

ImVoteNet: Boosting 3D Object Detection in Point Clouds with Image Votes

Charles R. Qi^{*†}Xinlei Chen^{*1}
¹Facebook AIOr Litany^{1,2}
²Stanford UniversityLeonidas J. Guibas^{1,2}

Abstract

3D object detection has seen quick progress thanks to advances in deep learning on point clouds. A few recent works have even shown state-of-the-art performance with just point clouds input (e.g. VOTENET). However, point cloud data have inherent limitations. They are sparse, lack color information and often suffer from sensor noise. Images, on the other hand, have high resolution and rich texture. Thus they can complement the 3D geometry provided by point clouds. Yet how to effectively use image information to assist point cloud based detection is still an open question. In this work, we build on top of VOTENET and propose a 3D detection architecture called IMVOTENET specialized for RGB-D scenes. IMVOTENET is based on fusing 2D votes in images and 3D votes in point clouds. Compared to prior work on multi-modal detection, we explicitly extract both geometric and semantic features from the 2D images. We leverage camera parameters to lift these features to 3D. To improve the synergy of 2D-3D feature fusion, we also propose a multi-tower training scheme. We validate our model on the challenging SUN RGB-D dataset, advancing state-of-the-art results by **5.7 mAP**. We also provide rich ablation studies to analyze the contribution of each design choice.

1. Introduction

Recognition and localization of objects in a 3D environment is an important first step towards full scene understanding. Even such low dimensional scene representation can serve applications like autonomous navigation and augmented reality. Recently, with advances in deep networks for point cloud data, several works [33, 56, 41] have shown state-of-the-art 3D detection results with point cloud as the *only* input. Among them, the recently proposed VOTENET [33] work by Qi *et al.*, taking 3D geometry input only, showed remarkable improvement for indoor object recognition compared with previous works that exploit

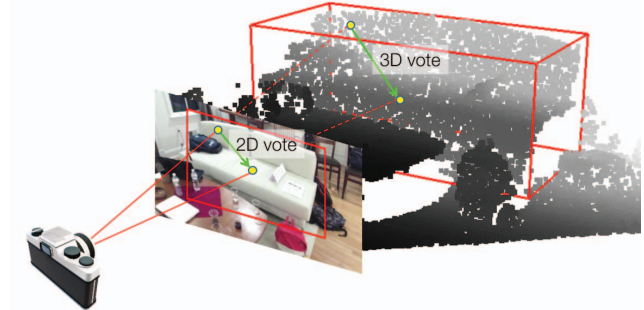


Figure 1. Voting using both an image and a point cloud from an indoor scene. The 2D vote reduces the search space of the 3D object center to a ray while the color texture in image provides a strong semantic prior. Motivated by the observation, our model lifts the 2D vote to 3D to boost 3D detection performance.

all RGB-D channels. This leads to an interesting research question: Is 3D geometry data (point clouds) sufficient for 3D detection, or is there any way RGB images can further boost current detectors?

By examining the properties of point cloud data and RGB image data (see for example Fig. 1), we believe the answer is clear: RGB images have value in 3D object detection. In fact, images and point clouds provide *complementary* information. RGB images have higher resolution than depth images or LiDAR point clouds and contain rich textures that are not available in the point domain. Additionally, images can cover “blind regions” of active depth sensors which often occur due to reflective surfaces. On the other hand, images are limited in the 3D detection task as they lack absolute measures of object depth and scale, which are exactly what 3D point clouds can provide. These observations, strengthen our intuition that images can help point cloud-based 3D detection.

However, how to make effective use of 2D images in a 3D detection pipeline is still an open problem. A naïve way is to directly append raw RGB values to the point clouds – since the point-pixel correspondence can be established through projection. But since 3D points are much sparser, in doing so we will lose the dense patterns from the image domain. In light of this, more advanced ways to fuse 2D and 3D data have been proposed recently. One line of

^{*}: equal contributions.

[†]: work done while at Facebook.

work [34, 52, 19] uses mature 2D detectors to provide initial proposals in the form of frustums. This limits the 3D search space for estimating 3D bounding boxes. However, due to its *cascaded* design, it does not leverage 3D point clouds in the initial detection. In particular, if an object is missed in 2D, it will be missed in 3D as well. Another line of work [45, 18, 48, 11] takes a more 3D-focused way to concatenate intermediate ConvNet features from 2D images to 3D voxels or points to enrich 3D features, before they are used for object proposal and box regression. The downside of such systems is that they do not use 2D images directly for localization, which can provide helpful guidance for detection objects in 3D.

In our work, we build upon the successful VOTENET architecture [33] and design a *joint* 2D-3D voting scheme for 3D object detection named IMVOTENET. It takes advantage of the more mature 2D detectors [38] but at the same time still reserves the ability to propose objects from the full point cloud itself – combining the best of both lines of work while avoiding the drawbacks of each. One motivation for our design is to leverage both geometric and semantic/textural cues in 2D images (Fig. 1). The geometric cues come from accurate 2D bounding boxes in images, such as the output by a 2D detector. Instead of solely relying on the 2D detection for object proposal [34], we defer the proposal process to 3D. Given a 2D box, we generate 2D votes on the image space, where each vote connects from the object pixel to the 2D amodal box center. To pass the 2D votes to 3D, we *lift* them by applying geometric transformations based on the camera intrinsic and pixel depth, so as to generate “pseudo” 3D votes. These pseudo 3D votes become extra features appended to seed points in 3D for object proposals. Besides geometric cues from the 2D votes, each pixel also passes semantic and texture cues to the 3D points, as either features extracted per-region, or ones extracted per-pixel.

After lifting and passing all the features from the images to 3D, we concatenate them with the 3D point features from a point cloud backbone network [35, 36]. Next, following the VOTENET pipeline, those points with the fused 2D and 3D features generate 3D Hough votes [12] – not limited by 2D boxes – toward object centers and aggregate the votes to produce the final object detections in 3D. As the seed features have both 2D and 3D information, they are intuitively more informative for recovering heavily truncated objects or objects with few points, as well as more confident in distinguishing geometrically similar objects.

In addition, we recognize that when fusing 2D and 3D sources, one has to carefully balance the information from two modalities to avoid one being dominated by the other. To this end, we further introduce a multi-towered network structure with gradient blending [49] to ensure our network makes the best use of both the 2D and 3D features. During testing, only the main tower that operates on the joint 2D-

3D features are used, minimizing the sacrifice on efficiency.

We evaluate IMVOTENET on the challenging SUN RGB-D dataset [43]. Our model achieves the state-of-the-art results while showing a significant improvement (+5.7 mAP) over the 3D geometry only VOTENET, validating the usefulness of image votes and 2D features. We also provide extensive ablation studies to demonstrate the importance of each individual component. Finally, we also explore the potential of using color to compensate for *sparsity* in depth points, especially for the case of lower quality depth sensors or for cases where depth is estimated from a moving monocular camera (SLAM), showing potential of our method to more broader use cases.

To summarize, the contributions of our work are:

1. A geometrically principled way to fuse 2D object detection cues into a point cloud based 3D detection pipeline.
2. The designed deep network IMVOTENET achieves state-of-the-art 3D object detection performance on SUN RGB-D.
3. Extensive analysis and visualization to understand various design choices of the system.

2. Related Work

Advances in 3D sensing devices have led to a surge of methods designed to identify and localize objects in a 3D scene. The most relevant lines of work are detection with point clouds and detection with full RGB-D data. We also briefly discuss a few additional relevant works in the area of multi-modal data fusion.

3D object detection with point clouds. To locate objects using purely geometric information, one popular line of methods is based on template matching using a collection of clean CAD models either directly [21, 28, 25], or through extracted features [44, 2]. More recent methods are based on point cloud deep nets [35, 56, 20, 41, 33]. In the context of 3D scene understanding, there have also been promising results on semantic and instance segmentation [53, 4, 9]. Most relevant to our work are PointRCNN [41] and Deep Hough Voting (VOTENET) [33] (related to earlier Hough voting ideas [47, 50, 16, 15]) which demonstrated state-of-the-art 3D object detection in outdoor and indoor scenes, respectively. Notably, these results are achieved *without* using the RGB input. To leverage this additional information, we propose a way to further boost detection performance in this work.

3D object detection with RGB-D data. Depth and color channels both contain useful information that can be useful for 3D object detection. Prior methods for fusing those two modalities broadly fall into three categories: 2D-driven,

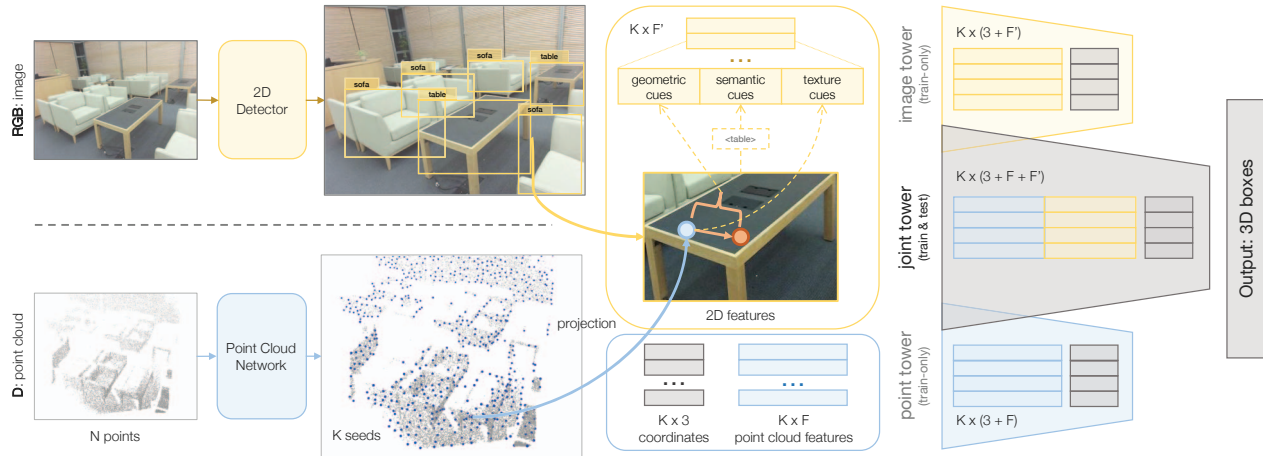


Figure 2. **3D object detection pipeline for IMVOTENET.** Given RGB-D input (with the depth image converted to a point cloud), the model initially have two separate branches: one for 2D object detection on the image and the other for point cloud feature extraction (with a PointNet++ [36] backbone) on the point clouds. Then we lift 2D image votes as well as semantic and texture cues to the 3D seed points (the fusion part). The seed points with concatenated image and point cloud features then generate votes towards 3D object centers and also propose 3D bounding boxes with its features (the joint tower). To push for more effective multi-modality fusion, we have two other towers that take image features only (image tower) and point cloud features only (point tower) for voting and box proposals.

3D-driven, and feature concatenation. The first type of method [19, 34, 6, 52] starts with object detections in the 2D image, which are then used to guide the search space in 3D. By 3D-driven, we refer to methods that first generate region proposals in 3D and then utilize 2D features to make a prediction, such as the Deep Sliding Shapes [45]. Recently more works focus on fusing 2D and 3D features earlier in the process such as Multi-modal Voxernet [48], AVOD [18], multi-sensor [22] and 3D-SIS [11]. However, all these mostly perform fusion through concatenation of 2D features to 3D features. Our proposed method is more closely related to the third type, but differs from it in two important aspects. First, we propose to make explicit use of geometric cues from the 2D detector and lift them to 3D in the form of pseudo 3D votes. Second, we use a multi-tower architecture [49] to balance features from both modalities, instead of simply training on the concatenated features.

Multi-modal fusion in learning. How to fuse signals from multiple modalities is an open research problem in other areas than 3D object detection. Semantic segmentation of 3D scenes often uses both RGB and depth data [10, 27]. Vision and language research develops ways to jointly reason over visual data and texts [7, 32, 54] for tasks like visual question answering [1, 14]. Another active area of research is video+sound [30, 8], where the additional sound track can either provide supervision signal [31], or propose interesting tasks to test joint understanding of both streams [55]. Targeting at all such tasks, a recent gradient blending approach [49] is proposed to make the multi-modal network more robust (to over-fitting and different convergence rates), which is adopted in our approach too.

3. ImVoteNet Architecture

We design a 3D object detection solution suited for RGB-D scenes, based on the recently proposed deep Hough voting framework (VOTENET [33]) by passing *geometric* and *semantic/texture* cues from 2D images to the voting process (as illustrated in Fig. 2). In this section, after a short summary of the original VOTENET pipeline, we describe how to build ‘2D votes’ with the assistance of 2D detectors on RGB and explain how the 2D information is lifted to 3D and passed to the point cloud to improve the 3D voting and proposal. Finally, we describe our multi-tower architecture for fusing 2D and 3D detection with gradient blending [49]. More implementation details are provided in supplement.

3.1. Deep Hough Voting

VOTENET [33] is a feed-forward network that consumes a 3D point cloud and outputs object proposals for 3D object detection. Inspired by the seminal work on the generalized Hough transform [3], VOTENET proposes an adaptation of the voting mechanism for object detection to a deep learning framework that is fully differentiable.

Specifically, it is comprised of a point cloud feature extraction module that enriches a subsampled set of scene points (called *seeds*) with high-dimensional features (bottom of Fig. 2 from $N \times 3$ input points to $K \times (3+F)$ seeds). These features are then pushed through a Multi-Layer-Perceptron (MLP) to generate *votes*. Every vote is both a point in the 3D space with its Euclidean coordinates (3-dim) supervised to be close to the object center, and a feature vector learned for the final detection task (F-dim). The

votes form a clustered point cloud near object centers and are then processed by another point cloud network to generate object proposals and classification scores. This process is equivalent to the pipeline in Fig. 2 with just the point tower and without the image detection and fusion.

VOTENET recently achieved state-of-the-art results on indoor 3D object detection in RGB-D [33]. Yet, it is solely based on point cloud inputs and neglects the image channels which, as we show in this work, are a very useful source of information. In IMVOTENET, we leverage the additional image information and propose a lifting module from 2D votes to 3D that improves detection performance. Next, we explain how to get 2D votes in images and how we lift its geometric cues to 3D together with semantic/texture cues.

3.2. Image Votes from 2D Detection

We generate image votes based on a set of candidate boxes from 2D detectors. An *image vote*, in its geometric part, is simply a vector connecting an image pixel and the center of the 2D object bounding box that pixel belongs to (see Fig. 1). Each image vote is also augmented with its semantic and texture cues from the features of its source pixel, such that each image vote has F' dimension in total as in the fusion block in Fig. 2.

To form the set of boxes given an RGB image, we apply an off-the-shelf 2D detector (e.g. Faster R-CNN [38]) pre-trained on color channels of the RGB-D dataset. The detector outputs the M most confident bounding boxes and their corresponding classes. We assign each pixel within a detected box a vote to the box center. Pixels inside multiple boxes are given multiple votes (corresponding 3D seed points are duplicated for each of them), and those outside of any box are padded with zeros. Next we go to details on how we derive geometric, semantic and texture cues.

Geometric cues: lifting image votes to 3D The translational 2D votes provide useful geometric cues for 3D object localization. Given the camera matrix, the 2D object center in the image plane becomes a ray in 3D space connecting the 3D object center and the camera optical center (Fig. 1). Adding this information to a seed point can effectively narrow down the 3D search space of the object center to 1D.

In details, as shown in Fig. 3, given an object in 3D with its detected 2D bounding box in the image plane, we denote the 3D object center as C and its projection onto the image as c . A point P on the object surface is associated with its projected point p in the image plane, hence knowing the 2D vote to the 2D object center c , we can reduce the search space for the 3D center to a 1D position on the ray OC . We now derive the computation we follow to pass the ray information to the 3D seed point. Defining $P=(x_1, y_1, z_1)$ in the camera coordinate, and $p=(u_1, v_1)$, $c=(u_2, v_2)$ in the image plane coordinate, we seek to recover the 3D object

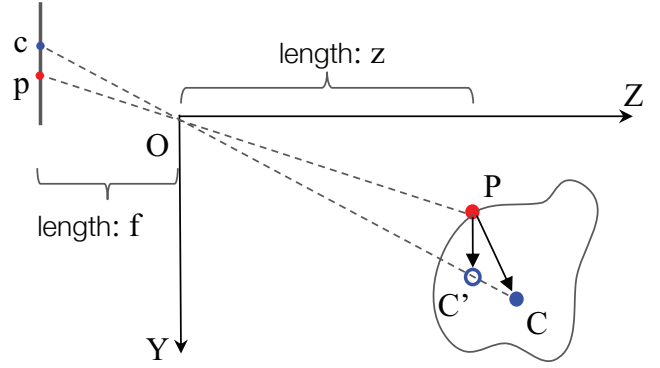


Figure 3. **Illustration of the pseudo 3D vote.** In the figure, P is a surface point in 3D, C is the unknown object center while p and c are their projections on the image plane respectively. C' is the pseudo 3D center and the vector PC' is the pseudo 3D vote.

center $C=(x_2, y_2, z_2)$ (voting target for the 3D point P). The true 3D vote from P to C is:

$$\overrightarrow{PC} = (x_2 - x_1, y_2 - y_1, z_2 - z_1). \quad (1)$$

The 2D vote, assuming a simple pin-hole camera¹ with focal length f , can be written as:

$$\begin{aligned} \overrightarrow{pc} &= (u_2 - u_1, v_2 - v_1) = (\Delta u, \Delta v) \\ &= (f(\frac{x_2}{z_2} - \frac{x_1}{z_1}), f(\frac{y_2}{z_2} - \frac{y_1}{z_1})). \end{aligned} \quad (2)$$

We further assume the depth of the surface point P is similar to the center point C . This is a reasonable assumption for most objects when they are not too close to the camera. Then, given $z_1 \approx z_2$, we compute $\overrightarrow{PC'}$,

$$\overrightarrow{PC'} = (\frac{\Delta u}{f} z_1, \frac{\Delta v}{f} z_1, 0), \quad (3)$$

which we refer to as a *pseudo 3D vote*, as C' lies on the ray OC and is in the proximity of C . This pseudo 3D vote provides information about where the 3D center is relative to the point surface point P .

To compensate for the error caused by the depth approximation ($z_1 \approx z_2$), we pass the ray direction as extra information to the 3D surface point. The error (along the X -axis) caused by the approximated depth, after some derivation, can be expressed by

$$\text{err}_x = \Delta x - \Delta x' = \frac{x_2}{z_2} (z_2 - z_1). \quad (4)$$

Hence, if we input the direction of the ray \overrightarrow{OC} : $(x_2/z_2, y_2/z_2)$, the network should have more information

¹See supplementary for more details on how to deal with a general camera model and camera-to-world transformations.

to estimate the true 3D vote by estimate the depth different $\Delta z = z_2 - z_1$. As we do not know the true 3D object center C , we can use the ray direction of $\overrightarrow{OC'}$ which aligns with \overrightarrow{OC} after all, where

$$\begin{aligned}\overrightarrow{OC'} &= \overrightarrow{OP} + \overrightarrow{PC'} \\ &= (x_1 + \frac{\Delta u}{f} z_1, y_1 + \frac{\Delta v}{f} z_1, z_1).\end{aligned}\quad (5)$$

Normalizing and concatenating with the pseudo vote, the image geometric features we pass to the seed point P are:

$$(\frac{\Delta u}{f} z_1, \frac{\Delta v}{f} z_1, \frac{\overrightarrow{OC'}}{\|\overrightarrow{OC'}\|}). \quad (6)$$

Semantic cues On top of the geometric features just discussed that just use the spatial coordinates of the bounding boxes, an important type of information RGB can provide is features that convey a semantic understanding of what's inside the box. This information often complements what can be learned from 3D point clouds and can help to distinguish between classes that are geometrically very similar (such as table *vs.* desk or nightstand *vs.* dresser).

In light of this, we provide additional *region-level* features extracted per bounding box as semantic cues for 3D points. For all the 3D seed points that are projected within the 2D box, we pass a vector representing that box to the point. If a 3D seed point falls into more than one 2D boxes (*i.e.*, when they overlap), we duplicate the seed point for each of the overlapping 2D regions (up to a maximum number of K). If a seed point is not projected to any 2D box, we simply pass an all-zero feature vector for padding.

It is important to note that the 'region features' here include but are *not* limited to features extracted from RoI pooling operations [38]. In fact, we find representing each box with a simple one-hot class vector (with a confidence score for that class) is already sufficient to cover the semantic information needed for disambiguation in 3D. It not only gives a light-weight input (*e.g.* 10-dim [46] *vs.* 1024-dim [23]) that performs well, but also generalizes to *all* other competitive (*e.g.* faster) 2D detectors [37, 26, 24] that do not explicitly use RoI but directly outputs classification scores. Therefore, we use this semantic cue by default.

Texture cues Different from the depth information that spreads sparsely in the 3D space, RGB images can capture high-resolution signals at a dense, *per-pixel* level in 2D. While region features can offer a high-level, semantic-rich representation per bounding box, it is complementary and equally important to use the low-level, texture-rich representations as another type of cues. Such cues can be passed to the 3D seed points via a simple mapping: a seed point gets pixel features from the corresponding pixel of its 2D

projection².

Although any learned, convolutional feature maps with spatial dimensions (height and width) can serve our purpose, by default we still use the simplest texture feature by feeding in the raw RGB pixel-values directly. Again, this choice is not only light-weight, but also makes our pipeline less sensitive to the image neural nets that may incur biases.

Experimentally, we show that even with such minimalist choice of both our semantic and texture cues, significant performance boost over geometric-only VOTENET can be achieved with our multi-tower training paradigm, which we discuss next.

3.3. Feature Fusion and Multi-tower Training

With lifted image votes and its corresponding semantic and texture cues ($K \times F'$ in the fusion block in Fig. 2) as well as the point cloud features with the seed points $K \times F$, each seed point can generate 3D votes and aggregate them to propose 3D bounding boxes (through a voting and proposal module similar to that in [33]). Yet it takes extra care to optimize the deep network to fully utilize cues from all modalities. As a recent paper [49] mentions, without a careful strategy, multi-modal training can actually result in degraded performance as compared to a single modality training. The reason is that different modalities may learn to solve the task at different rates so, without attention, certain features may dominate the learning and result in overfitting. In this work, we follow the gradient blending strategy introduced in [49] to weight the gradient for different modality towers (by weighting the loss functions).

In our multi-tower formulation, as shown in Fig. 2, we have three towers taking seed points with three sets of features: point cloud features only, image features only and joint features. Each tower has the same target task of detecting 3D objects – but they each have their separate 3D voting and box proposal network parameters as well as their separate losses. The final training loss is the weighted sum of three detection losses:

$$L = w_{\text{img}} L_{\text{img}} + w_{\text{point}} L_{\text{point}} + w_{\text{joint}} L_{\text{joint}}. \quad (7)$$

Within the image tower, while image features alone cannot localize 3D objects, we have leveraged surface point geometry and camera intrinsic to have pseudo 3D votes that are useful approximations to the true 3D votes. So combining this image geometric cue with other semantic/texture cues we can still localize objects in 3D with image features only.

Note that, although the multi-tower structure introduces extra parameters, at inference time we no longer need to compute for the point cloud only and the image only towers – therefore there is minimal computation overhead.

²If the coordinates after projection is fractional, bi-linear interpolation is used.

| methods | RGB | bathtub | bed | bookshelf | chair | desk | dresser | nightstand | sofa | table | toilet | mAP |
|------------------|-----|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| DSS [45] | ✓ | 44.2 | 78.8 | 11.9 | 61.2 | 20.5 | 6.4 | 15.4 | 53.5 | 50.3 | 78.9 | 42.1 |
| COG [39] | ✓ | 58.3 | 63.7 | 31.8 | 62.2 | 45.2 | 15.5 | 27.4 | 51.0 | 51.3 | 70.1 | 47.6 |
| 2D-driven [19] | ✓ | 43.5 | 64.5 | 31.4 | 48.3 | 27.9 | 25.9 | 41.9 | 50.4 | 37.0 | 80.4 | 45.1 |
| PointFusion [52] | ✓ | 37.3 | 68.6 | 37.7 | 55.1 | 17.2 | 23.9 | 32.3 | 53.8 | 31.0 | 83.8 | 45.4 |
| F-PointNet [34] | ✓ | 43.3 | 81.1 | 33.3 | 64.2 | 24.7 | 32.0 | 58.1 | 61.1 | 51.1 | 90.9 | 54.0 |
| VOTENET [33] | ✗ | 74.4 | 83.0 | 28.8 | 75.3 | 22.0 | 29.8 | 62.2 | 64.0 | 47.3 | 90.1 | 57.7 |
| +RGB | ✓ | 70.0 | 82.8 | 27.6 | 73.1 | 23.2 | 27.2 | 60.7 | 63.7 | 48.0 | 86.9 | 56.3 |
| +region feature | ✓ | 71.7 | 86.1 | 34.0 | 74.7 | 26.0 | 34.2 | 64.3 | 66.5 | 49.7 | 88.4 | 59.6 |
| IMVOTENET | ✓ | 75.9 | 87.6 | 41.3 | 76.7 | 28.7 | 41.4 | 69.9 | 70.7 | 51.1 | 90.5 | 63.4 |

Table 1. **3D object detection results on SUN RGB-D v1 val set.** Evaluation metric is average precision with 3D IoU threshold 0.25 as proposed by [43]. Note that both COG [39] and 2D-driven [19] use room layout context to boost performance. The evaluation is on the SUN RGB-D v1 data for fair comparisons.

4. Experiments

In this section, we first compare our model with previous state-of-the-art methods on the challenging SUN RGB-D dataset (Sec. 4.1). Next, we provide visualizations of detection results showing how image information helps boost the 3D recognition (Sec. 4.2). Then, we present an extensive set of analytical experiments to validate our design choices (Sec. 4.3). Finally, we test our method in the conditions of very sparse depth, and demonstrate its robustness (Sec. 4.4) in such scenarios.

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

Benchmark dataset. We use SUN RGB-D [42, 13, 51, 43] as our benchmark for evaluation, which is a single-view³ RGB-D dataset for 3D scene understanding. It consists of ~10K RGB-D images, with ~5K for training. Each image is annotated with amodal oriented 3D bounding boxes. In total, 37 object categories are annotated. Following standard evaluation protocol [45], we only train and report results on the 10 most common categories. To feed the data to the point cloud backbone network, we convert the depth images to point clouds using the provided camera parameters. The RGB image is aligned to the depth channel and is used to query corresponding image regions from scene 3D points.

Methods in comparison. We compare IMVOTENET with previous methods that use both geometry and RGB. Moreover, since previous state-of-the-art (VOTENET [33]) used only geometric information, to better appreciate the improvement due to our proposed fusion and gradient blending modules we add two more strong baselines by extending the basic VOTENET with additional features from image.

Among the previous methods designed for RGB-D, 2D-driven [19], PointFusion [52] and F-PointNet [34] are all cascaded systems that rely on 2D detectors to provide proposals for 3D. Deep Sliding Shapes [45] designs a Faster

R-CNN [38] style 3D CNN network to generate 3D proposals from voxel input and then combines 3D and 2D RoI features for box regression and classification. COG [39] is a sliding shape based detector using 3D HoG like feature extracted from RGB-D data.

As for the variations of VOTENET [33], the first one, ‘+RGB’, directly appends the the RGB values as a three-dimensional vector to the point cloud features (of the seed points). For the second one (‘+region feature’), we use the same pre-trained Faster R-CNN (as in our model) to obtain the region-level one-hot class confidence feature, and concatenate it to the seed points inside that 2D box frustum. These two variations can also be viewed as ablated versions of our method.

Results. Table 1 shows the per-class 3D object detection results on SUN RGB-D. We can see that our model outperforms all previous methods by large margins. Especially, it improves upon the previously best model VOTENET by **5.7** mAP, showing effectiveness of the lifted 2D image votes. It gets better results on nearly all categories and has the biggest improvements on object categories that are often occluded (+12.5 AP for bookshelves) or geometrically similar to the others (+11.6 AP for dressers and +7.7 AP for nightstands).

Compared to the variations of the VOTENET that also uses RGB data, our method also shows significant advantages. Actually we find that naively appending RGB values to the point features resulted in *worse* performance, likely due to the over-fitting on RGB values. Adding region features as a one-hot score vector helps a bit but is still inferior compared to our method that more systematically leverage image votes.

4.2. Qualitative Results and Discussion

In Fig. 4, we highlight detection results of both the original VOTENET [33] (with only point cloud input) and our IMVOTENET with point cloud plus image input, to show how image information can help 3D detection in various ways. The first example shows how 2D object localization

³We do not evaluate on the ScanNet dataset [5] as in VOTENET because ScanNet involves *multiple* 2D views for each reconstructed scene – thus requires extra handling to merge multi-view features.

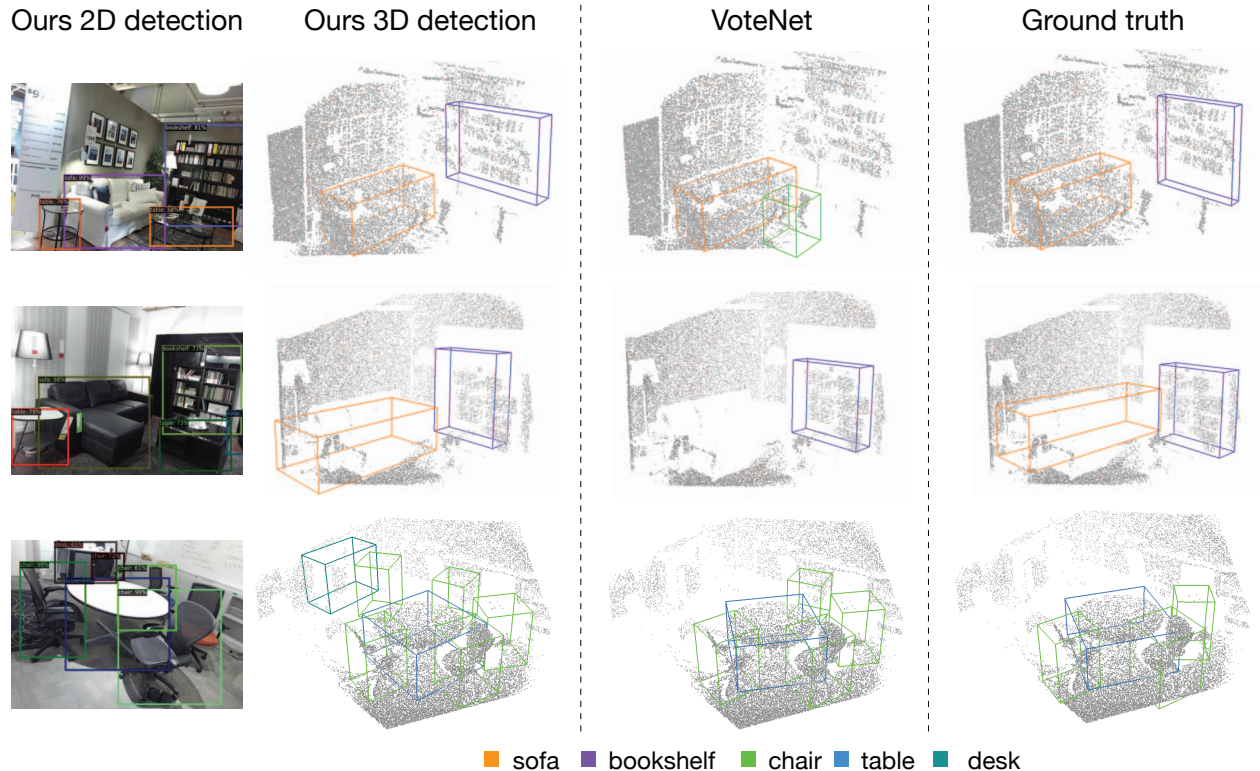


Figure 4. **Qualitative results showing how image information helps.** First row: the bookshelf is detected by IMVOTENET thanks to the cues from the 2D detector; Second row: the black sofa has barely any depth points due to its material, but leveraging images, we can detect it; Third row: with 2D localization cues and semantics, we detect the desk and chairs in the back which are even missed by ground truth annotations. Best viewed in color with zoom in.

and semantic help. We see a cluttered bookshelf that was missed by the VOTENET but thanks to the 2D detection in the images, we have enough confidence to recognize it in our network. The image semantics also help our network to avoid the false positive chair as that in the VOTENET output (coffee table and candles confused the network there). The second example shows how images can compensate depth sensor limitations. Due to the color and material of the black sofa, there is barely any depth point captured for it. While VOTENET completely misses the sofa, our network is able to pick it up. The third example shows how image cues can push the limit of 3D detection performance, by recovering far away objects (the desk and chairs in the back) that are even missed in the ground truth annotations.

4.3. Analysis Experiments

In this subsection, we show extensive ablation studies on our design choices and discuss how different modules affect the model performance. For all experiments we report mAP@0.25 on SUN RGB-D as before.

Analysis on geometric cues. To validate that geometric cues lifted from 2D votes help, we ablate geometric features (as in Eq. 6) passed to the 3D seed points in Table 2a.

We see that from row 1 to row 3, not using any 2D geometric cue results in a 2.2 point drop. On the other hand, not using the ray angle resulted in a 1.2 point drop, indicating the ray angle helps provide corrective cue to the pseudo 3D votes.

Analysis on semantic cues. Table 2b shows how different types of region features from the 2D images affect 3D detection performance. We see that the one-hot class score vector (probability score for the detected class, other classes set to 0), though simple, leads to the best result. Directly using the 1024-dim ROI features from the Faster R-CNN network actually got the worst number likely due to the optimization challenge to fuse this high-dim feature with the rest point features. Reducing the 1024-dim feature to 64-dim helps but is still inferior to the simple one-hot score feature.

Analysis on texture cues. Table 2c shows how different low-level image features (texture features) affect the end detection performance. It is clear that the raw RGB features are already effective while the more sophisticated per-pixel CNN features (from feature pyramids [23] of the Faster R-CNN detector) actually hurts probably due to over-fitting. More details are in the supplementary material.

| geometric cues | | mAP |
|----------------|-----------|------|
| 2D vote | ray angle | |
| ✓ | ✓ | 63.4 |
| ✓ | ✗ | 62.2 |
| ✗ | ✗ | 61.2 |

(a) Ablation studies on 2D **geometric** cues. 2D vote means the lifted 2D vote (2-dim) as in Eq. 6 and ray angle means the direction of $\overrightarrow{OC'}$ (3-dim). Both geometric cues helped our model.

| semantic cues | | mAP |
|----------------|--------|------|
| region feature | # dims | |
| one-hot score | 10 | 63.4 |
| RoI [38] | 64 | 62.4 |
| | 1024 | 59.5 |
| ✗ | - | 58.9 |

(b) Ablation studies on 2D **semantic** cues. Different region features are experimented. This includes simple one-hot class score vector and rich RoI features. The former (default) works best.

| texture cues | | mAP |
|-----------------|--------|------|
| pixel feature | # dims | |
| RGB | 3 | 63.4 |
| FPN- P_2 [23] | 256 | 62.0 |
| FPN- P_3 | 256 | 62.0 |
| ✗ | - | 62.4 |

(c) Ablation studies on 2D **texture** cues. We experiment with different pixel-level features including RGB values (default) and learned representations from the feature pyramid.

Table 2. **Ablation analysis** on 2D cues. We provide detailed analysis on all types of features from 2D (see Sec. 3.2 for detailed descriptions).

| tower weights | | | mAP | | |
|------------------|--------------------|--------------------|-------------|-------------|-------------|
| w_{img} | w_{point} | w_{joint} | image | point cloud | joint |
| - | - | - | 46.8 | 57.4 | 62.1 |
| 0.1 | 0.8 | 0.1 | 46.9 | 57.8 | 62.7 |
| 0.8 | 0.1 | 0.1 | 46.8 | 58.2 | 63.3 |
| 0.1 | 0.1 | 0.8 | 46.1 | 56.8 | 62.7 |
| 0.3 | 0.3 | 0.4 | 46.6 | 57.9 | 63.4 |

Table 3. Analysis on **multi-tower training**. In the first block we show performance *without* blending in gray. Then we show the setting where each of the tower dominates (0.8) the overall training. Finally we show our default setting where weights are more balanced.

Gradient blending. Table 3 studies how tower weights affect the gradient blending training. We ablate with a few sets of representative weights ranging from single tower training (the first row), dominating weights for each of the tower (2nd to 4th rows) and our best set up. It is interesting to note that even with just the image features (the 1st row, 4th column) i.e. the pseudo votes and semantic/texture cues from the images, we can already outperform several previous methods (see Table 1), showing the power of our fusion and voting design.

4.4. Detection with Sparse Point Clouds

While depth images provide dense point clouds for a scene (usually 10k to 100k points), there are other scenarios that only sparse points are available. One example is when the point cloud is computed through visual odometry [29] or Structure from Motion (SfM) [17] where 3D point positions are triangulated by estimating poses of a monocular camera in multiple views. With such sparse data, it is valuable to have a system that can still achieve decent detection performance.

To analyze the potential of our model with sparse point clouds, we simulate scans with much less points through two types of point sub-sampling: uniformly random sub-sampling (remove existing points with a uniform distribution) and ORB [40] key-point based sub-sampling (sam-

| point cloud settings | | mAP | | |
|----------------------|----------|-------------|-------|----------|
| sampling method | # points | point cloud | joint | Δ |
| random uniform | 20k | 57.7 | 63.4 | +5.7 |
| | 5k | 56.2 | 61.7 | +5.5 |
| | 1k | 49.6 | 58.5 | +8.9 |
| ORB [40] | 5k | 32.4 | 49.9 | +16.5 |
| | 1k | 27.9 | 47.1 | +19.2 |

Table 4. **Sparse point cloud** experiment, where we sub-sample the number of points in the cloud either via *random uniform* sampling or with *ORB key points* [40]. In such cases, our IMVOTENET significantly outperforms purely geometry based VOTENET.

ple ORB key points on the image and only keep 3D points that project close to those 2D key points). In Table 4, we present detection results with different distribution and density of point cloud input. We see that in the column of “point cloud”, with decreased number of points, 3D detection performance quickly drops. On the other hand, we see including image cues significantly improves performance. This improvement is most significant when the sampled points are from ORB key points that are more non-uniformly distributed.

5. Conclusion

In this work we have explored how image data can assist a voting-based 3D detection pipeline. The VOTENET detector we build upon relies on a voting mechanism to effectively aggregate geometric information in point clouds. We have demonstrated that our new network, IMVOTENET, can leverage extant image detectors to provide both geometric and semantic/texture information about an object in a format that can be integrated into the 3D voting pipeline. Specifically, we have shown how to lift 2D geometric information to 3D, using knowledge of the camera parameters and pixel depth. IMVOTENET significantly boosts 3D object detection performance exploiting multi-modal training with gradient blending, especially in settings when the point cloud is sparse or unfavorably distributed.

References

- [1] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *ICCV*, 2015. 3
- [2] Armen Avetisyan, Manuel Dahnert, Angela Dai, Manolis Savva, Angel X. Chang, and Matthias Nießner. Scan2cad: Learning cad model alignment in rgb-d scans. In *CVPR*, 2019. 2
- [3] Dana H Ballard. Generalizing the hough transform to detect arbitrary shapes. *Pattern recognition*, 13(2):111–122, 1981. 3
- [4] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. *arXiv preprint arXiv:1904.08755*, 2019. 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 *CVPR*, 2017. 6
- [6] Zhuo Deng and Longin Jan Latecki. Amodal detection of 3d objects: Inferring 3d bounding boxes from 2d ones in rgb-depth images. In *CVPR*, volume 2, 2017. 3
- [7] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. *arXiv preprint arXiv:1606.01847*, 2016. 3
- [8] Ruohan Gao and Kristen Grauman. 2.5 d visual sound. In *CVPR*, 2019. 3
- [9] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *CVPR*, 2018. 2
- [10] Alexander Hermans, Georgios Floros, and Bastian Leibe. Dense 3d semantic mapping of indoor scenes from rgb-d images. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2631–2638. IEEE, 2014. 3
- [11] Ji Hou, Angela Dai, and Matthias Nießner. 3d-sis: 3d semantic instance segmentation of rgb-d scans. In *CVPR*, 2019. 2, 3
- [12] Paul VC Hough. Machine analysis of bubble chamber pictures. In *Conf. Proc.*, 1959. 2
- [13] Allison Janoch, Sergey Karayev, Yangqing Jia, Jonathan T Barron, Mario Fritz, Kate Saenko, and Trevor Darrell. A category-level 3-d object dataset: Putting the kinect to work. In *ICCV Workshops (ICCV Workshops)*, 2011. 6
- [14] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, 2017. 3
- [15] 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
- [16] Jan Knopp, Mukta Prasad, Geert Willems, Radu Timofte, and Luc Van Gool. Hough transform and 3d surf for robust three dimensional classification. In *ECCV*. 2010. 2
- [17] Jan J Koenderink and Andrea J Van Doorn. Affine structure from motion. *JOSA A*, 8(2):377–385, 1991. 8
- [18] Jason Ku, Melissa Mozifian, Jungwook Lee, Ali Harakeh, and Steven L Waslander. Joint 3d proposal generation and object detection from view aggregation. In *IROS*. IEEE, 2018. 2, 3
- [19] Jean Lahoud and Bernard Ghanem. 2d-driven 3d object detection in rgb-d images. In *CVPR*, pages 4622–4630, 2017. 2, 3, 6
- [20] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, 2019. 2
- [21] 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
- [22] Ming Liang, Bin Yang, Shenlong Wang, and Raquel Urtasun. Deep continuous fusion for multi-sensor 3d object detection. In *ECCV*, pages 641–656, 2018. 3
- [23] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 5, 7, 8
- [24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 5
- [25] 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
- [26] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *ECCV*, pages 21–37. Springer, 2016. 5
- [27] John McCormac, Ankur Handa, Andrew Davison, and Stefan Leutenegger. Semanticfusion: Dense 3d semantic mapping with convolutional neural networks. In *2017 IEEE International Conference on Robotics and automation (ICRA)*, pages 4628–4635. IEEE, 2017. 3
- [28] Liangliang Nan, Ke Xie, and Andrei Sharf. A search-classify approach for cluttered indoor scene understanding. *ACM TOG*, 31(6), 2012. 2
- [29] David Nistér, Oleg Naroditsky, and James Bergen. Visual odometry. In *CVPR*, 2004. 8
- [30] Andrew Owens, Phillip Isola, Josh McDermott, Antonio Torralba, Edward H Adelson, and William T Freeman. Visually indicated sounds. In *CVPR*, 2016. 3
- [31] Andrew Owens, Jiajun Wu, Josh H McDermott, William T Freeman, and Antonio Torralba. Ambient sound provides supervision for visual learning. In *ECCV*. Springer, 2016. 3
- [32] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *AAAI*, 2018. 3
- [33] Charles R. Qi, Or Litany, Kaiming He, and Leonidas J. Guibas. Deep hough voting for 3d object detection in point clouds. 2019. 1, 2, 3, 4, 5, 6
- [34] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *CVPR*, 2018. 2, 3, 6
- [35] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *CVPR*, 2017. 2

- [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. 2, 3
- [37] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 5
- [38] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 2, 4, 5, 6, 8
- [39] Zhile Ren and Erik B Sudderth. Three-dimensional object detection and layout prediction using clouds of oriented gradients. In *CVPR*, 2016. 6
- [40] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary R Bradski. Orb: An efficient alternative to sift or surf. In *ICCV*, 2011. 8
- [41] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointcnn: 3d object proposal generation and detection from point cloud. *arXiv preprint arXiv:1812.04244*, 2018. 1, 2
- [42] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. *ECCV*, pages 746–760, 2012. 6
- [43] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *CVPR*, 2015. 2, 6
- [44] Shuran Song and Jianxiong Xiao. Sliding shapes for 3d object detection in depth images. In *ECCV*. Springer, 2014. 2
- [45] Shuran Song and Jianxiong Xiao. Deep sliding shapes for amodal 3d object detection in rgb-d images. In *CVPR*, 2016. 2, 3, 6
- [46] Minhuk Sung, Vladimir G Kim, Roland Angst, and Leonidas Guibas. Data-driven structural priors for shape completion. *ACM TOG*, 2015. 5
- [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] Chen Wang, Danfei Xu, Yuke Zhu, Roberto Martín-Martín, Cewu Lu, Li Fei-Fei, and Silvio Savarese. Densefusion: 6d object pose estimation by iterative dense fusion. In *CVPR*, pages 3343–3352, 2019. 2, 3
- [49] Weiyao Wang, Du Tran, and Matt Feiszli. What makes training multi-modal networks hard? *arXiv preprint arXiv:1905.12681*, 2019. 2, 3, 5
- [50] Oliver J Woodford, Minh-Tri Pham, Atsuto Maki, Frank Perbet, and Björn Stenger. Demisting the hough transform for 3d shape recognition and registration. *IJCV*, 106(3):332–341, 2014. 2
- [51] Jianxiong Xiao, Andrew Owens, and Antonio Torralba. Sun3d: A database of big spaces reconstructed using sfm and object labels. In *ICCV*. IEEE, 2013. 6
- [52] Danfei Xu, Dragomir Anguelov, and Ashesh Jain. Pointfusion: Deep sensor fusion for 3d bounding box estimation. In *CVPr*, pages 244–253, 2018. 2, 3, 6
- [53] Li Yi, Wang Zhao, He Wang, Minhuk Sung, and Leonidas Guibas. Gspn: Generative shape proposal network for 3d instance segmentation in point cloud. *arXiv preprint arXiv:1812.03320*, 2018. 2
- [54] Zhou Yu, Jun Yu, Chenchao Xiang, Jianping Fan, and Dacheng Tao. Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering. *TNN*, 2018. 3
- [55] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *ECCV*, 2018. 3
- [56] Yin Zhou and Oncel Tuzel. Voxnet: End-to-end learning for point cloud based 3d object detection. In *CVPR*, 2018. 1, 2