

Deep Hough Voting for 3D Object Detection in Point Clouds

Charles R. Qi¹ Or Litany¹ Kaiming He¹ Leonidas J. Guibas^{1,2}

¹Facebook AI Research ²Stanford University

Abstract

Current 3D object detection methods are heavily influenced by 2D detectors. In order to leverage architectures in 2D detectors, they often convert 3D point clouds to regular grids (i.e., to voxel grids or to bird’s eye view images), or rely on detection in 2D images to propose 3D boxes. Few works have attempted to directly detect objects in point clouds. In this work, we return to first principles to construct a 3D detection pipeline for point cloud data and as generic as possible. However, due to the sparse nature of the data – samples from 2D manifolds in 3D space – we face a major challenge when directly predicting bounding box parameters from scene points: a 3D object centroid can be far from any surface point thus hard to regress accurately in one step. To address the challenge, we propose VoteNet, an end-to-end 3D object detection network based on a synergy of deep point set networks and Hough voting. Our model achieves state-of-the-art 3D detection on two large datasets of real 3D scans, ScanNet and SUN RGB-D with a simple design, compact model size and high efficiency. Remarkably, VoteNet outperforms previous methods by using purely geometric information without relying on color images.

1. Introduction

The goal of 3D object detection is to localize and recognize objects in a 3D scene. More specifically, in this work, we aim to estimate oriented 3D bounding boxes as well as semantic classes of objects from point clouds.

Compared to images, 3D point clouds provide accurate geometry and robustness to illumination changes. On the other hand, point clouds are irregular, thus typical CNNs are not well suited to process them directly.

To avoid processing irregular point clouds, current 3D detection methods heavily rely on 2D-based detectors in various aspects. For example, [42, 12] extend 2D detection frameworks such as the Faster/Mask R-CNN [37, 11] to 3D. They voxelize the irregular point clouds to regular 3D grids and apply 3D CNN detectors, which fails to leverage sparsity in the data and suffer from high computation cost due to expensive 3D convolutions. Alternatively, [4, 55] project

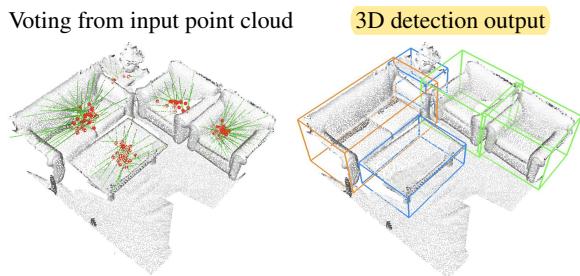


Figure 1. **3D object detection in point clouds with a deep Hough voting model.** Given a point cloud of a 3D scene, our VoteNet votes to object centers and then groups and aggregates the votes to predict 3D bounding boxes and semantic classes of objects. Our code is open sourced at <https://github.com/facebookresearch/votenet>

points to regular 2D bird’s eye view images and then apply 2D detectors to localize objects. This, however, sacrifices geometric details which may be critical in cluttered indoor environments. More recently, [20, 34] proposed a cascaded two-step pipeline by firstly detecting objects in front-view images and then localizing objects in frustum point clouds extruded from the 2D boxes, which however is strictly dependent on the 2D detector and will miss an object entirely if it is not detected in 2D.

In this work we introduce a *point cloud focused* 3D detection framework that directly processes raw data and does not depend on any 2D detectors neither in architecture nor in object proposal. Our detection network, *VoteNet*, is based on recent advances in 3D deep learning models for point clouds, and is inspired by the generalized Hough voting process for object detection [23].

We leverage PointNet++ [36], a hierarchical deep network for point cloud learning, to mitigate the need to convert point clouds to regular structures. By directly processing point clouds not only do we avoid information loss by a quantization process, but we also take advantage of the sparsity in point clouds by only computing on sensed points.

While PointNet++ has shown success in object classification and semantic segmentation [36], few research study how to detect 3D objects in point clouds with such architectures. A naïve solution would be to follow common practice

in 2D detectors and perform dense object proposal [29, 37], i.e. to propose 3D bounding boxes directly from the sensed points (with their learned features). However, the inherent sparsity of point clouds makes this approach unfavorable. In images there often exists a pixel near the object center, but it is often not the case in point clouds. As depth sensors only capture surfaces of objects, 3D object centers are likely to be in empty space, far away from any point. As a result, point based networks have difficulty aggregating scene context in the vicinity of object centers. Simply increasing the receptive field does not solve the problem because as the network captures larger context, it also causes more inclusion of nearby objects and clutter.

To this end, we propose to endow point cloud deep networks with a *voting mechanism similar to the classical Hough voting*. By *voting* we essentially generate new points that lie close to objects centers, which can be *grouped and aggregated* to generate box *proposals*.

In contrast to traditional Hough voting with multiple separate modules that are difficult to optimize jointly, *VoteNet* is end-to-end optimizable. Specifically, after passing the input point cloud through a backbone network, we sample a set of seed points and generate votes from their features. Votes are targeted to reach object centers. As a result, vote clusters emerge near object centers and in turn can be aggregated through a learned module to generate box proposals. The result is a powerful 3D object detector that is purely geometric and can be applied directly to point clouds.

We evaluate our approach on two challenging 3D object detection datasets: SUN RGB-D [40] and ScanNet [5]. On both datasets VoteNet, *using geometry only*, significantly outperforms prior arts that use both RGB and geometry or even multi-view RGB images. Our study shows that the voting scheme supports more effective context aggregation, and verifies that VoteNet offers the largest improvements when object centers are far from the object surface (e.g. tables, bathtubs, etc.).

In summary, the contributions of our work are:

- A reformulation of Hough voting in the context of deep learning through an end-to-end differentiable architecture, which we dub VoteNet.
- State-of-the-art 3D object detection performance on SUN RGB-D and ScanNet.
- An in-depth analysis of the importance of voting for 3D object detection in point clouds.

2. Related Work

3D object detection. Many previous methods were proposed to detect 3D bounding boxes of objects. Examples include: [27] where a pair-wise semantic context potential helps guide the proposals objectness score; template-

based methods [26, 32, 28]; Sliding-shapes [41] and its deep learning-based successor [42]; Clouds of Oriented Gradients (COG) [38]; and the recent 3D-SIS [12].

Due to the complexity of directly working in 3D, especially in large scenes, many methods resort to some type of projection. For example in MV3D [4] and VoxelNet [55], the 3D data is first reduced to a bird’s-eye view before proceeding to the rest of the pipeline. A reduction in search space by first processing a 2D input was demonstrated in both Frustum PointNets [34] and [20]. Similarly, in [16] a segmentation hypothesis is verified using the 3D map. More recently, deep networks on point clouds are used to exploit sparsity of the data by GSPN [54] and PointRCNN [39].

Hough voting for object detection. Originally introduced in the late 1950s, the Hough transform [13] translates the problem of detecting simple patterns in point samples to detecting peaks in a parametric space. The Generalized Hough Transform [2] further extends this technique to image patches as indicators for the existence of a complex object. Examples of using Hough voting include the seminal work of [24] which introduced the implicit shape model, planes extraction from 3D point clouds [3], and 6D pose estimation [44] to name a few.

Hough voting has also been previously combined with advanced learning techniques. In [30] the votes were assigned with weights indicating their importance, which were learned using a max-margin framework. Hough forests for object detection were introduced in [8, 7]. More recently, [15] demonstrated improved voting-based 6D pose estimation by using deep features extracted to build a code-book. Similarly [31] learned deep features to build code-books for segmentation in MRI and ultrasounds images. In [14] the classical Hough algorithm was used to extract circular patterns in car logos, which were then input to a deep classification network. [33] proposed the semi-convolutional operator for 2D instance segmentation in images, which is also related to Hough voting.

There have also been works using Hough voting for 3D object detection [50, 18, 47, 19], which adopted a similar pipeline as in 2D detectors.

Deep learning on point clouds. Recently we see a surge of interest in designing deep network architectures suited for point clouds [35, 36, 43, 1, 25, 9, 48, 45, 46, 22, 17, 53, 52, 49, 51], which showed remarkable performance in 3D object classification, object part segmentation, as well as scene segmentation. In the context of 3D object detection, VoxelNet [55] learn voxel feature embeddings from points in voxels, while in [34] PointNets are used to localize objects in a frustum point cloud extruded from a 2D bounding box. However, few methods studied how to directly propose and detect 3D objects in raw point cloud representation.

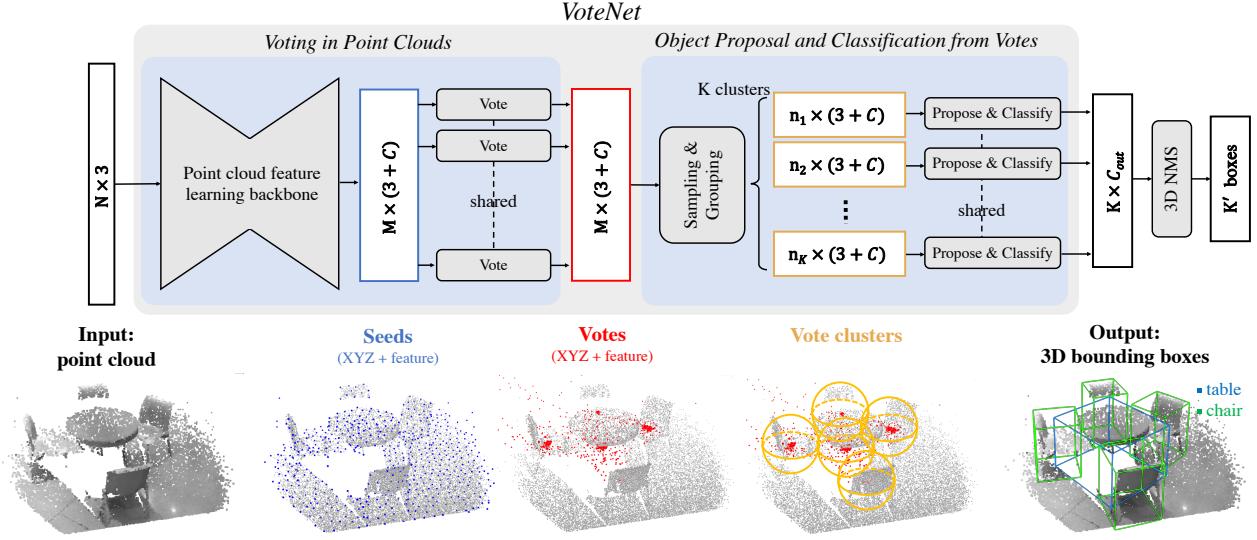


Figure 2. **Illustration of the VoteNet architecture** for 3D object detection in point clouds. Given an input point cloud of N points with XYZ coordinates, a backbone network (implemented with PointNet++ [36] layers) subsamples and learns deep features on the points and outputs a subset of M points but extended by C -dim features. This subset of points are considered as seed points. Each seed independently generates a vote through a voting module. Then the votes are grouped into clusters and processed by the proposal module to generate the final proposals. The classified and NMSed proposals become the final 3D bounding boxes output. Image best viewed in color.

3. Deep Hough Voting

A traditional Hough voting 2D detector [24] comprises an offline and an online step. First, given a collection of images with annotated object bounding boxes, a codebook is constructed with stored mappings between image patches (or their features) and their offsets to the corresponding object centers. At inference time, interest points are selected from the image to extract patches around them. These patches are then compared against patches in the codebook to retrieve offsets and compute votes. As *object* patches will tend to vote in agreement, clusters will form near object centers. Finally, the object boundaries are retrieved by tracing cluster votes back to their generating patches.

We identify two ways in which this technique is well suited to our problem of interest. First, voting-based detection is more compatible with sparse sets than region-proposal networks (RPN) [37] is. For the latter, the RPN has to generate a proposal near an object center which is likely to be in an empty space, causing extra computation. Second, it is based on a bottom-up principle where small bits of partial information are accumulated to form a confident detection. Even though neural networks can potentially aggregate context from a large receptive field, it may be still beneficial to aggregate in the vote space.

However, as traditional Hough voting comprises multiple separated modules, integrating it into state-of-the-art point cloud networks is an open research topic. To this end, we propose the following adaptations to the different pipeline ingredients.

Interest points are described and selected by deep neural networks instead of depending on hand-crafted features.

Vote generation is learned by a network instead of using a codebook. Leveraging larger receptive fields, voting can be made less ambiguous and thus more effective. In addition, a vote *location* can be augmented with a *feature vector* allowing for better aggregation.

Vote aggregation is realized through point cloud processing layers with trainable parameters. Utilizing the vote features, the network can potentially filter out low quality votes and generate improved proposals.

Object proposals in the form of: location, dimensions, orientation and even semantic classes can be directly generated from the aggregated features, mitigating the need to trace back votes’ origins.

In what follows, we describe how to combine all the aforementioned ingredients into a single end-to-end trainable network named as VoteNet.

4. VoteNet Architecture

Fig. 2 illustrates our end-to-end detection network (*VoteNet*). The entire network can be split into two parts: one that processes *existing* points to generate votes; and the other part that operates on *virtual* points – the votes – to propose and classify objects.

4.1. Learning to Vote in Point Clouds

From an input point cloud of size $N \times 3$, with a 3D coordinate for each of the N points, we aim to generate M votes,

where each vote has both a 3D coordinate and a high dimensional feature vector. There are two major steps: point cloud feature learning through a backbone network and learned Hough voting from seed points.

Point cloud feature learning. Generating an accurate vote requires geometric reasoning and contexts. Instead of relying on hand-crafted features, we leverage recently proposed deep networks [36, 9, 43, 25] on point clouds for point feature learning. While our method is not restricted to any point cloud network, we adopt PointNet++ [36] as our backbone due to its simplicity and demonstrated success on tasks ranging from normal estimation [10], semantic segmentation [21] to 3D object localization [34].

The backbone network has several set-abstraction layers and feature propagation (upsampling) layers with skip connections, which outputs a subset of the input points with XYZ and an enriched C -dimensional feature vector. The results are M *seed points* of dimension $(3 + C)$. Each seed point generates one vote¹.

Hough voting with deep networks. Compared to traditional Hough voting where the votes (offsets from local keypoints) are determined by look ups in a pre-computed codebook, we generate votes with a deep network based voting module, which is both more efficient (without kNN look ups) and more accurate as it is trained jointly with the rest of the pipeline.

Given a set of seed points $\{s_i\}_{i=1}^M$ where $s_i = [x_i; f_i]$ with $x_i \in \mathbb{R}^3$ and $f_i \in \mathbb{R}^C$, a shared *voting module* generates votes from each seed independently. Specifically, the voting module is realized with a multi-layer perceptron (MLP) network with fully connected layers, ReLU and batch normalization. The MLP takes seed feature f_i and outputs the Euclidean space offset $\Delta x_i \in \mathbb{R}^3$ and a feature offset $\Delta f_i \in \mathbb{R}^C$ such that the vote $v_i = [y_i; g_i]$ generated from the seed s_i has $y_i = x_i + \Delta x_i$ and $g_i = f_i + \Delta f_i$.

The predicted 3D offset Δx_i is explicitly supervised by a regression loss

$$L_{\text{vote-reg}} = \frac{1}{M_{\text{pos}}} \sum_i \|\Delta x_i - \Delta x_i^*\| \mathbf{1}[s_i \text{ on object}], \quad (1)$$

where $\mathbf{1}[s_i \text{ on object}]$ indicates whether a seed point s_i is on an object surface and M_{pos} is the count of total number of seeds on object surface. Δx_i^* is the ground truth displacement from the seed position x_i to the bounding box center of the object it belongs to.

Votes are the same as seeds in tensor representation but are no longer grounded on object surfaces. A more essential difference though is their position – votes generated from seeds on the same object are now closer to each other than the seeds are, which makes it easier to combine cues from different parts of the object. Next we will take advantage of

¹The case of more than one vote is discussed in the appendix.

this semantic-aware locality to aggregate vote features for object proposal.

4.2. Object Proposal and Classification from Votes

The votes create canonical “meeting points” for context aggregation from different parts of the objects. After clustering these votes we aggregate their features to generate object proposals and classify them.

Vote clustering through sampling and grouping. While there can be many ways to cluster the votes, we opt for a simple strategy of uniform sampling and grouping according to spatial proximity. Specifically, from a set of votes $\{v_i = [y_i; g_i] \in \mathbb{R}^{3+C}\}_{i=1}^M$, we sample a subset of K votes using farthest point sampling based on $\{y_i\}$ in 3D Euclidean space, to get $\{v_{i_k}\}$ with $k = 1, \dots, K$. Then we form K clusters by finding neighboring votes to each of the v_{i_k} ’s 3D location: $\mathcal{C}_k = \{v_i^{(k)} | \|v_i - v_{i_k}\| \leq r\}$ for $k = 1, \dots, K$. Though simple, this clustering technique is easy to integrate into an end-to-end pipeline and works well in practice.

Proposal and classification from vote clusters. As a vote cluster is in essence a set of high-dim points, we can leverage a generic point set learning network to aggregate the votes in order to generate object proposals. Compared to the back-tracing step of traditional Hough voting for identifying the object boundary, this procedure allows to propose *amodal* boundaries even from partial observations, as well as predicting other parameters like orientation, class, etc.

In our implementation, we use a *shared* PointNet [35] for vote aggregation and proposal in clusters. Given a vote cluster $\mathcal{C} = \{w_i\}$ with $i = 1, \dots, n$ and its cluster center w_j , where $w_i = [z_i; h_i]$ with $z_i \in \mathbb{R}^3$ as the vote location and $h_i \in \mathbb{R}^C$ as the vote feature. To enable usage of local vote geometry, we transform vote locations to a local normalized coordinate system by $z'_i = (z_i - z_j)/r$. Then an object proposal for this cluster $p(\mathcal{C})$ is generated by passing the set input through a PointNet-like module:

$$p(\mathcal{C}) = \text{MLP}_2 \left\{ \max_{i=1, \dots, n} \{\text{MLP}_1([z'_i; h_i])\} \right\} \quad (2)$$

where votes from each cluster are independently processed by a MLP_1 before being max-pooled (channel-wise) to a single feature vector and passed to MLP_2 where information from different votes are further combined. We represent the proposal p as a multidimensional vector with an objectness score, bounding box parameters (center, heading and scale parameterized as in [34]) and semantic classification scores.

Loss function. The loss functions in the proposal and classification stage consist of objectness, bounding box estimation, and semantic classification losses.

We supervise the objectness scores for votes that are located either close to a ground truth object center (within

0.3 meters) or far from any center (by more than 0.6 meters). We consider proposals generated from those votes as *positive* and *negative* proposals, respectively. Objectness predictions for other proposals are not penalized. Objectness is supervised via a cross entropy loss normalized by the number of non-ignored proposals in the batch. For positive proposals we further supervise the bounding box estimation and class prediction according to the closest ground truth bounding box. Specifically, we follow [34] which decouples the box loss to center regression, heading angle estimation and box size estimation. For semantic classification we use the standard cross entropy loss. In all regression in the detection loss we use the Huber (smooth- L_1 [37]) loss. Further details are provided in the appendix.

4.3. Implementation Details

Input and data augmentation. Input to our detection network is a point cloud of N points randomly sub-sampled from either a popped-up depth image ($N = 20k$) or a 3D scan (mesh vertices, $N = 40k$). In addition to XYZ coordinates, we also include a height feature for each point indicating its distance to the floor. The floor height is estimated as the 1% percentile of all points' heights. To augment the training data, we randomly sub-sample the points from the scene points on-the-fly. We also randomly flip the point cloud in both horizontal direction, randomly rotate the scene points by Uniform[$-5^\circ, 5^\circ$] around the upright-axis, and randomly scale the points by Uniform[0.9, 1.1].

Network architecture details. The backbone feature learning network is based on PointNet++ [36], which has four set abstraction (SA) layers and two feature propagation/upsampling (FP) layers, where the SA layers have increasing receptive radius of 0.2, 0.4, 0.8 and 1.2 in meters while they sub-sample the input to 2048, 1024, 512 and 256 points respectively. The two FP layers up-sample the 4th SA layer's output back to 1024 points with 256-dim features and 3D coordinates (more details in the appendix).

The voting layer is realized through a multi-layer perceptron with FC output sizes of 256, 256, 259, where the last FC layer outputs XYZ offset and feature residuals.

The proposal module is implemented as a set abstraction layer with a post processing MLP₂ to generate proposals after the max-pooling. The SA uses radius 0.3 and MLP₁ with output sizes of 128, 128, 128. The max-pooled features are further processed by MLP₂ with output sizes of 128, 128, $5 + 2NH + 4NS + NC$ where the output consists of 2 objectness scores, 3 center regression values, $2NH$ numbers for heading regression (NH heading bins) and $4NS$ numbers for box size regression (NS box anchors) and NC numbers for semantic classification.

Training the network. We train the entire network end-to-end and from scratch with an Adam optimizer, batch size 8 and an initial learning rate of 0.001. The learning rate is

decreased by $10\times$ after 80 epochs and then decreased by another $10\times$ after 120 epochs. Training the model to convergence on one Volta Quadro GP100 GPU takes around 10 hours on SUN RGB-D and less than 4 hours on ScanNetV2.

Inference. Our VoteNet is able to take point clouds of the entire scenes and generate proposals in one forward pass. The proposals are post-processed by a 3D NMS module with an IoU threshold of 0.25. The evaluation follows the same protocol as in [42] using mean average precision.

5. Experiments

In this section, we firstly compare our Hough voting based detector with previous state-of-the-art methods on two large 3D indoor object detection benchmarks (Sec. 5.1). We then provide analysis experiments to understand the importance of voting, the effects of different vote aggregation approaches and show our method's advantages in its compactness and efficiency (Sec. 5.2). Finally we show qualitative results of our detector (Sec. 5.3). More analysis and visualizations are provided in the appendix.

5.1. Comparing with State-of-the-art Methods

Dataset. SUN RGB-D [40] is a single-view RGB-D dataset for 3D scene understanding. It consists of $\sim 5K$ RGB-D training images annotated with amodal oriented 3D bounding boxes for 37 object categories. To feed the data to our network, we firstly convert the depth images to point clouds using the provided camera parameters. We follow a standard evaluation protocol and report performance on the 10 most common categories.

ScanNetV2 [5] is a richly annotated dataset of 3D reconstructed meshes of indoor scenes. It contains $\sim 1.2K$ training examples collected from hundreds of different rooms, and annotated with semantic and instance segmentation for 18 object categories. Compared to partial scans in SUN RGB-D, scenes in ScanNetV2 are more complete and cover larger areas with more objects on average. We sample vertices from the reconstructed meshes as our input point clouds. Since ScanNetV2 does not provide amodal or oriented bounding box annotation, we aim to predict axis-aligned bounding boxes instead, as in [12].

Methods in comparison. We compare with a wide range of prior art methods. Deep sliding shapes (DSS) [42] and 3D-SIS [12] are both 3D CNN based detectors that combine geometry and RGB cues in object proposal and classification, based on the Faster R-CNN [37] pipeline. Compared with DSS, 3D-SIS introduces a more sophisticated sensor fusion scheme (back-projecting RGB features to 3D voxels) and therefore is able to use multiple RGB views to improve performance. 2D-driven [20] and F-PointNet [34] are 2D-based 3D detectors that rely on object detection in 2D images to reduce the 3D detection search space. Cloud of gra-

	Input	bathtub	bed	bookshelf	chair	desk	dresser	nightstand	sofa	table	toilet	mAP
DSS [42]	Geo + RGB	44.2	78.8	11.9	61.2	20.5	6.4	15.4	53.5	50.3	78.9	42.1
COG [38]	Geo + RGB	58.3	63.7	31.8	62.2	45.2	15.5	27.4	51.0	51.3	70.1	47.6
2D-driven [20]	Geo + RGB	43.5	64.5	31.4	48.3	27.9	25.9	41.9	50.4	37.0	80.4	45.1
F-PointNet [34]	Geo + RGB	43.3	81.1	33.3	64.2	24.7	32.0	58.1	61.1	51.1	90.9	54.0
VoteNet (ours)	Geo only	74.4	83.0	28.8	75.3	22.0	29.8	62.2	64.0	47.3	90.1	57.7

Table 1. **3D object detection results on SUN RGB-D val set.** Evaluation metric is average precision with 3D IoU threshold 0.25 as proposed by [40]. Note that both COG [38] and 2D-driven [20] use room layout context to boost performance. To have fair comparison with previous methods, the evaluation is on the SUN RGB-D V1 data.

	Input	mAP@0.25	mAP@0.5
DSS [42, 12]	Geo + RGB	15.2	6.8
MRCNN 2D-3D [11, 12]	Geo + RGB	17.3	10.5
F-PointNet [34, 12]	Geo + RGB	19.8	10.8
GSPN [54]	Geo + RGB	30.6	17.7
3D-SIS [12]	Geo + 1 view	35.1	18.7
3D-SIS [12]	Geo + 3 views	36.6	19.0
3D-SIS [12]	Geo + 5 views	40.2	22.5
3D-SIS [12]	Geo only	25.4	14.6
VoteNet (ours)	Geo only	58.6	33.5

Table 2. **3D object detection results on ScanNetV2 val set.** DSS and F-PointNet results are from [12]. Mask R-CNN 2D-3D results are from [54]. GSPN and 3D-SIS results are up-to-date numbers provided by the original authors.

dients [38] is a sliding window based detector using a newly designed 3D HoG-like feature. MRCNN 2D-3D is a naïve baseline that directly projects Mask-RCNN [11] instance segmentation results into 3D to get a bounding box estimation. GSPN [54] is a recent instance segmentation method using a generative model to propose object instances, which is also based on a PointNet++ backbone.

Results are summarized in Table 1 and 2. VoteNet outperforms all previous methods by at least **3.7** and **18.4** mAP increase in SUN RGB-D and ScanNet respectively. Notably, we achieve such improvements when we *use geometric input (point clouds) only* while they used both geometry and RGB images. Table 1 shows that in the category “chair” with the most training samples, our method improves upon previous state of the art by more than **11 AP**. Table 2 shows that when taking geometric input only, our method outperforms 3D CNN based method 3D-SIS by more than **33 AP**. A per-category evaluation for ScanNet is provided in the appendix. Importantly, the same set of network hyperparameters was used in both datasets.

5.2. Analysis Experiments

To Vote or Not To Vote? A straightforward baseline to VoteNet is a network that directly proposes boxes from sam-

Method	mAP@0.25	
	SUN RGB-D	ScanNet
BoxNet (ours)	53.0	45.4
VoteNet (ours)	57.7	58.6

Table 3. **Comparing VoteNet with a no-vote baseline.** Metric is 3D object detection mAP. VoteNet estimate object bounding boxes from vote clusters. BoxNet proposes boxes directly from seed points on object surfaces without voting.

pled scene points. Such a baseline – which we refer to as *BoxNet* – is essential to distill the improvement due to voting. The BoxNet has the same backbone as the VoteNet but instead of voting, it directly generates boxes from the seed points (more details in appendix). Table 3 shows voting boosts the performance by a significant margin of ~ 5 mAP on SUN RGB-D and >13 mAP on ScanNet.

In what ways, then, does voting help? We argue that since in sparse 3D point clouds, existing scene points are often far from object centroids, a direct proposal may have lower confidence and inaccurate amodal boxes. Voting, instead, brings closer together these lower confidence points and allows to reinforce their hypothesis through aggregation. We demonstrate this phenomenon in Fig. 3 on a typical ScanNetV2 scene. We overlay the scene with only those seed points which, if sampled, would generate an accurate proposal. As can be seen, VoteNet (right) offers a much broader coverage of “good” seed points compared to BoxNet (left), showing its robustness brought by voting.

We proceed with a second analysis in Fig. 4 showing on the same plot (in separate scales), for each SUN RGB-D category: (in blue dots) gains in mAP between VoteNet and BoxNet, and (in red squares) closest distances between object points (on their surfaces) and their amodal box centers, averaged per category and normalized by the mean class size (a large distance means the object center is usually far from its surface). Sorting the categories according to the former, we see a strong correlation. Namely, when object points tend to be further from the amodal box center, voting helps much more.

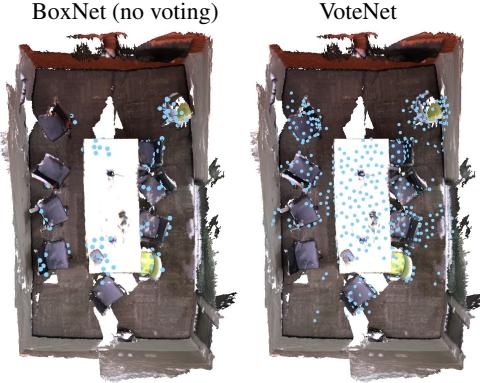


Figure 3. Voting helps increase detection contexts. Seed points that generate good boxes (BoxNet), or good votes (VoteNet) which in turn generate good boxes, are overlaid (in blue) on top of a representative ScanNet scene. As the voting step effectively increases context, VoteNet demonstrates a much denser cover of the scene, therefore increasing the likelihood of accurate detection.

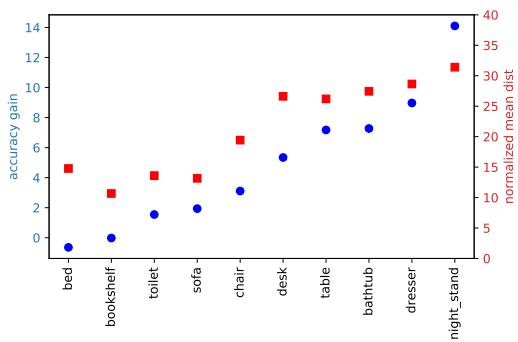


Figure 4. Voting helps more in cases where object points are far from object centers. We show for each category: voting accuracy gain (in blue dots) of VoteNet w.r.t our direct proposal baseline BoxNet; and (in red squares) average object-center distance, normalized by the mean class size.

Effect of Vote Aggregation Aggregation of votes is an important component in VoteNet as it allows communication between votes. Hence, it is useful to analyze how different aggregation schemes influence performance.

In Fig. 5 (right), we show that vote aggregation with a learned Pointnet and max pooling achieves far better results than manually aggregating the vote features in the local regions due to the existence of clutter votes (i.e. votes from non-object seeds). We test 3 types of those aggregations (first three rows): max, average, and RBF weighting (based on vote distance to the cluster center). In contrast to aggregation with Pointnet (Eq. 2), the vote features are directly pooled, e.g. for avg. pooling: $p = \text{MLP}_2 \{\text{AVG}\{h_i\}\}$.

In Fig. 5 (left), we show how vote aggregation radius affects detection (tested with Pointnet using max pooling). As the aggregation radius increases, VoteNet improves until it peaks at around 0.2 radius. Attending to a larger region

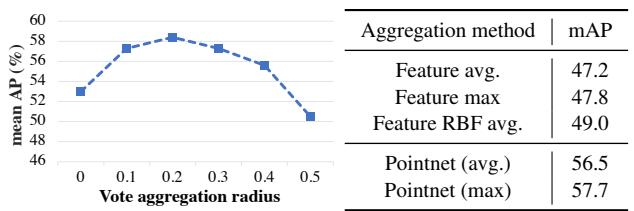


Figure 5. Vote aggregation analysis. *Left:* mAP@0.25 on SUN RGB-D for varying aggregation radii when aggregating via Pointnet (max). *Right:* Comparisons of different aggregation methods (radius = 0.3 for all methods). Using a learned vote aggregation is far more effective than manually pooling the features in a local neighborhood.

though introduces more clutter votes thus contaminating the good votes and results in decreased performance.

Model Size and Speed Our proposed model is very efficient since it leverages sparsity in point clouds and avoids search in empty space. Compared to previous best methods (Table 4), our model is more than $4\times$ smaller than F-PointNet (the prior art on SUN RGB-D) in size and more than $20\times$ times faster than 3D-SIS (the prior art on ScanNetV2) in speed. Note that the ScanNetV2 processing time by 3D-SIS is computed as averaged time in offline batch mode while ours is measured with sequential processing which can be realized in online applications.

5.3. Qualitative Results and Discussion

Fig. 6 and Fig. 7 show several representative examples of VoteNet detection results on ScanNet and SUN RGB-D scenes, respectively. As can be seen, the scenes are quite diverse and pose multiple challenges including clutter, partiality, scanning artifacts, etc. Despite these challenges, our network demonstrates quite robust results. See for example in Fig. 6, how the vast majority of chairs were correctly detected in the top scene. Our method was able to nicely distinguish between the attached sofa-chairs and the sofa in the bottom left scene; and predicted the complete bounding box of the much fragmented and cluttered desk at the bottom right scene.

There are still limitations in our method though. Com-

Method	Model size	SUN RGB-D	ScanNetV2
F-PointNet [34]	47.0MB	0.09s	-
3D-SIS [12]	19.7MB	-	2.85s
VoteNet (ours)	11.2MB	0.10s	0.14s

Table 4. Model size and processing time (per frame or scan). Our method is more than $4\times$ more compact in model size than [34] and more than $20\times$ faster than [12].

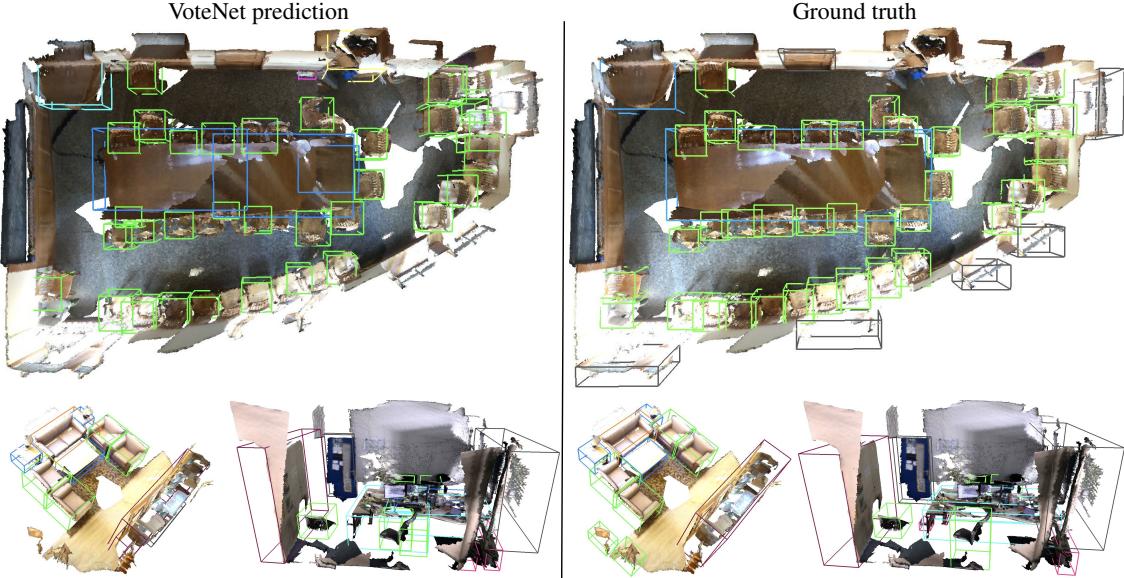


Figure 6. **Qualitative results of 3D object detection in ScanNetV2.** Left: our VoteNet, Right: ground-truth. See Section 5.3 for details.

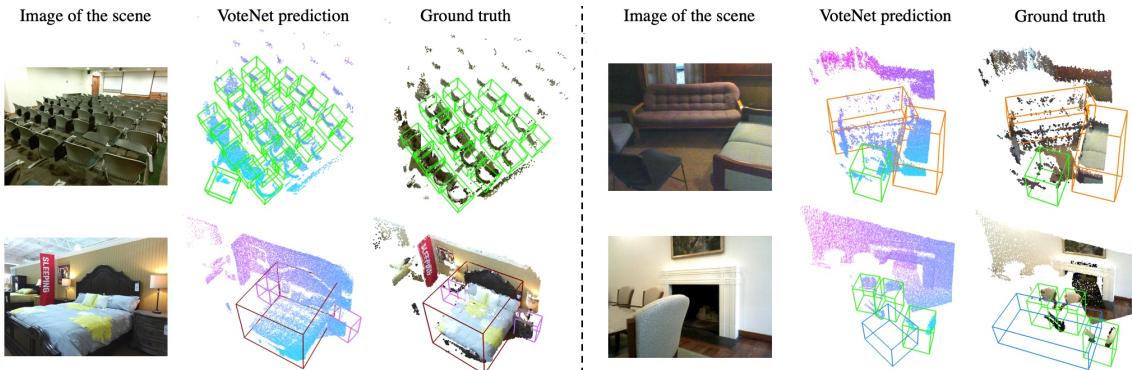


Figure 7. **Qualitative results on SUN RGB-D.** Both left and right panels show (from left to right): an image of the scene (not used by our network), 3D object detection by VoteNet, and ground-truth annotations. See Section 5.3 for details.

mon failure cases include misses on very thin objects like doors, windows and pictures denoted in black bounding boxes in the top scene (Fig. 6). As we do not make use of RGB information, detecting these categories is almost impossible. Fig. 7 on SUN RGB-D also reveals the strengths of our method in partial scans with single-view depth images. For example, it detected more chairs in the top-left scene than were provided by the ground-truth. In the top-right scene we can see how VoteNet can nicely hallucinate the amodal bounding box despite seeing only part of the sofa. A less successful amodal prediction is shown in the bottom right scene where an extremely partial observation of a very large table is given.

6. Conclusion

In this work we have introduced VoteNet: a simple, yet powerful 3D object detection model inspired by Hough voting. The network learns to vote to object centroids directly

from point clouds and learns to aggregate votes through their features and local geometry to generate high-quality object proposals. Using only 3D point clouds, the model showed significant improvements over previous methods that utilize both depth and colored images.

In future work we intend to explore how to incorporate RGB images into our detection framework and to utilize our detector in downstream application such as 3D instance segmentation. We believe that the synergy of Hough voting and deep learning can be generalized to more applications such as 6D pose estimation, template based detection etc. and expect to see more future research along this line.

Acknowledgements. This work was supported in part by ONR MURI grant N00014-13-1-0341, NSF grant IIS-1763268 and a Vannevar Bush Faculty Fellowship. We thank Daniel Huber, Justin Johnson, Georgia Gkioxari and Jitendra Malik for valuable discussions and feedback.

References

- [1] Matan Atzmon, Haggai Maron, and Yaron Lipman. Point convolutional neural networks by extension operators. *arXiv preprint arXiv:1803.10091*, 2018. [2](#)
- [2] Dana H Ballard. Generalizing the hough transform to detect arbitrary shapes. *Pattern recognition*, 13(2):111–122, 1981. [2](#)
- [3] Dorit Borrmann, Jan Elseberg, Kai Lingemann, and Andreas Nüchter. The 3d hough transform for plane detection in point clouds: A review and a new accumulator design. *3D Research*, 2(2):3, 2011. [2](#)
- [4] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *IEEE CVPR*, 2017. [1, 2](#)
- [5] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5828–5839, 2017. [2, 5](#)
- [6] Haoqiang Fan, Hao Su, and Leonidas J Guibas. 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*, pages 605–613, 2017. [11](#)
- [7] Juergen Gall and Victor Lempitsky. Class-specific hough forests for object detection. In *Decision forests for computer vision and medical image analysis*, pages 143–157. Springer, 2013. [2](#)
- [8] Juergen Gall, Angela Yao, Nima Razavi, Luc Van Gool, and Victor Lempitsky. Hough forests for object detection, tracking, and action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 33(11):2188–2202, 2011. [2](#)
- [9] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9224–9232, 2018. [2, 4](#)
- [10] Paul Guerrero, Yanir Kleiman, Maks Ovsjanikov, and Niloy J Mitra. Pcpnet learning local shape properties from raw point clouds. In *Computer Graphics Forum*, volume 37, pages 75–85. Wiley Online Library, 2018. [4](#)
- [11] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. *arXiv preprint arXiv:1703.06870*, 2017. [1, 6](#)
- [12] Ji Hou, Angela Dai, and Matthias Nießner. 3D-SIS: 3d semantic instance segmentation of rgb-d scans. *arXiv preprint arXiv:1812.07003*, 2018. [1, 2, 5, 6, 7, 13](#)
- [13] Paul VC Hough. Machine analysis of bubble chamber pictures. In *Conf. Proc.*, volume 590914, pages 554–558, 1959. [2](#)
- [14] Li Huan, Qin Yujian, and Wang Li. Vehicle logo retrieval based on hough transform and deep learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 967–973, 2017. [2](#)
- [15] Wadim Kehl, Fausto Milletari, Federico Tombari, Slobodan Ilic, and Nassir Navab. Deep learning of local rgb-d patches for 3d object detection and 6d pose estimation. In *European Conference on Computer Vision*, pages 205–220. Springer, 2016. [2](#)
- [16] Byung-soo Kim, Shili Xu, and Silvio Savarese. Accurate localization of 3d objects from rgbd data using segmentation hypotheses. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3182–3189, 2013. [2](#)
- [17] Roman Klokov and Victor Lempitsky. Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 863–872, 2017. [2](#)
- [18] Jan Knopp, Mukta Prasad, and Luc Van Gool. Orientation invariant 3d object classification using hough transform based methods. In *Proceedings of the ACM workshop on 3D object retrieval*, pages 15–20. ACM, 2010. [2](#)
- [19] Jan Knopp, Mukta Prasad, and Luc Van Gool. Scene cut: Class-specific object detection and segmentation in 3d scenes. In *2011 International Conference on 3D Imaging, Modeling, Processing, Visualization and Transmission*, pages 180–187. IEEE, 2011. [2](#)
- [20] Jean Lahoud and Bernard Ghanem. 2d-driven 3d object detection in rgbd images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4622–4630, 2017. [1, 2, 5, 6](#)
- [21] Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4558–4567, 2018. [4](#)
- [22] Truc Le and Ye Duan. Pointgrid: A deep network for 3d shape understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9204–9214, 2018. [2](#)
- [23] Bastian Leibe, Aleš Leonardis, and Bernt Schiele. Combined object categorization and segmentation with an implicit shape model. In *Workshop on statistical learning in computer vision, ECCV*, volume 2, page 7, 2004. [1](#)
- [24] Bastian Leibe, Aleš Leonardis, and Bernt Schiele. Robust object detection with interleaved categorization and segmentation. *International journal of computer vision*, 77(1–3):259–289, 2008. [2, 3](#)
- [25] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. In *Advances in Neural Information Processing Systems*, pages 828–838, 2018. [2, 4](#)
- [26] Yangyan Li, Angela Dai, Leonidas Guibas, and Matthias Nießner. Database-assisted object retrieval for real-time 3d reconstruction. In *Computer Graphics Forum*, volume 34. Wiley Online Library, 2015. [2](#)
- [27] Dahua Lin, Sanja Fidler, and Raquel Urtasun. Holistic scene understanding for 3d object detection with rgbd cameras. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1417–1424, 2013. [2](#)
- [28] Or Litany, Tal Remez, Daniel Freedman, Lior Shapira, Alex Bronstein, and Ran Gal. Asist: automatic semantically invariant scene transformation. *CVIU*, 157:284–299, 2017. [2](#)
- [29] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C

- Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. [2](#)
- [30] Subhransu Maji and Jitendra Malik. Object detection using a max-margin hough transform. 2009. [2](#)
- [31] Fausto Milletari, Seyed-Ahmad Ahmadi, Christine Kroll, Annika Plate, Verena Rozanski, Juliana Maiostre, Johannes Levin, Olaf Dietrich, Birgit Ertl-Wagner, Kai Bötzsel, et al. Hough-cnn: deep learning for segmentation of deep brain regions in mri and ultrasound. *Computer Vision and Image Understanding*, 164:92–102, 2017. [2](#)
- [32] Liangliang Nan, Ke Xie, and Andrei Sharf. A search-classify approach for cluttered indoor scene understanding. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia 2012)*, 31(6), 2012. [2](#)
- [33] David Novotny, Samuel Albanie, Diane Larlus, and Andrea Vedaldi. Semi-convolutional operators for instance segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 86–102, 2018. [2](#)
- [34] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgbd data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 918–927, 2018. [1, 2, 4, 5, 6, 7, 11](#)
- [35] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2017. [2, 4](#)
- [36] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413*, 2017. [1, 2, 3, 4, 5, 11](#)
- [37] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. [1, 2, 3, 5](#)
- [38] Zhile Ren and Erik B Sudderth. Three-dimensional object detection and layout prediction using clouds of oriented gradients. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1525–1533, 2016. [2, 6](#)
- [39] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointrcnn: 3d object proposal generation and detection from point cloud. *arXiv preprint arXiv:1812.04244*, 2018. [2](#)
- [40] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 567–576, 2015. [2, 5, 6](#)
- [41] Shuran Song and Jianxiong Xiao. Sliding shapes for 3d object detection in depth images. In *Computer Vision–ECCV 2014*, pages 634–651. Springer, 2014. [2](#)
- [42] Shuran Song and Jianxiong Xiao. Deep sliding shapes for amodal 3d object detection in rgb-d images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 808–816, 2016. [1, 2, 5, 6](#)
- [43] Hang Su, Varun Jampani, Deqing Sun, Subhransu Maji, Evangelos Kalogerakis, Ming-Hsuan Yang, and Jan Kautz. Splatnet: Sparse lattice networks for point cloud processing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2530–2539, 2018. [2, 4](#)
- [44] Min Sun, Gary Bradski, Bing-Xin Xu, and Silvio Savarese. Depth-encoded hough voting for joint object detection and shape recovery. In *European Conference on Computer Vision*, pages 658–671. Springer, 2010. [2](#)
- [45] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. *arXiv preprint arXiv:1703.09438*, 2017. [2](#)
- [46] Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-Yi Zhou. Tangent convolutions for dense prediction in 3d. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3887–3896, 2018. [2](#)
- [47] Alexander Velizhev, Roman Shapovalov, and Konrad Schindler. Implicit shape models for object detection in 3d point clouds. In *International Society of Photogrammetry and Remote Sensing Congress*, volume 2, 2012. [2](#)
- [48] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-cnn: Octree-based convolutional neural networks for 3d shape analysis. *ACM Transactions on Graphics (TOG)*, 36(4):72, 2017. [2](#)
- [49] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *arXiv preprint arXiv:1801.07829*, 2018. [2](#)
- [50] Oliver J Woodford, Minh-Tri Pham, Atsuto Maki, Frank Perbet, and Björn Stenger. Demisting the hough transform for 3d shape recognition and registration. *International Journal of Computer Vision*, 106(3):332–341, 2014. [2](#)
- [51] Saining Xie, Sainan Liu, Zeyu Chen, and Zhuowen Tu. Attentional shapecontextnet for point cloud recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4606–4615, 2018. [2](#)
- [52] Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. Spidercnn: Deep learning on point sets with parameterized convolutional filters. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 87–102, 2018. [2](#)
- [53] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 206–215, 2018. [2](#)
- [54] Li Yi, Wang Zhao, He Wang, Minhyuk Sung, and Leonidas Guibas. Gspn: Generative shape proposal network for 3d instance segmentation in point cloud. *arXiv preprint arXiv:1812.03320*, 2018. [2, 6](#)
- [55] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4490–4499, 2018. [1, 2](#)

A. Appendix

This appendix provides additional details on the network architectures and loss functions (Sec. A.1), more analysis experiment results (Sec. A.2), per-category results on ScanNet (Sec. A.3), and finally more visualizations (Sec. A.4).

A.1. Details on Architectures and Loss Functions

VoteNet architecture details. As mentioned in the main paper, the VoteNet architecture composes of a backbone point feature learning network, a voting module and a proposal module.

The backbone network, based on the PointNet++ architecture [36], has four set abstraction layers and two feature up-sampling layers. The detailed layer parameters are shown in Table 5. Each set abstraction (SA) layer has a receptive field specified by a ball-region radius r , a MLP network for point feature transform $MLP[c_1, \dots, c_k]$ where c_i is output channel number of the i -th layer in the MLP. The SA layer also subsamples the input point cloud with farthest point sampling to n points. Each SA layer is specified by $(n, r, [c_1, \dots, c_k])$ as shown in the Table 5. Compared to [36], we also normalize the XYZ scale of points in each local region by the region radius.

Each set feature propagation (FP) layer upsamples the point features by interpolating the features on input points to output points (each output point's feature is weighted average of its three nearest input points' features). It also combines the skip-linked features through a MLP (interpolated features and skip-linked features are concatenated before fed into the MLP). Each FP layer is specified by $[c_1, \dots, c_k]$ where c_i is the output of the i -th layer in the MLP.

The voting module as mentioned in the main paper is a MLP that transforms seeds' features to votes including a XYZ offset and a feature offset. The seed points are outputs of the fp2 layer. The voting module MLP has output sizes of 256, 256, 259 for its fully connected layers. The last fully connected layer does not have ReLU or BatchNorm.

The proposal module as mentioned in the main paper is a SA layer followed by another MLP after the max-pooling in each local region. We follow [34] on how to parameterize the oriented 3D bounding boxes. The layer's output has $5 + 2NH + 4NS + NC$ channels where NH is the number of heading bins (we predict a classification score for each heading bin and a regression offset for each bin – relative to the bin center and normalized by the bin size), NS is the number of size templates (we predict a classification score for each size template and 3 scale regression offsets for height, width and length) and NC is the number of semantic classes. In SUN RGB-D: $NH = 12, NS = NC = 10$, in ScanNet: $NH = 12, NS = NC = 18$. In the first 5 channels, the first two are for objectness classification and the rest three are for center regression (relative to the vote cluster center).

VoteNet loss function details. The network is trained end-to-end with a multi-task loss including a voting loss, an objectness loss, a 3D bounding box estimation loss and a semantic classification loss. We weight the losses such that they are in similar scales with $\lambda_1 = 0.5$, $\lambda_2 = 1$ and $\lambda_3 = 0.1$.

$$L_{\text{VoteNet}} = L_{\text{vote-reg}} + \lambda_1 L_{\text{obj-cls}} + \lambda_2 L_{\text{box}} + \lambda_3 L_{\text{sem-cls}} \quad (3)$$

Among the losses, the vote regression loss is as defined in the main paper (with L1 distance). For ScanNet we compute the ground truth votes as offset from the mesh vertices of an instances to the centers of the axis-aligned tight bounding boxes of the instances. Note that since the bounding box is not amodal, they can vary in sizes due to scan completeness (e.g. a chair may have a floating bounding box if its leg is not recovered from the reconstruction). For SUN RGB-D since the dataset does not provide instance segmentation annotations but only amodal bounding boxes, we cannot compute a ground truth vote directly (as we don't know which points are on objects). Instead, we consider any point inside an annotated bounding box as an object point (required to vote) and compute its offset to the box center as the ground truth. In cases that a point is in multiple ground truth boxes, we keep a set of up to three ground truth votes, and consider the minimum distance between the predicted vote and any ground truth vote in the set during vote regression on this point.

The objectness loss is just a cross-entropy loss for two classes and the semantic classification loss is also a cross-entropy loss of NC classes.

The box loss follows [34] (but without the corner loss regularization for simplicity) and is composed of center regression, heading estimation and size estimation sub-losses. In all regression in the box loss we use the robust L_1 -smooth loss. Both the box and semantic losses are only computed on positive vote clusters and normalized by the number of positive clusters. We refer readers to [34] for more details.

$$\begin{aligned} L_{\text{box}} = & L_{\text{center-reg}} + 0.1 L_{\text{angle-cls}} + L_{\text{angle-reg}} \\ & + 0.1 L_{\text{size-cls}} + L_{\text{size-reg}} \end{aligned} \quad (4)$$

One difference though is that, instead of a naive regression loss, we use a *Chamfer loss* [6] for $L_{\text{center-reg}}$ (between regressed centers and ground truth box centers). It requires that each positive proposal is close to a ground truth object and each ground truth object center has a proposal near it. The latter part also influences the voting in the sense that it encourages non-object seed points near the object to also vote for the center of the object, which helps further increase contexts in detection.

layer name	input layer	type	output size	layer params
sa1	raw point cloud	SA	(2048,3+128)	(2048,0.2,[64,64,128])
sa2	sa1	SA	(1024,3+256)	(1024,0.4,[128,128,256])
sa3	sa2	SA	(512,3+256)	(512,0.8,[128,128,256])
sa4	sa3	SA	(256,3+256)	(256,1.2,[128,128,256])
fp1	sa3, sa4	FP	(512,3+256)	[256,256]
fp2	sa2, sa3	FP	(1024,3+256)	[256,256]

Table 5. Backbone network architecture: layer specifications.

BoxNet architecture details. Our baseline network without voting, BoxNet, shares most parts with the VoteNet. They share the same backbone architecture. But instead of voting from seeds, the BoxNet directly proposes bounding boxes and classifies object classes from seed points’ features. To make the BoxNet and VoteNet have similar capacity we also include a SA layer for the proposal in BoxNet. However this SA layer takes *seed clusters* instead of *vote clusters* i.e. it samples seed points and then combines neighboring seeds with MLP and max-pooling. This SA layer has exactly the same layer parameters with that in the VoteNet, followed by the same MLP_2 .

BoxNet loss function details. BoxNet has the same loss function as VoteNet, except it is not supervised by vote regression. There is also a slight difference in how objectness labels (used to supervise objectness classification) are computed. As seed points (on object surfaces) are often far from object centroids, it no longer works well to use the distances between seed points and object centroids to compute the objectness labels. In BoxNet, we assign positive objectness labels to seed points that are on objects (those belonging to the semantic categories we consider) and negative labels to all the other seed points on clutter (e.g. walls, floors).

$$L_{\text{BoxNet}} = \lambda_1 L_{\text{obj-cls}} + \lambda_2 L_{\text{box}} + \lambda_3 L_{\text{sem-cls}} \quad (5)$$

A.2. More Analysis Experiments

Average precision and recall plots Fig. 8 shows how average precision (AP) and average recall (AR) change as we increase the number of proposals. The AP and AR are both averaged across 10 categories on SUN RGB-D. We report two ways of using the proposals: joint and per-class. For the joint proposal we propose K objects’ bounding boxes for all the 10 categories, where we consider each proposal as the semantic class it has the largest confidence in, and use their objectness scores to rank them. For the per-class proposal we duplicate the K proposal 10 times thus have K proposals per class where we use the multiplication of semantic probability for that class and the objectness probability to rank them. The latter way of using proposals gives us a slight improvement on AP and a big boost on AR.

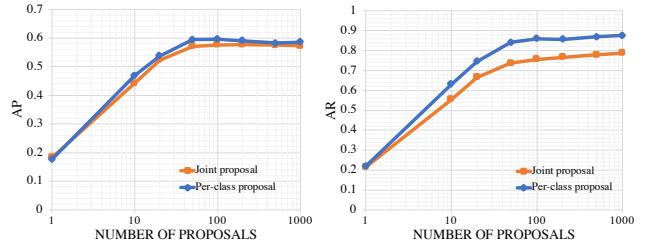


Figure 8. Number of proposals per scene v.s. Average Precision (AP) and Average Recall (AR) on SUN RGB-D. The AP and AR are averaged across the 10 classes. The recall is maximum recall given a fixed number of detection per scene. The “joint proposal” means that we assign each proposal to a single class (the class with the highest classification score); The “per-class proposal” means that we assign each proposal to all the 10 classes (the objectness score is multiplied by the semantic classification probability).

We see that with as few as 10 proposals our VoteNet can achieve a decent AP of around 45% while having 100 proposals already pushes the AP to above 57%. With a thousand proposals, our network can achieve around 78.7% recall with joint proposal and around 87.7% recall with per-class proposal.

Context of voting One difference of a deep Hough voting scheme with the traditional Hough voting is that we can take advantage of deep features, which can provide more context knowledge for voting. In Table 8 we show how features from different levels of the PointNet++ affects detection performance (from SA2 to FP3, the network has increasing contexts for voting). FP3 layer is extended from the FP2 with a MLP of output sizes 256 and 256 with 2048 output points (the same set of XYZ as that output by SA1).

It is surprising to find that voting from even SA2 can achieve reasonable detection results (mAP 51.2%) while voting from FP2 achieves the best performance. Having larger context (e.g. FP3) than FP2 does not show further improvements on the performance.

Multiple votes per seed In default we just generate one vote per seed since we find that with large enough context there is little need to generate more than one vote to resolve ambiguous cases. However, it is still possible to generate

	cab	bed	chair	sofa	tabl	door	wind	bkshf	pic	cntr	desk	curt	fridg	showr	toil	sink	bath	ofurn	mAP
3DSIS 5views [12]	19.76	69.71	66.15	71.81	36.06	30.64	10.88	27.34	0.00	10.00	46.93	14.06	53.76	35.96	87.60	42.98	84.30	16.20	40.23
3DSIS Geo [12]	12.75	63.14	65.98	46.33	26.91	7.95	2.79	2.30	0.00	6.92	33.34	2.47	10.42	12.17	74.51	22.87	58.66	7.05	25.36
VoteNet ours	36.27	87.92	88.71	89.62	58.77	47.32	38.10	44.62	7.83	56.13	71.69	47.23	45.37	57.13	94.94	54.70	92.11	37.20	58.65

Table 6. 3D object detection scores per category on the ScanNetV2 dataset, evaluated with mAP@0.25 IoU.

	cab	bed	chair	sofa	tabl	door	wind	bkshf	pic	cntr	desk	curt	fridg	showr	toil	sink	bath	ofurn	mAP
3DSIS 5views [12]	5.73	50.28	52.59	55.43	21.96	10.88	0.00	13.18	0.00	0.00	23.62	2.61	24.54	0.82	71.79	8.94	56.40	6.87	22.53
3DSIS Geo [12]	5.06	42.19	50.11	31.75	15.12	1.38	0.00	1.44	0.00	0.00	13.66	0.00	2.63	3.00	56.75	8.68	28.52	2.55	14.60
VoteNet (ours)	8.07	76.06	67.23	68.82	42.36	15.34	6.43	28.00	1.25	9.52	37.52	11.55	27.80	9.96	86.53	16.76	78.87	11.69	33.54

Table 7. 3D object detection scores per category on the ScanNetV2 dataset, evaluated with mAP@0.5 IoU.

Seed layer	SA2	SA3	SA4	FP1	FP2	FP3
mAP	51.2	56.3	55.1	56.6	57.7	57.1

Table 8. Effects of seed context for 3D detection. Evaluation metric is mAP@0.25 on SUN RGB-D.

Vote factor	1	3	3
Random number	N	N	Y
mAP	57.7	55.8	55.8

Table 9. Effects of number of votes per seed. Evaluation metric is mAP@0.25 on SUN RGB-D. If random number is on, we concatenate a random number to the seed feature before voting, which helps break symmetry in the case of multiple votes per seed.

more than one vote with our network architecture. Yet to break the symmetry in multiple vote generation, one has to introduce some bias to different votes to prevent them from pointing to the same place.

In experiments, we find that one vote per seed achieves the best results, as shown in Table 9. We ablate by using a vote factor of 3, where the voting module generates 3 votes per seed with a MLP layer spec: [256, 256, 259 * 3]). In computing the vote regression loss on a seed point, we consider the minimum distance between any predicted votes to the ground truth vote (in case of SUN RGB-D where we may have a set of ground truth votes for a seed, we compute the minimum distance among any pair of predicted vote and ground truth vote).

To break symmetry, we generate 3 random numbers and inject them to the second last features from the MLP layer. We show results both with and without this procedure which shows no observable difference.

On proposal sampling In the proposal step, to generate K proposals from the votes, we need to select K vote clusters. How to select those clusters is a design choice we study here (each cluster is simply a group of votes near a center

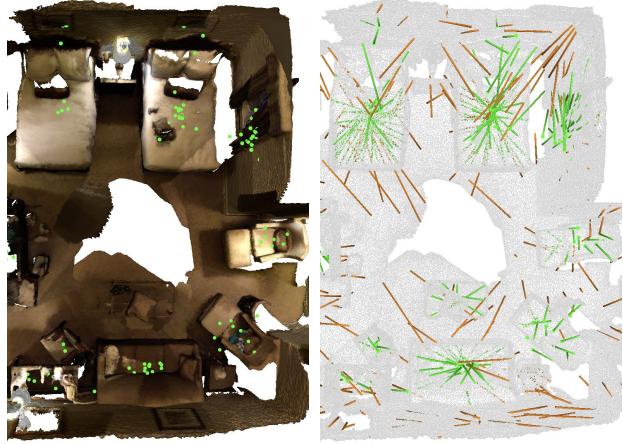


Figure 9. Vote meeting point. *Left:* ScanNet scene with votes coming from object points. *Right:* vote offsets from source seed-points to target-votes. Object votes are colored green, and non-object ones are colored red. See how object points from all-parts of the object vote to form a cluster near the center. Non-object points, however, either vote “nowhere” and therefore lack structure, or are near object and have gathered enough context to also vote properly.

vote). In Table 10, we report mAP results on SUN RGB-D with 256 proposals (joint proposal) using cluster sampling strategies of vote FPS, seed FPS and random sampling, where FPS means farthest point sampling. From 1024 vote clusters, vote FPS samples K clusters based on votes’ XYZ. Seed FPS firstly samples on seed XYZ and then finds the votes corresponding to the sampled seeds – it enables a direct comparison with BoxNet as it uses the same sampling scheme, making the two techniques similar up to the space in which the points are grouped: VoteNet groups votes according to vote XYZ, while BoxNet groups seeds according to seed XYZ. Random sampling simply selects a random set of K votes and take their neighborhoods for proposal generation. Note that the results from Table 10 are from the same model trained with vote FPS to select proposals.

Proposal sampling	mAP
Random sampling	57.5
Farthest point sampling on votes	57.2
Farthest point sampling on seeds	57.7

Table 10. **Effects of proposal sampling.** Evaluation metric is mAP@0.25 on SUN RGB-D. 256 proposals are used for all evaluations. Our method is not sensitive to how we choose centers for vote groups/clusters.

Dataset	with height	without height
SUN RGB-D	57.7	57.0
ScanNet	58.6	58.1

Table 11. **Effects of the height feature.** Evaluation metric is mAP@0.25 on both datasets.

We can see that while seed FPS gets the best number in mAP, the difference caused by different sampling strategies is small, showing the robustness of our method.

Effects of the height feature In point clouds from indoor scans, point height is a useful feature in recognition. As mentioned in the main paper, we can use 1% of the Z values (Z -axis is up-right) of all points from a scan as an approximate as the floor height z_{floor} , and then compute the a point (x, y, z) 's height as $z - z_{\text{floor}}$. In Table 11 we show how this extra height feature affect detection performance. We see that adding the height feature consistently improves performance in both SUN RGB-D and ScanNet.

A.3. ScanNet Per-class Evaluation

Table 6 and Table 7 report per-class average precision on 18 classes of ScanNetV2 with 0.25 and 0.5 box IoU thresholds respectively. Relying on purely geometric data, our method excels (esp. with mAP@0.25) in detecting objects like bed, chair, table, desk etc. where geometry is a strong cue for recognition; and struggles with objects best recognized by texture and color like pictures.

A.4. Visualization of Votes

Fig. 9 shows (a subset of) votes predicted from our VoteNet in a typical ScanNet scene. We clearly see that seed points on objects (bed, sofa etc.) vote to object centers while clutter points vote either to object center as well (if the clutter point is close to the object) or to nowhere due to lack of structure in the clutter area (e.g. a wall).