

S-VoteNet: Deep Hough Voting with Spherical Proposal for 3D Object Detection

Yanxian Chen¹, Huimin Ma^{1*}, Xi Li², Xiong Luo¹

¹University of Science and Technology Beijing ²Tsinghua University

Motivation

Indoor Scene Object Detection:





(x, y, z, l, w, h, orientation) + class

2D Object Detection (x, y, w, h) + class

Challenge of Indoor 3D Object Detection:

- Indoor scenes usually contain many objects, and the categories of objects are variety.
- There are large size differences between different categories of objects, which has a great impact on object localization.
- There are complex spatial relationships between objects, some clustered in groups, and some stacked on other objects.
- These factors make it difficult for 3D detectors to accurately predict the position and size of objects in a regression head.

Contribution:

- We show that with a proper structure, the decoupling of 3D bounding box regression can effectively improve the performance of the 3D detector.
- The proposed spherical center loss further considers the geometric distance between proposal and ground truth, which achieves higher 3D localization accuracy.
- Our S-VoteNet achieves state-of-the-art 3D object detection performance on SUN RGB-D dataset by only using point cloud as input.

Related Work

3D Box Encoding:

- Axis Aligned 3D box: no orientation, 6-dim vector, (x, y, z, dx, dy, dz)
- Oriented 3D box: original, 7-dim vector, (*x*, *y*, *z*, *l*, *w*, *h*, *orientation*)

8-corners, 24-dim vector, $(x_i, y_i, z_i), i \in [1, 8]$

4-corners, 10-dim vector, $h_{top}, h_{bottom}, (x_i, y_i), i \in [1, 4]$

Object Location Loss:

- *l*2 center loss: The Euclidean distance between proposal and ground truth is used as supervision.
- IoU loss: The intersection over union between proposal and ground truth is used as supervision.

Methodology



Spherical encoding: 4-dim vector, (x, y, z, r)

Based on spherical encoding, 3D object detection task can be decoupled into object location task, size and orientation prediction task.

For object location, we use spherical center loss to constrain the prediction result. For size and orientation, we adopt the method of F-PointNet.

Spherical Center Loss:

 $L_{center} = d$ $L_{spherical-center} = \frac{d}{d + r_{pro} + r_{gt}}$

- The distance between ground truth and proposal is the same for objects of different sizes, but the IoU is different.
- Compared to l2 center loss, spherical center loss introduces object size into object center prediction, achieving higher 3D localization accuracy.

Geometric Information of Point Cloud:



From the geometric information of point cloud on different stage, seeds are suitable for object size and orientation prediction, while votes are fit for object location prediction.

Overall Structure of S-VoteNet:



- S-VoteNet is built on the basis of VoteNet, which introduces spherical proposal to decouple the 3D object detection task.
- To align the object center predictions with the size and orientation predictions, we use votes indices to find the corresponding seeds.

Result

Performance on SUN RGB-D Val Set:

methods	input	bathtub	bed	bookshelf	chair	desk	dresser	nightstand	sofa	table	toilet mAP
DSS COG 2D-driven F-PointNet VoteNet + region feature ImVoteNet	Geo + RGB Geo + RGB Geo + RGB Geo + RGB Geo + RGB Geo + RGB	44.2 58.3 43.5 44.3 71.7 75.9	78.8 63.7 64.5 81.1 86.1 87.6	11.9 31.8 31.4 33.3 34.0 41.3	61.2 62.2 48.3 64.2 74.7 76.7	20.5 45.2 27.9 24.7 26.0 28.7	6.4 15.5 25.9 32.0 34.2 41.4	15.4 27.4 41.9 58.1 64.3 69.9	53.5 51.0 50.4 61.1 66.5 70.7	50.3 51.3 37.0 51.1 49.7 51.1	78.9 42.1 70.1 47.6 80.4 45.1 90.9 54.0 88.4 59.6 90.5 63.4
VoteNet S-VoteNet (ours)	Geo only	74.4	83.0 85.6	28.8 35.9	75.3	22.0 26.9	29.8 31.7	62.2 65.5	64.0 67.6	47.3 47.5	90.1 57.7

 S-VoteNet advances the baseline by 2.6% mAP, which achieves performance second only to ImVoteNet without the use of RGB information.

Ablation Study:

methods	use spherical center loss		use seed		mAP	•
BoxNet	×	1	\checkmark	1	53.0	-
VoteNet	×		×		57.7	
VoteNet*	×		×		58.0	•
VoteNet**	\checkmark		×		59.5	
S-VoteNet	, √		\checkmark		60.3	

BoxNet is the baseline of VoteNet, which generates proposals without the voting module.

VoteNet* is a variant of VoteNet, which decouples 3D object detection task without spherical encoding.

Qualitative Results:



📕 bed 🔳 table 📕 chair 📕 desk 🔲 dresser 📕 nightstand 📕 bookshelf

The distance between ground truth