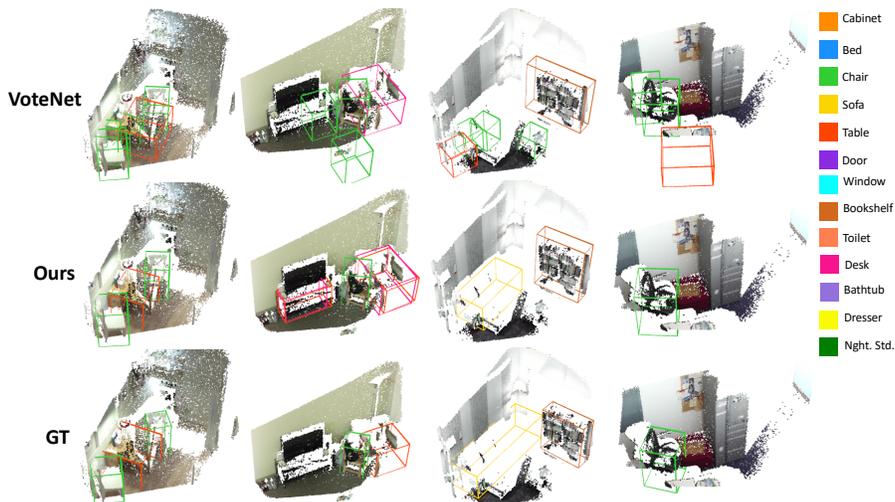


Fig. 6: Qualitative baseline comparisons on SUN RGB-D.

object proposals. Such choices have profound influences on the detected objects. As illustrated in Figure 7, when only using BB edge primitives, we can see that objects with prominent edge features, i.e., window, possess accurate predictions. In contrast, objects with dense face/surface features, such as shower curtain, exhibit relative low prediction accuracy. However, these objects can be easily detected by activating BB face primitives. H3DNet, which combines BB centers, BB edge centers, and BB face centers, adds the strength of their generalization behaviors together. The resulting performance gains are salient when compared to using a single set of geometric primitives.

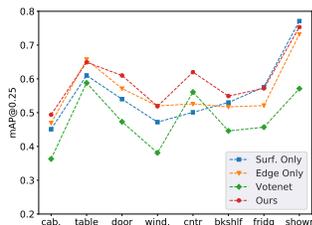


Fig. 7: Quantitative comparisons between VoteNet, our approach, ours with only face primitive and ours with only edge primitive, across sampled categories for ScanNet.

Table 3: Quantitative results without refining predicted center, size, semantic or object existence score for ScanNet, and without refining predicted angle for SUN RGB-D and differences compared with refining all.

	mAP@0.25		mAP@0.5	
w\o center	66.9	-0.3	46.3	-1.8
w\o size	65.4	-1.8	44.2	-3.9
w\o semantic	66.2	-1.0	47.3	-0.8
w\o existence	65.2	-1.8	45.1	-3.0
w\o angle	58.6	-1.5	36.6	-2.4

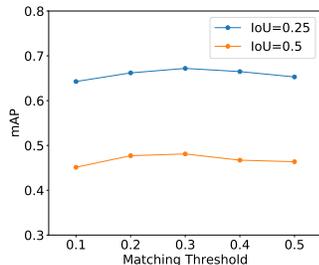


Fig. 8: Quantitative comparisons between different truncation threshold δ for ScanNet.

Table 4: Quantitative comparisons between different number of descriptor computation towers, among our approach and VoteNet, for ScanNet and SUN RGB-D.

	# of Towers	mAP@0.25	mAP@0.5
Ours	1	64.4	43.4
	2	65.4	46.2
	3	66.0	47.7
	4	67.2	48.3
Vote	4 (Scan)	60.11	37.12
	4 (SUN)	57.5	32.1

Effects of proposal refinement. During object proposal refinement, object center, size, heading angle, semantic and existence are all optimized. As shown in Table 3, without fine-tuning any of the geometric parameters of the detected objects, the performance drops, which shows the importance of this sub-module.

Effect of different truncation threshold As shown in Figure 8, with different truncation values of δ , results with mAP@0.25 and mAP@0.5 remain stable. It shows that our model is robust to different truncation threshold δ .

Effect of multiple descriptor computation towers. One hyper-parameter of H3DNet is the number of descriptor computation towers. Table 4 shows that adding more descriptor computation towers leads to better results, yet the performance gain of adding more descriptor computation towers quickly drops. Moreover, the performance gain of H3DNet from VoteNet comes from the hybrid set of geometric primitives and object proposal matching and refinement. For example, replacing the descriptor computation tower of VoteNet by the four descriptor computation towers of H3DNet only results in modest and no performance gains on ScanNet and SUN RGB-D, respectively (See Table 4).

5 Conclusions and Future Work

In this paper, we have introduced a novel 3D object detection approach that takes a 3D scene as input and outputs a collection of labeled and oriented bounding boxes. The key idea of our approach is to predict a hybrid and overcomplete set of geometric primitives and then fit the detected objects to these primitives and their associated features. Experimental results demonstrate the advantages of this approach on ScanNet and SUN RGB-D. In the future, we would like to apply this approach to other 3D scene understanding tasks such as instance segmentation and CAD model reconstruction. Another future direction is to integrate more geometric primitives, like BB corners, for 3D object detection.

Acknowledgement. We would like to acknowledge the support from NSF DMS-1700234, a Gift from Snap Research, and a hardware donation from NVIDIA.

References

1. Akata, Z., Malinowski, M., Fritz, M., Schiele, B.: Multi-cue zero-shot learning with strong supervision. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 59–68 (2016)
2. Atzmon, M., Maron, H., Lipman, Y.: Point convolutional neural networks by extension operators. arXiv preprint arXiv:1803.10091 (2018)
3. Baxter, J.: A bayesian/information theoretic model of learning to learn via multiple task sampling. *Machine learning* **28**(1), 7–39 (1997)
4. Chen, X., Ma, H., Wan, J., Li, B., Xia, T.: Multi-view 3d object detection network for autonomous driving. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1907–1915 (2017)
5. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Niessner, M.: ScanNet: Richly-annotated 3d reconstructions of indoor scenes. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (July 2017)
6. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 605–613 (2017)
7. Graham, B., Engelcke, M., van der Maaten, L.: 3d semantic segmentation with sub-manifold sparse convolutional networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 9224–9232 (2018)
8. Guibas, L.J., Huang, Q., Liang, Z.: A condition number for joint optimization of cycle-consistent networks. In: Advances in Neural Information Processing Systems. pp. 1007–1017 (2019)
9. Hermosilla, P., Ritschel, T., Vázquez, P.P., Vinacua, À., Ropinski, T.: Monte carlo convolution for learning on non-uniformly sampled point clouds. In: SIGGRAPH Asia 2018 Technical Papers. p. 235. ACM (2018)
10. Hou, J., Dai, A., Niessner, M.: 3d-sis: 3d semantic instance segmentation of rgb-d scans. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
11. Hua, B.S., Tran, M.K., Yeung, S.K.: Pointwise convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 984–993 (2018)
12. Kendall, A., Cipolla, R.: Geometric loss functions for camera pose regression with deep learning. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017. pp. 6555–6564 (2017). <https://doi.org/10.1109/CVPR.2017.694>
13. Kendall, A., Gal, Y., Cipolla, R.: Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7482–7491 (2018)
14. Kim, B., Xu, S., Savarese, S.: Accurate localization of 3d objects from RGB-D data using segmentation hypotheses. In: 2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013. pp. 3182–3189 (2013). <https://doi.org/10.1109/CVPR.2013.409>, <https://doi.org/10.1109/CVPR.2013.409>
15. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings (2015), <http://arxiv.org/abs/1412.6980>
16. Klokov, R., Lempitsky, V.: Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 863–872 (2017)

17. Lahoud, J., Ghanem, B.: 2d-driven 3d object detection in RGB-D images. In: IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017. pp. 4632–4640 (2017). <https://doi.org/10.1109/ICCV.2017.495>, <https://doi.org/10.1109/ICCV.2017.495>
18. Lahoud, J., Ghanem, B.: 2d-driven 3d object detection in rgb-d images. In: The IEEE International Conference on Computer Vision (ICCV) (Oct 2017)
19. Lahoud, J., Ghanem, B., Pollefeys, M., Oswald, M.R.: 3d instance segmentation via multi-task metric learning. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 9256–9266 (2019)
20. Lepetit, V., Moreno-Noguer, F., Fua, P.: Epnp: An accurate $O(n)$ solution to the pnp problem. *International Journal of Computer Vision* **81**(2), 155–166 (2009). <https://doi.org/10.1007/s11263-008-0152-6>, <https://doi.org/10.1007/s11263-008-0152-6>
21. Li, Y., Bu, R., Sun, M., Wu, W., Di, X., Chen, B.: Pointcnn: Convolution on x-transformed points. In: *Advances in Neural Information Processing Systems*. pp. 820–830 (2018)
22. Li, Y., Wang, G., Ji, X., Xiang, Y., Fox, D.: Deepim: Deep iterative matching for 6d pose estimation. In: *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VI*. pp. 695–711 (2018). https://doi.org/10.1007/978-3-030-01231-1_42, https://doi.org/10.1007/978-3-030-01231-1_42
23. Liang, M., Yang, B., Chen, Y., Hu, R., Urtasun, R.: Multi-task multi-sensor fusion for 3d object detection. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 7345–7353 (2019)
24. Lin, D., Fidler, S., Urtasun, R.: Holistic scene understanding for 3d object detection with RGBD cameras. In: *IEEE International Conference on Computer Vision, ICCV 2013, Sydney, Australia, December 1-8, 2013*. pp. 1417–1424 (2013). <https://doi.org/10.1109/ICCV.2013.179>, <https://doi.org/10.1109/ICCV.2013.179>
25. Luvizon, D.C., Picard, D., Tabia, H.: 2d/3d pose estimation and action recognition using multitask deep learning. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 5137–5146 (2018)
26. Pavlakos, G., Zhou, X., Chan, A., Derpanis, K.G., Daniilidis, K.: 6-dof object pose from semantic keypoints. In: *2017 IEEE International Conference on Robotics and Automation, ICRA 2017, Singapore, Singapore, May 29 - June 3, 2017*. pp. 2011–2018 (2017). <https://doi.org/10.1109/ICRA.2017.7989233>, <https://doi.org/10.1109/ICRA.2017.7989233>
27. Peng, S., Liu, Y., Huang, Q., Zhou, X., Bao, H.: Pvnet: Pixel-wise voting network for 6dof pose estimation. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. pp. 4561–4570 (2019). <https://doi.org/10.1109/CVPR.2019.00469>
28. Pham, Q.H., Nguyen, T., Hua, B.S., Roig, G., Yeung, S.K.: Jsis3d: Joint semantic-instance segmentation of 3d point clouds with multi-task pointwise networks and multi-value conditional random fields. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 8827–8836 (2019)
29. Qi, C.R., Litany, O., He, K., Guibas, L.J.: Deep hough voting for 3d object detection in point clouds. *arXiv preprint arXiv:1904.09664* (2019)
30. Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3d object detection from rgb-d data. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 918–927 (2018)

31. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 652–660 (2017)
32. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in neural information processing systems. pp. 5099–5108 (2017)
33. Ren, Z., Sudderth, E.B.: Three-dimensional object detection and layout prediction using clouds of oriented gradients. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016)
34. Ruder, S.: An overview of multi-task learning in deep neural networks. CoRR **abs/1706.05098** (2017), <http://arxiv.org/abs/1706.05098>
35. Sener, O., Koltun, V.: Multi-task learning as multi-objective optimization. In: Advances in Neural Information Processing Systems. pp. 527–538 (2018)
36. Shi, S., Wang, X., Li, H.: Pointcnn: 3d object proposal generation and detection from point cloud. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 770–779 (2019)
37. Song, C., Song, J., Huang, Q.: Hybridpose: 6d object pose estimation under hybrid representations. CoRR **abs/2001.01869** (2020), <http://arxiv.org/abs/2001.01869>
38. Song, S., Lichtenberg, S.P., Xiao, J.: Sun rgb-d: A rgb-d scene understanding benchmark suite. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 567–576 (2015)
39. Song, S., Xiao, J.: Sliding shapes for 3d object detection in depth images. In: ECCV (2014)
40. Song, S., Xiao, J.: Deep sliding shapes for amodal 3d object detection in rgb-d images. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016)
41. Su, H., Jampani, V., Sun, D., Maji, S., Kalogerakis, E., Yang, M.H., Kautz, J.: Splatnet: Sparse lattice networks for point cloud processing. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2530–2539 (2018)
42. Tatarchenko, M., Dosovitskiy, A., Brox, T.: Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2088–2096 (2017)
43. Tatarchenko, M., Park, J., Koltun, V., Zhou, Q.Y.: Tangent convolutions for dense prediction in 3d. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3887–3896 (2018)
44. Wang, N., Zhou, W., Tian, Q., Hong, R., Wang, M., Li, H.: Multi-cue correlation filters for robust visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4844–4853 (2018)
45. Wang, P.S., Liu, Y., Guo, Y.X., Sun, C.Y., Tong, X.: O-cnn: Octree-based convolutional neural networks for 3d shape analysis. ACM Transactions on Graphics (TOG) **36**(4), 72 (2017)
46. Wang, X., Liu, S., Shen, X., Shen, C., Jia, J.: Associatively segmenting instances and semantics in point clouds. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4096–4105 (2019)
47. Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M.: Dynamic graph cnn for learning on point clouds. arXiv preprint arXiv:1801.07829 (2018)
48. Xie, S., Liu, S., Chen, Z., Tu, Z.: Attentional shapecontextnet for point cloud recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4606–4615 (2018)

49. Xu, Y., Fan, T., Xu, M., Zeng, L., Qiao, Y.: Spidercnn: Deep learning on point sets with parameterized convolutional filters. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 87–102 (2018)
50. Yang, Y., Feng, C., Shen, Y., Tian, D.: Foldingnet: Point cloud auto-encoder via deep grid deformation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 206–215 (2018)
51. Yang, Z., Yan, S., Huang, Q.: Extreme relative pose network under hybrid representations. CoRR **abs/1912.11695** (2019), <http://arxiv.org/abs/1912.11695>
52. Yi, L., Zhao, W., Wang, H., Sung, M., Guibas, L.J.: Gspn: Generative shape proposal network for 3d instance segmentation in point cloud. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3947–3956 (2019)
53. Zhang, Y., Yang, Q.: A survey on multi-task learning. CoRR **abs/1707.08114** (2017), <http://arxiv.org/abs/1707.08114>
54. Zhang, Z., Liang, Z., Wu, L., Zhou, X., Huang, Q.: Path-invariant map networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 11084–11094 (June 2019)
55. Zhou, Y., Tuzel, O.: Voxelnet: End-to-end learning for point cloud based 3d object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4490–4499 (2018)
56. Zou, Y., Luo, Z., Huang, J.B.: Df-net: Unsupervised joint learning of depth and flow using cross-task consistency. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 36–53 (2018)