# Single Image 3D Object Estimation with Primitive Graph Networks

Qian He[1,2,3], Desen Zhou[4], Bo Wan[5], Xuming He[1,6*]

[1]School of Information Science and Technology, ShanghaiTech University
[2]Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences
[3]University of Chinese Academy of Sciences
[4]Department of Computer Vision Technology (VIS), Baidu Inc.
[5]Department of Electrical Engineering (ESAT), KU Leuven
[6]Shanghai Engineering Research Center of Intelligent Vision and Imaging
{heqian,wanbo,hexm}@shanghaitech.edu.cn,zhoudesen@baidu.com

## ABSTRACT

Reconstructing 3D object from a single image (RGB or depth) is a fundamental problem in visual scene understanding and yet remains challenging due to its ill-posed nature and complexity in real-world scenes. To address those challenges, we adopt a primitive-based representation for 3D object, and propose a two-stage graph network for primitive-based 3D object estimation, which consists of a sequential proposal module and a graph reasoning module. Given a 2D image, our proposal module first generates a sequence of 3D primitives from input image with local feature attention. Then the graph reasoning module performs joint reasoning on a primitive graph to capture the global shape context for each primitive. Such a framework is capable of taking into account rich geometry and semantic constraints during 3D structure recovery, producing 3D objects with more coherent structure even under challenging viewing conditions. We train the entire graph neural network in a stage-wise strategy and evaluate it on three benchmarks: Pix3D, ModelNet and NYU Depth V2. Extensive experiments show that our approach outperforms the previous state of the arts with a considerable margin.

## CCS CONCEPTS

• **Computing methodologies** → **Scene understanding**; **Reconstruction**.

## KEYWORDS

3D object estimation; part-based object representation; graph neural networks
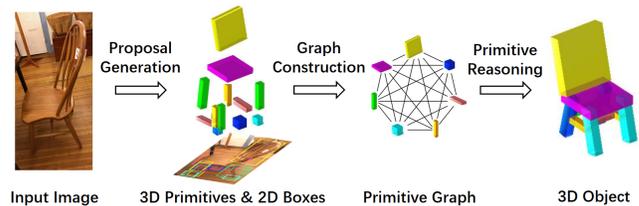
---

---

**Figure 1: Given an input image, our network first generates primitive proposals attending to local features, then builds a graph on the proposals for primitive reasoning, and finally outputs a 3D object shape in primitive representation.**

## 1 INTRODUCTION

As a fundamental problem in visual scene understanding, 3D object recovery from a monocular image (RGB or depth) plays an indispensable role in many vision tasks when multi-view sensors are difficult to deploy. It has also gained increasing attention in multimedia community [15, 26, 27, 41, 50, 52, 53, 55] due to its great potential in VR/AR content generation and robot-scene interaction. However, unlike 3D reconstruction in multi-view settings, single-image 3D object estimation remains highly challenging due to the real-world imaging process. A variety of factors, such as (self-)occlusion, varying viewing angles and large variation in object appearance/shapes, lead to a complex mapping from input observation to 3D structure.

Extensive effort has been made recently to address this problem via learning a deep neural network to map the input image to an output 3D shape. One dominating strategy represents the 3D object shapes in voxel [4, 47], point cloud [7] or mesh [19, 42]. However, such holistic geometric representations are sensitive to noisy inputs and lack the capacity of encoding higher-level regularity in object shapes, such as part symmetry or functionality.

The primitive-based methods, which exploits the compositionality in object shapes [1], provide an alternative solution that has potential to overcome the above limitations. They typically employ deep networks with an encoder-decoder architecture, which first extract a holistic feature representation of the input and then

decode it into a collection of shape primitives. For instance, 3D-PRNN [56] and PQ-Net [49] predicts a sequence of 3D rectangles, while Im2Struct [31] generates a tree-structured representation of object parts with a recursive neural network. Most existing primitive-based approaches, however, suffer from several limitations in practice: First, their inference process typically generates primitives in a progressive manner, and hence are prone to error propagation in part estimation. Additionally, their decoders rely on a holistic image representation, which is sensitive to occlusion. Furthermore, part semantics are largely ignored in geometric modeling, often leading to unreasonable 3D object configurations.

In this paper, we propose a novel primitive-based 3D object recovery framework to address the afore-mentioned limitations. Our main idea is to use a two-stage strategy, in which we first generate a pool of primitive proposals from the input image and then build a fully-connected primitive graph to perform joint reasoning on object configurations (see Fig. 1). Our graph-based primitive inference allows us to better capture the global shape context, mitigate the errors in proposal generation and predict the primitives coherently. Besides, each primitive proposal can attend to its local features and overall the object recovery is less susceptible to occlusion. Finally, with additional semantic part annotations, our method can easily incorporate semantic part relations into the primitive reasoning, which enforces stronger structure constraints for 3D estimation.

To achieve this, we develop a deep graph neural network that takes a monocular RGB or depth image of an object as input and generate its 3D primitive representation. Our model consists of two main modules, a primitive proposal network and a primitive reasoning network. The proposal module first computes a feature map by a deep Convnet and then adopts a recurrent network to lift the 2D feature into a sequence of 3D rectangular primitive proposals. The reasoning module builds a fully-connected graph network on the proposals, where each node embeds the geometry and semantic features of a primitive and each edge captures the relative geometric and semantic relations. The graph network then refines the primitive features by performing message-passing to capture the global context of object shape. Based on those refined representations, we finally predict the confidence scores (foreground vs. background or part label) and the geometric parameters of all the primitives. By exploiting its modular architecture, we train the entire primitive model with a stage-wise strategy, which simplifies the learning procedure and yet works effectively in practice.

We evaluate our method on three benchmarks: a subset of Pix3D [37], ModelNet [56] and NYU Depth V2 [36]. The empirical results and ablation study show that our method consistently outperforms the prior state of the art [49, 56] with a considerable margin. The main contributions of our work are three-folds:

- We propose a new primitive-based method for 3D object recovery from a single image, achieving state-of-the-art performance on three benchmarks.
- We introduce a modular graph neural network that effectively captures 3D shape constraints and performs joint reasoning on the primitives for coherent 3D estimation.
- We enrich the primitive representation by its conv features for robustness in occlusion and optionally by its semantic part cues to enforce stronger constraints on object structure.

## 2 RELATED WORK

### 2.1 Single-view 3D Object Estimation

Recovering 3D objects from single images is an ill-posed problem. Previous works leverage rich CAD models and tackle this problem by joint analysis of images and 3D shapes [16, 24], exploring local correspondence between images and 3D shape [21], or applying deformation from a mean shape [17, 18]. More recently, extensive approaches for single image 3D object estimation explore various shape representation, such as key-point [14, 22, 38, 46], voxel [4, 47], point cloud [7], mesh [19, 42] and implicit functions [3, 28, 32]. Those holistic representations, however, cannot capture part-level object structure, which lack fine-grained shape constrains and are sensitive to occlusion.

Primitive-based methods address those challenges by adopting a part-based representation. [56] generates a primitive sequence directly from encoded image feature. [31] explicitly encodes shape structure like adjacency and symmetry based on a binary tree of 3D boxes. [33] also learns to recover a binary tree of part-based shape decomposition. Nevertheless, these methods generate object parts in a progressive manner and hence are prone to error propagation. More recently, [5] represents shapes via a learnable convex decomposition and [2] proposes to generate shapes via binary shape partitioning, both utilizing a fully-connected network to predict all object parts at once. By contrast, we explicitly model part relations with a non-local graph neural network [43], which enables us to incorporate a fine-grained shape regularity constraint, and to explore the guidance of part semantic relations in shape recovery with additional supervision.

Most existing methods adopt an encoder-decoder network architecture, which either directly encode the entire image into a feature representation [10], or resort to 2.5D intermediate representations such as silhouette, depth and normal map [45, 48] for better input features, and then map it to 3D shape output. However, such holistic features are susceptible to occlusion. In contrast, we propose to use attention mechanism to extract local features based on 2D bounding box predictions for object parts.

### 2.2 Part-based 3D Shape Representation

Representing objects with parts and discovering object structure is a long-standing research problem [12, 29, 44, 54]. Recent approaches learn a part-based object representation by mapping a shape volume to a set of primitive boxes [40], or superquadrics [34], or 3D Gaussians [8, 9], but largely ignore object part relations and high-level shape regularity. [23] designs a recursive autoencoder in a binary tree to generate object parts, which explicitly captures part symmetry and adjacency. StructureNet [30] enriches this tree structure into an n-ary graph and introduces part semantics in shape generation. [39] represents 3D volumes with programs [35] and further combines semantic structure to capture shape regularity and stability of real-world objects. [49] proposes a sequence-to-sequence network for part-based object generation. However, those approaches typically generate part-based shapes in a progressive manner and hence lack global refinement on the 3D object shapes. Moreover, they often use a simple encoder-decoder strategy without grounding part-based object structure into input image, and have difficulty in parsing complex realistic scenes.
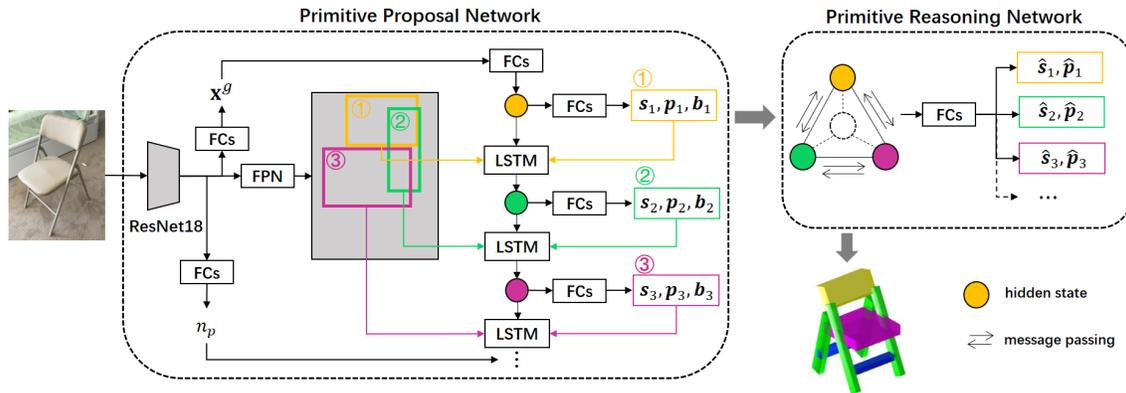
**Figure 2: Model framework.** Our model consists of two main modules: Primitive Proposal Network generates a sequence of proposals from the input image with LSTMs, and Primitive Reasoning Network is a fully-connected graph network built on the proposals to jointly reason the 3D shape configuration[1].

## 3 METHOD

### 3.1 Problem Setting

Given an input RGB or depth image $I$ of an object $O$, our goal is to estimate its full 3D shape in a form of a set of 3D geometric primitives $\mathcal{P} = \{\mathbf{p}_1, ..., \mathbf{p}_{N_O}\}$, where $N_O$ is the number of its primitives. In this work, we adopt a box representation for each primitive $\mathbf{p}_i = [l_x, l_y, l_z, t_x, t_y, t_z, \theta_x, \theta_y, \theta_z] \in \mathbb{R}^9$, where $[l_x, l_y, l_z]$ measures the length of three orthogonal edges of the box, $[t_x, t_y, t_z]$ represents the translation of the box center and $[\theta_x, \theta_y, \theta_z]$ denotes the rotated Euler angles along each axis.

In order to incorporate object part semantics, we also consider the problem setting of 3D recovery with additional semantic part annotations for each object category $c$, in which a semantic part label $s_i \in \{1, ..., M_c\}$ is assigned to each primitive $\mathbf{p}_i$. In the rest of the paper, we refer to this part-augmented scenario as the *semantic-aware setting* while denoting the more general scenario without such annotations (i.e., $M_c = 1$) as the *semantic-agnostic setting*.

### 3.2 Model Overview

To tackle the problem of 3D primitive estimation from a monocular image, we adopt a hypothesize-and-refine strategy, in which we first generate a set of 3D primitive proposals from the image and then jointly infer their 3D parameters and semantic part labels. We instantiate this two-stage strategy by developing a modular graph neural network consisting of two main modules: a primitive proposal network and a primitive reasoning network. The front-end proposal network extracts the set of initial primitive proposals from the input image, while the next reasoning network captures global context of the entire object and performs the refinement of primitive features, which are used for final prediction on their geometry and semantic labels. Fig. 2 shows an overview of our model architecture.

In the remaining of this section, we will introduce our network design in details. We start from the semantic-aware setting and first

introduce the primitive generation network in Sec. 3.3, followed by the primitive reasoning network in Sec 3.4. We then discuss our stage-wise model learning in Sec. 3.5. Finally, we generalize our method to the semantic-agnostic setting in Sec. 3.6.

### 3.3 Primitive Proposal Network

Our first step is to generate a set of 3D box primitives from the input image. Generating each primitive independently from a 2D input is a challenging task and hence we consider a structural approach to exploit the 3D shape prior of objects. To this end, we introduce a CNN-RNN model that predicts a sequence of primitive proposals each time. Such a sequence proposal generation mechanism captures the dependency between neighboring primitives, and is capable of producing more stable 3D candidate boxes.

Specifically, we propose a Primitive Proposal Network which consists of two submodules: a Backbone ConvNet and a Recurrent Generator. The Backbone ConvNet computes image features and estimates the number of primitives. Then a Recurrent Generator utilizes recurrent units to predicts primitive proposals, including their semantic labels and geometric parameters. Below we describe the details of those two submodules.

*3.3.1 Backbone ConvNet.* We use a ResNet18+FPN [25] as the backbone network to compute a 2D feature map $\mathbf{\Gamma}$ for later local conv feature extraction. Then we attach two sets of fully-connected layers (FCs) to the ResNet18, separately, to obtain a global image representation $\mathbf{x}^g \in \mathbb{R}^{256}$ and to directly regress the number of object primitives $N_O$ for this instance. The output of the regressor is denoted as $n_p$ after rounding.

*3.3.2 Recurrent Generator.* We adopt a light-weight two-layer LSTM network to generate primitive proposals in a sequential manner. At each step $i$, the LSTM updates its hidden state $\mathbf{h}_i \in \mathbb{R}^{D_h}$ and output three properties of the current primitive, including its 3D parameters $\mathbf{p}_i$, semantic label distribution $\mathbf{s}_i \in \mathbb{R}^{M_c}$, and its 2D bounding box location $\mathbf{b}_i \in \mathbb{R}^4$ in image.

The initial hidden state is computed from the global image feature $\mathbf{x}^g$ with an MLP as $\mathbf{h}_1 = f_h(\mathbf{x}^g)$, and the LSTM unit, denoted

---

[1]We only show one LSTM sequence in our proposal generation for clarity. Note that we also use global image feature $\mathbf{x}^g$ as input to LSTM at each step of proposal generation.

as $F_{\text{LSTM}}$, updates the hidden state as follows,

$$\mathbf{h}_{i+1} = F_{\text{LSTM}}(\mathbf{h}_i, \mathbf{x}_i, \mathbf{s}_i, \mathbf{p}_i), \quad 1 \le i < n_p \qquad (1)$$

where $\mathbf{x}_i$ is a set of input features including the global image feature $\mathbf{x}^g$, the bounding box location $\mathbf{b}_i$ and a local conv feature $\mathbf{v}_i$ extracted with ROI-Align [13] on the conv feature map $\mathbf{\Gamma}$.

Given the hidden state $\mathbf{h}_i$, the LSTM unit outputs three predictions on the current primitive with MLPs as follows,

$$\mathbf{s}_i = f_s(\mathbf{h}_i), \quad \mathbf{b}_i = f_b(\mathbf{h}_i), \quad \mathbf{p}_i = f_p(\mathbf{h}_i, \hat{s}_i), \quad \hat{s}_i = \arg\max(\mathbf{s}_i) \quad (2)$$

where $f_s(\cdot)$ is a single FC followed by a Softmax function and $f_b(\cdot)$, $f_p(\cdot)$ are multiple FCs with ReLU as activation function. We use a different MLP for each semantic class $s_i$ to predict primitives.

As generating 3D primitives from 2D images is particularly challenging, we build two separate LSTM networks to improve the recall rate of our proposal set. Specifically, the networks generate two primitive sequences in opposite directions for each input image (see Sec. 4.2). We discard the ordering information in two sequences, and take their union as the output of the Primitive Proposal Network. We denote the final proposal set as $Q = \{(\mathbf{p}_i, \mathbf{s}_i, \mathbf{h}_i)\}_{i=1}^{2n_p}$.

## 3.4 Primitive Reasoning Network

Given the generated proposal set $Q$, our second module aims to refine the proposals jointly by taking into account the global context of the entire object shape and to generate the final prediction of 3D primitives. To this end, we propose a fully-connected graph neural network [43] to update the primitive representations, based on which we then predict the parameters of 3D primitives.

Specifically, we first build a primitive graph $\mathcal{G}$ whose nodes are primitive proposals and edges are fully connected. We then associate the $i$-th primitive to its corresponding node and denote its feature embedding as $\mathbf{z}_i \in \mathbb{R}^{D_z}$. The initial feature embedding is computed from the proposal features:

$$\mathbf{z}_i = [g_h(\mathbf{h}_i); g_s(\mathbf{s}_i); g_p(\mathbf{p}_i)], \quad 1 \le i \le 2n_p \qquad (3)$$

where $g_h(\cdot), g_s(\cdot),$ and $g_p(\cdot)$ are a FC layer followed by a ReLU function, respectively. $[;]$ denotes concatenation.

Our primitive reasoning network uses a one-step message passing to update the primitive feature embedding for capturing the global context of each object instance. Concretely, given a graph node $\mathbf{z}_i$, we aggregate information from all the remaining primitives with an attention mechanism and update the node feature with a residual block:

$$\alpha_{ij} = (\mathbf{U}\mathbf{z}_i)^\top (\mathbf{V}\mathbf{z}_j), \quad \mathbf{y}_i = \mathbf{z}_i + \frac{1}{2n_p - 1} \sum_{\forall j \ne i} \alpha_{ij} \mathbf{W}\mathbf{z}_j \qquad (4)$$

where $\mathbf{U}, \mathbf{V}, \mathbf{W} \in \mathbb{R}^{D_z \times D_z}$ are embedding matrices and $\alpha_{ij}$ is the importance weight of each input message, and $\mathbf{y}_i$ is the updated primitive representation[2].

*3.4.1 Model Prediction.* Given the updated primitive representations, we now refine the semantic label prediction and geometric parameters of all the proposals. To remove the false positive primitives, we augment the semantic label space with a background class 0 so that the semantic class of the $i$-th primitive $s_i \in \{0, 1, ..., M_c\}$.

The final prediction of the semantic label and the primitive parameters are estimated based on two FC networks as follows:

$$\hat{\mathbf{s}}_i = f_s^o(\mathbf{y}_i), \quad \hat{\mathbf{p}}_i = f_p^o(\mathbf{y}_i, \hat{s}_i), \quad \hat{s}_i = \arg\max \hat{\mathbf{s}}_i \wedge \hat{s}_i > 0 \quad (5)$$

where $f_s^o(\cdot)$ is a single FC layer followed by a softmax function and $f_p^o(\cdot)$ is a set of MLPs with ReLU as activation function, one for each part semantic class.

During testing, we apply non-maximum suppression on the final primitive predictions. Specifically, we first match all primitives into pairs according to their L1 distance. Each time we match a pair with minimum distance in the remaining unpaired primitives. After the matching process, we take the primitive with higher class confidence score in each pair to be our final outputs. We use pairing for NMS as we have two LSTM models to generate proposals in two directions. Each LSTM typically generates non-overlapping primitives due to sequential dependency, and we only need to remove potential redundant proposals from two LSTMs.

## 3.5 Model Learning

Based on our modular architecture, we adopt a stage-wise strategy in our model learning. Below we will introduce the training process of the Primitive Proposal Network and the Primitive Reasoning Network in turn. For notation clarity, we define our training loss on an input image of object $O$.

*3.5.1 Training Primitive Proposal Network.* For our backbone network, we fix the ResNet18 with pre-trained parameters and finetune the parameters of the FPN module with the rest of the proposal network. The training of the regressor and the recurrent generator are decoupled due to the fixed base network.

We first train the regressor with an $l_1$ loss on the predicted number of primitives, which is then applied to the next stage for training primitive reasoning module. For the recurrent generator, we assume the primitives are generated with a predefined ordering and length $N_O$ as specified in the ground-truth (see Sec. 4.2). The overall loss for the proposal network consists of three terms: a semantic label classification loss, a primitive regression loss and a box regression loss, which are defined as follows:

$$L_p = \sum_{i=1}^{2N_O} L_{ce}(\mathbf{s}_i, s_i^*) + \lambda_{rp} \sum_{i=1}^{2N_O} L_{l_1}(\mathbf{p}_i, \mathbf{p}_i^*) + \lambda_{rb} \sum_{i=1}^{2N_O} L_{sm}(\mathbf{b}_i, \mathbf{b}_i^*)$$
$$(6)$$

where $s_i^*, \mathbf{p}_i^*$ and $\mathbf{b}_i^*$ are the ground-truth semantic label, primitive and box location, $L_{ce}, L_{l_1}$ and $L_{sm}$ are Cross Entropy loss, L1 distance and Smooth-L1 loss, and $\lambda_{rp}$ and $\lambda_{rb}$ are the corresponding weights to primitive regression loss and box regression loss, respectively.

*3.5.2 Training Primitive Reasoning Network.* We train the primitive reasoning network with a multi-task loss similar to the object detection systems. Specifically, given $2n_p$ final predictions, we match them to $2N_O$ ground-truth primitives (doubled) in a greedy manner: Each time we first compute the $l_1$ distance between all remained prediction and ground-truth primitive pairs and choose the pair with minimum distance. Then we remove the matched pair from remained primitives and repeat matching, until no predictions or no ground-truth primitives are left. Finally, if any predictions remain unmatched, we assign them to background class. Empirically, we

---

[2]We note that while it is straightforward to extend the message passing to multiple iterations, we found empirically that single iteration achieves the best performance.

do not include rotation in this $l_1$ distance computation, as rotation predictions are relatively less reliable than primitive translation and edge length.

The total loss for the Primitive Reasoning Network has two terms, including a semantic classification loss and a primitive regression loss, which can be written as:

$$L_r = \sum_{i=1}^{2n_p} L_{ce}(\hat{s}_i, \hat{s}_i^*) + \lambda_{rp} \sum_{i=1}^{2n_p} L_{l_1}(\hat{p}_i, \hat{p}_i^*) [\![\hat{s}_i^* > 0]\!] \tag{7}$$

Where $\hat{s}_i^*$ and $\hat{p}_i^*$ are the ground-truth semantic label and primitive respectively, $[\![\cdot]\!]$ is the indicator function, and the weight $\lambda_{rp}$ of primitive regression loss is the same as in Eq. 6.

## 3.6 Semantic-agnostic Model

For the semantic-agnostic setting, it is straightforward to extend our model formulation by setting $M_c = 1$, which indicates the semantic part annotations are unavailable.

In addition, we simplify our model slightly to speed up model inference and learning. Specifically, in the Primitive Proposal Network, as $M_c = 1$, we set the output function for the semantic label $f_s$ as a constant, i.e., $f_s(\cdot) \equiv 1$ and remove the semantic label classification loss in Eq. 6 in the training stage. Similarly, in the Primitive Reasoning Network, we set the embedding function for the semantic label $g_s$ as a constant, i.e., $g_s(\cdot) \equiv 1$. We note that the rest of our model pipeline and training process are the same as the semantic-aware setting.

## 4 EXPERIMENT

We train and evaluate our model on Pix3D dataset [37] with real-world RGB images and ModelNet dataset [51] with synthetic depth images in both semantic-aware and semantic-agnostic setting. In addition, we directly apply the models trained on ModelNet to NYU Depth V2 [36] dataset to demonstrate the robustness of our method. On each dataset, we compare with the state-of-the-art methods 3D-PRNN [56], Tulsiani et al. [40], and PQ-Net [49].

Below we first introduce the datasets and evaluation protocol in Sec. 4.1 and implementation details in Sec. 4.2. Then we present results comparing to other methods on three benchmarks in Sec. 4.3. To illustrate the effectiveness of different components of our model design, we conduct comprehensive ablation study on Pix3D *chair* class in Sec. 4.4. Finally, we demonstrate further analysis of our method in supplementary, regarding our proposal quality, robustness of our method, and our failure cases.

### 4.1 Datasets and Metrics

**Pix3D** Pix3D [37] is a challenging benchmark consisting of pixel-level aligned 3d shapes and real-world images with various occlusion and background clutter. We conduct our experiments on its three major categories: *chair*, *table* and *sofa*, which contain 3839, 1870 and 1947 images, as well as 221, 63 and 20 unique 3D models, respectively. We take all models in voxel of $32 \times 32 \times 32$ and manually label them into well-aligned primitives with corresponding semantic class. We assign class label to each part based on its functionality, such as cushion to sit on and legs to support. Given camera pose and primitive annotations, we obtain 2D part bounding boxes from reprojected 3D primitives. We divide all parts of

*chair* into six semantic classes, *table* into four, and *sofa* into four. The maximum number of primitives in an object for *chair*, *table*, and *sofa* are 15, 8, and 8. We randomly sample and fix 20% of images to be our test set and use the rest as our training set. Note that our test and train splits have no shared 3D model annotations.

Our semi-automatic manual labeling of primitives consists of two steps: First, we utilize the same preprocessing toolbox as in the 3D-PRNN to generate a coarse primitive (oriented bounding box) configuration for each object volume in $32 \times 32 \times 32$ from original ground truth volume annotations. The toolbox is based on an energy minimization procedure that resembles Iterative Closest Point (ICP). In the second step, we then slightly adjust the parameters of those primitives manually by visually inspecting their alignment with the ground truth volume through a graphical user interface.

**ModelNet** We follow 3D-PRNN and apply our method on three categories of ModelNet [51]: *chair*, *table* and *night stand*. For each category, we utilize 889, 392 and 200 models for training and 100, 100 and 82 models for testing, respectively. We adopt the same data synthesis strategy as in 3D-PRNN to get input depth maps of size $64 \times 64$. For each object, we sample five views on a unit sphere by rejection-sampling, bounded within 20 degrees of the equator, and compute depth by projection of meshes. The primitive annotations are from 3D-PRNN and we augment each primitive with a semantic class label. We divide all parts of *chair* into six classes, *table* into four, and *night stand* into four. The maximum number of primitives in an object for *chair*, *table*, and *night stand* are 16, 11, and 10.

**NYU Depth V2** NYU Depth V2 [36] is a very challenging indoor scene dataset consisting of real-world RGBD images. We use annotations from [11], with 30 CAD models of six furniture categories: *chair*, *table*, *desk*, *bed*, *bookshelf* and *sofa*, and extruded polygons for other objects. As primitive annotations are not available, we directly apply the models trained on ModelNet to the test split of NYU Depth V2 and compare the results in voxel.

**Evaluation Protocol** For our semantic-aware setting, we train a class-specific model with semantic supervision for each object category. While for our semantic-agnostic setting, we only train one class-agnostic model with no semantics on all object categories.

We compare with the state-of-the-art methods in two different settings. Firstly, we compare our method with 3D-PRNN [56] and Tulsiani et al. [40], for estimating Oriented Bounding Box (OBB) primitives from single images. Note that Tulsiani et al. [40] uses mesh as supervision. Secondly, we compare our method with PQ-Net [49] in recovering Axis-Aligned Bounding Box (AABB) primitives, which outperforms other methods, such as StructureNet [30].

Specifically, we first compare our method in the semantic-agnostic setting (**Ours Agn**) with 3D-PRNN and Tulsiani, in recovering OBB. To show the benefits of semantic relations in both primitive reasoning and class-specific model setting, we also compare our semantic-aware model (**Ours Sem**) to the semantic-agnostic version. Moreover, we compare our semantic-aware model (**Ours Sem***) with PQ-Net, which is also trained in class-specific setting, in recovering AABB.

We evaluate our method using three metrics: *Hausdorff error* (HErr), *thresholded accuracy* (TAcc$^\delta$ [%]) and *intersection over union* (IoU [%]). To evaluate the quality of our estimated primitives, we adopt HErr and TAcc$^\delta$ ($\delta$ is the threshold) from Im2Struct [31],

**Table 1: Performance on Pix3D [37] dataset. Results for Tulsiani, 3D-PRNN, Ours Agn, and Ours Sem are for OBB estimation, while results for PQ-Net and Ours Sem\* are for AABB. $IoU_p$ indicates our comparison to the voxelized primitive ground truth, focusing on primitive quality instead of surface details.**

| Methods | HErr | $TAcc^{0.1}$ | $TAcc^{0.2}$ | $TAcc^{0.3}$ | $IoU_p$ |
|---|---|---|---|---|---|
| ***Chair*** | | | | | |
| Tulsiani | - | - | - | - | 19.2 |
| 3D-PRNN | 0.242 | 5.8 | 26.1 | 50.6 | 25.9 |
| Ours Agn | 0.207 | 6.7 | 31.2 | 55.6 | 30.0 |
| Ours Sem | **0.187** | **10.5** | **42.5** | **64.6** | **38.1** |
| PQ-Net | *0.215* | *6.8* | *27.3* | *48.2* | *36.6* |
| Ours Sem\* | *0.186* | *10.9* | *43.8* | *66.9* | *40.3* |
| ***Table*** | | | | | |
| Tulsiani | - | - | - | - | 8.4 |
| 3D-PRNN | 0.295 | 3.8 | 29.1 | 56.0 | 12.0 |
| Ours Agn | 0.245 | 8.2 | 41.9 | 64.5 | 19.6 |
| Ours Sem | **0.233** | **12.8** | **43.0** | **70.5** | **20.5** |
| PQ-Net | *0.247* | *8.3* | *38.1* | *62.2* | *19.0* |
| Ours Sem\* | *0.230* | *12.8* | *43.8* | *72.7* | *20.6* |
| ***Sofa*** | | | | | |
| Tulsiani | - | - | - | - | 55.4 |
| 3D-PRNN | 0.073 | 73.5 | 91.1 | 94.4 | 74.1 |
| Ours Agn | 0.061 | 76.7 | 95.9 | 97.3 | 76.6 |
| Ours Sem | **0.055** | **78.7** | **98.8** | **99.3** | **79.6** |
| PQ-Net | *0.065* | *67.7* | *96.6* | *98.2* | *73.2* |
| Ours Sem\* | *0.054* | *78.7* | *98.8* | *99.5* | *79.6* |

which are based on Hausdorff distance of primitive pairs. Since the primitives are well-aligned in voxel space, we also voxelize our primitive predictions and compare them with the ground truth voxels using IoU, to evaluate the quality of each object as a whole.

Given the test set of $T$ samples, we compute Hausdorff error $HErr = \frac{1}{2T} \sum_i^T (D(\mathcal{S}_i, \mathcal{S}_i^*) + D(\mathcal{S}_i^*, \mathcal{S}_i))$, where $\mathcal{S}_i$ is a predicted shape consisting of a set of primitives and $\mathcal{S}_i^*$ is its corresponding ground truth. For each primitive, we represent it as a set $\mathcal{V}$, consisting of its eight corners. $D(\mathcal{S}_1, \mathcal{S}_2) = \frac{1}{n} \sum_{\mathcal{V}_1^j \in \mathcal{S}_1} \min_{\mathcal{V}_2^k \in \mathcal{S}_2} H(\mathcal{V}_1^j, \mathcal{V}_2^k)$ is the averaged minimum Hausdorff distance from primitives in $\mathcal{S}_1$ to those in $\mathcal{S}_2$, where $n$ is the number of primitives in $\mathcal{S}_1$. $H(\mathcal{V}_1, \mathcal{V}_2) = \max_{\mathbf{q}_1 \in \mathcal{V}_1} \min_{\mathbf{q}_2 \in \mathcal{V}_2} ||\mathbf{q}_1 - \mathbf{q}_2||$ is the Hausdorff distance between two vertex set $\mathcal{V}_1$ and $\mathcal{V}_2$.

Thresholded accuracy $TAcc^\delta$ is the percentage of predicted primitives that $H(\mathcal{V}_i, \mathcal{V}_i^*)/L(\mathcal{V}_i^*) < \delta$, where $\mathcal{V}_i$ is a primitive in recovered shape, $\mathcal{V}_i^*$ is its nearest primitive in ground truth shape in Hausdorff distance, $L(\mathcal{V}_i^*)$ is the length of the diagonal of primitive $\mathcal{V}_i^*$, and $\delta$ is the threshold.

## 4.2 Implementation Details

We normalize all primitive parameters to have zero mean and standard deviation. Following 3D-PRNN, we assume all objects have bilateral symmetrical primitive pairs. For such symmetrical pairs, we only model the left half in our network and generate the right half by mirroring w.r.t. to the symmetry plane, to improve learning efficiency. For each object, we pre-sort all primitives based on the height of their centers, and use the proposal network to generate the primitive sequence in both bottom-up and top-down orderings separately. The dimension of LSTM hidden state $D_h$ and graph node $D_z$ is 800 and 1024, respectively. The weights to primitive regression loss and box regression loss $\lambda_{rp}$ and $\lambda_{rb}$ are both 10.

We re-implement 3D-PRNN and apply it to Pix3D and ModelNet datasets following their released codebase. For experiments of 3D-PRNN on Pix3D [37], we replace their original depth encoder with the same ResNet18 backbone as in our method, to encode RGB image features. To obtain AABB ground truth and predictions for **Ours Sem\***, we directly transform OBB ground truth and predictions from **Ours Sem**, by simply taking the tightest axis-aligned bounding box of each OBB. We train PQ-Net on AABB supervisions with their released code.

Please find more implementation details in supplementary.[3]

## 4.3 Results

We present comparison results with previous methods in two different settings: recovering Oriented Bounding Box (OBB) and Axis-Aligned Bounding Box (AABB). Below we introduce the detailed results on three benchmarks.

*4.3.1 Pix3D.* As shown in Table 1, for OBB estimation, **Ours Agn** outperforms 3D-PRNN in all metrics on all object categories. Our method achieves larger improvement on *chair* and *table*, whose shapes have more flexibility and are harder to recover compared to *sofa*. This demonstrates that our model can capture stronger shape regularity and improve recovered shape quality by utilizing local features and joint primitive reasoning. Additionally, **Ours Sem** achieves better performance on all categories than **Ours Agn**, especially in challenging category *chair* with a large margin (8.1% higher in $IoU_p$ and 0.020 lower in HErr). This is because *chair* in Pix3D has much more geometric and semantic part variations than other categories, which shows that our model can utilize rich semantic relations to help handle large geometric variations. For AABB estimation, **Ours Sem\*** also achieves better results than PQ-Net on all three categories.

*4.3.2 ModelNet.* As shown in Table 2, for OBB estimation, **Ours Agn** achieves better performance than 3D-PRNN on all categories in ModelNet. The performance gap is especially large on *night stand*, which has the least data for training. This shows that our model encodes stronger shape regularity prior with limited training data and thus can handle the problem of data imbalance more robustly. For *table* class, our $IoU_v$ is close to 3D-PRNN, but we achieves higher $IoU_p$ by 4.4% and lower HErr by 0.025, which is more related to the quality of our primitive prediction instead of detailed geometric surface. Moreover, **Ours Sem** also outperforms **Ours Agn** on all

---

[3]Our code is available at https://github.com/hailieqh/3D-Object-Primitive-Graph.

categories, but the gaps are not as large as those on Pix3D. This is because the primitives in ModelNet have larger distribution variation than Pix3D, hence they cannot be easily captured in our designed semantic space with a few labels. For AABB estimation, **Ours Sem\*** significantly outperforms PQ-Net in Hausdorff metrics HErr and $TAcc^\delta$. The performance gaps are smaller in $IoU_{p/v}$ for *night stand*, which is mainly because erroneous primitives have less impact on voxel, especially for object categories with large volume and regular shapes. The results on ModelNet verifies the capability of our method in modeling complex primitive dependency, as objects in ModelNet have much larger variations of primitive combinations than objects in Pix3D. (See also Fig. 3 for visual examples.)

*4.3.3 NYU Depth V2.* We train our model on ModelNet synthetic data and directly applied to NYU Depth V2 real-world depth images. Despite the large domain gap between training and testing, we outperform 3D-PRNN on all categories for OBB estimation, as shown in Table 3. We notice that the performance gaps on *chair* and *table* are relatively small. This is mainly because the sequential model cannot provide high quality proposals when heavy occlusion occurs, which further leads to suboptimal performance of our subsequent primitive reasoning module. Note that our method still achieves better performance on *night stand* with a considerable margin, which shows the robustness of our model in handling categories with limited training data. Furthermore, **Ours Sem** outperforms **Ours Agn** with a considerable margin, demonstrating the benefits of semantic-aware models for extremely challenging cases in real-world scenarios. For AABB estimation, **Ours Sem\*** also outperforms PQ-Net on *chair* and *table*.

*4.3.4 Qualitative Results.* We also visually compare the results of Tulsiani, PQ-Net, 3D-PRNN, **Ours Agn** and **Ours Sem**, as shown in Fig. 3. For Pix3D real-world images, our models can more robustly handle occlusion, truncation, unusual viewpoints and challenging novel instances, compared to Tulsiani, PQ-Net and 3D-PRNN. For ModelNet synthetic data, our recovered shapes are also better in terms of both detailed shape consistency with images and global shape regularity. Note that on categories with limited training data such as *night stand*, our models are able to correctly recognize object categories and to recover shapes matched to input images, while 3D-PRNN tends to generate an instance from the dominant category (*chair*). This shows that our primitive reasoning can learn to capture stronger shape structure prior during training and thus can handle the problem of training data imbalance more robustly. For more visualization, please see supplementary.

## 4.4 Ablation Study

In this section, we conduct several experiments to show the efficacy of the components of our full model **Ours Sem** on Pix3D *chair* class, which has large instance variations in complex real-world environment, as shown in Table 4.

**LSTM**: Our bidirectional sequential proposal network captures strong primitive dependency within an object and generates better primitive proposals. With only one **LSTM** unit, our model performance drops by 0.003 higher in HErr and 0.8% lower in $IoU_p$. Moreover, we also explored using a Faster R-CNN backbone instead of **LSTM**. Specifically, we first train a Faster R-CNN on 2D

**Table 2: Performance on ModelNet [51]. In addition to** $IoU_p$**, we also evaluate primitives in** $IoU_v$**, where** $v$ **denotes we compare our voxelized primitive predictions to original ground truth voxel. We show both the originally reported results and our re-implementation of 3D-PRNN before and after '/'.**

| Methods | HErr | $TAcc^{0.1}$ | $TAcc^{0.2}$ | $TAcc^{0.3}$ | $IoU_p$ | $IoU_v$ |
|---|---|---|---|---|---|---|
| *Chair* | | | | | | |
| Tulsiani | - | - | - | - | 30.4 | 26.8 |
| [51] | - | - | - | - | - | 25.3 |
| 3D-PRNN | -/0.176 | -/5.0 | -/26.2 | -/55.1 | -/40.1 | 23.8/28.6 |
| Ours Agn | 0.160 | 6.6 | 29.5 | 54.6 | 44.4 | 32.2 |
| Ours Sem | **0.157** | **8.7** | **36.2** | **59.9** | **46.5** | **33.6** |
| PQ-Net | *0.210* | *1.4* | *11.7* | *35.4* | *34.7* | *22.0* |
| Ours Sem\* | *0.156* | *11.0* | *39.2* | *64.2* | *47.9* | *31.2* |
| *Table* | | | | | | |
| Tulsiani | - | - | - | - | 30.3 | 22.4 |
| [51] | - | - | - | - | - | 25.0 |
| 3D-PRNN | -/0.190 | -/12.4 | -/38.0 | -/66.6 | -/34.0 | 26.3/28.9 |
| Ours Agn | 0.165 | 16.8 | 37.7 | 59.6 | 38.4 | 29.7 |
| Ours Sem | **0.160** | **20.1** | **42.9** | **64.7** | **41.9** | **33.3** |
| PQ-Net | *0.227* | *8.6* | *20.7* | *37.5* | *27.0* | *19.7* |
| Ours Sem\* | *0.160* | *20.7* | *45.2* | *68.2* | *41.5* | *33.0* |
| *Night stand* | | | | | | |
| Tulsiani | - | - | - | - | 42.5 | **42.9** |
| [51] | - | - | - | - | - | 29.5 |
| 3D-PRNN | -/0.279 | -/4.2 | -/24.5 | -/48.7 | -/32.8 | 26.6/31.9 |
| Ours Agn | 0.244 | 7.6 | 33.2 | 56.4 | 41.2 | 39.9 |
| Ours Sem | **0.235** | **13.1** | **42.3** | **63.9** | **43.1** | 41.5 |
| PQ-Net | *0.285* | *1.9* | *16.8* | *41.0* | *39.6* | *38.5* |
| Ours Sem\* | *0.234* | *13.1* | *43.0* | *64.8* | *43.4* | *41.9* |

**Table 3: Performance on NYU Depth V2 [36]. We evaluate in** $IoU_v$ **as there are no primitive annotations and we compare our prediction to ground truth voxel.**

| Category | 3D-PRNN | Ours Agn | Ours Sem | PQ-Net | Ours Sem\* |
|---|---|---|---|---|---|
| Chair | 13.8 | 14.6 | **17.2** | *13.1* | *16.9* |
| Table | 5.2 | 6.0 | **8.1** | *6.8* | *8.2* |
| Night stand | 8.6 | 18.7 | **28.0** | *31.7* | *28.3* |

bounding boxes to generate 2D proposals and then extract local conv features from 2D proposals to regress corresponding 3D primitive proposal parameters. In this case, model performance drops significantly by 0.016 higher in HErr and 4.9% lower in $IoU_p$.

**Sem**: Semantic class supervision and semantic-aware reasoning helps reduce primitive output search space and provide more guidance for object structure recovery. Without **Sem**, we observe a performance drop, which is 0.005 higher in HErr and 1.4% lower in $IoU_p$ than **Ours Sem**.

**BConv**: Local conv features extracted from box regions in an image can help handle occlusion and improve the robustness of our model. Without **BConv**, the performance of our model drops by 0.008 higher in HErr and 2.3% lower in $IoU_p$.
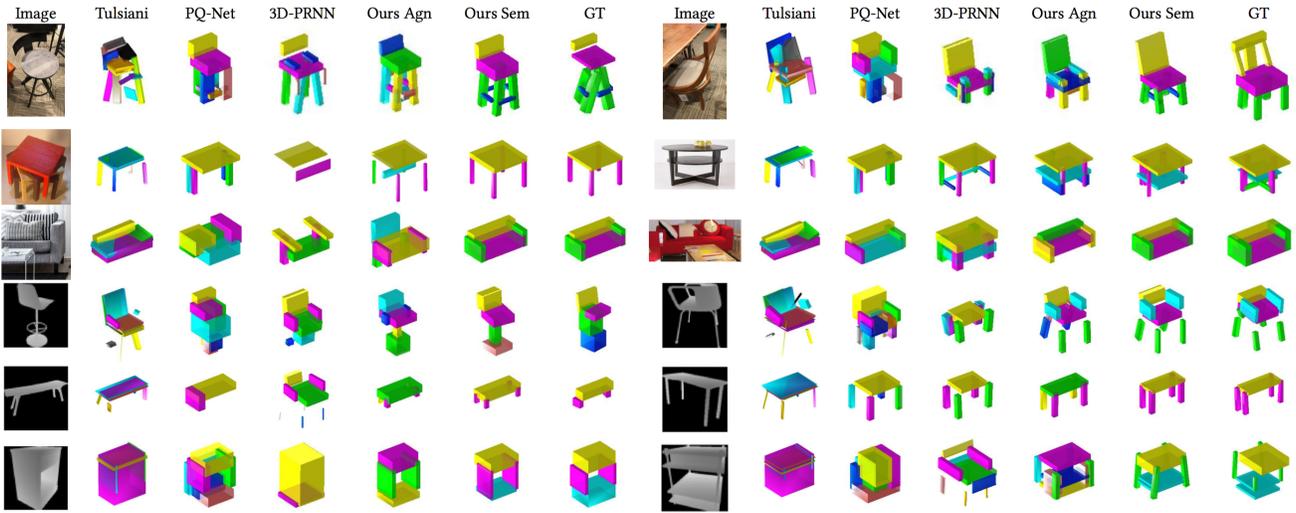
**Figure 3: Qualitative results on Pix3D [37] and ModelNet [51]. The color codings of Tulsiani, PQ-Net, 3D-PRNN, and Ours Agn are based on the order of drawing primitives and have no semantic meaning. The color codings of Ours Sem and GT are based on the semantics of primitives. Row 1-3 (Pix3D RGB):** *Chair, Table, Sofa*; **Row 4-6 (ModelNet Depth):** *Chair, Table, Night stand.*

**Table 4: Ablation study on Pix3D [37] chair. 1 or 2 under LSTM denotes the number of LSTM sequence models used to generate primitive proposals. We conduct the experiments in a drop-one-out manner.** $\text{TAcc}^{\delta}$ **and** $\text{IoU}_p$ **are in %.**

| Methods | LSTM | Module | | Graph | HErr ↓ | Metric | | | |
| | | Sem | BConv | | | $\text{TAcc}^{0.1}$ ↑ | $\text{TAcc}^{0.2}$ ↑ | $\text{TAcc}^{0.3}$ ↑ | $\text{IoU}_p$ ↑ |
|---|---|---|---|---|---|---|---|---|---|
| Ours Sem | 2 | ✓ | ✓ | Dense | **0.187** | **10.5** | **42.5** | **64.6** | **38.1** |
| | 2 | ✓ | ✓ | Star-like | 0.195 | 10.2 | 40.5 | 62.6 | 36.5 |
| | 2 | ✓ | ✓ | K-nearest | 0.197 | 10.5 | 39.6 | 61.7 | 35.8 |
| | 2 | ✓ | ✓ | - | 0.199 | 9.4 | 35.5 | 56.0 | 36.6 |
| | 2 | ✓ | - | Dense | 0.195 | 10.5 | 39.2 | 60.9 | 35.8 |
| | 2 | - | ✓ | Dense | 0.192 | 9.1 | 38.6 | 60.9 | 36.7 |
| | 1 | ✓ | ✓ | Dense | 0.190 | 9.7 | 38.0 | 63.1 | 37.3 |
| | Faster R-CNN | ✓ | ✓ | Dense | 0.203 | 9.2 | 35.5 | 58.7 | 33.2 |
| Baseline | 2 | - | - | - | 0.211 | 6.4 | 28.7 | 48.4 | 34.0 |

**Graph**: Our graph network helps capture the global shape context and mitigate errors in proposal generation. Without **Graph**, we fully train the primitive proposal network and apply NMS on final predictions to remove duplicate primitives. In this case, model performance deteriorates by 0.012 higher in HErr and 1.5% lower in $\text{IoU}_p$. Note that the performance drops more significantly in HErr and $\text{TAcc}^{\delta}$ than in $\text{IoU}_p$. This is mainly because the fine-level primitive quality is much lower without joint reasoning and global refinement, while the coarse-level shapes after voxelization are less affected by erratic primitives. To investigate the impact of graph structure, we also conduct experiments on two variants of graph design: (1) a distance graph in which each node is connected to K-nearest (K=3) nodes based on primitive center distance; (2) a star-like graph in which every node is connected to the central cushion nodes. Our densely-connected graph outperforms those two sparse graphs in all the metrics.

## 5 CONCLUSION

In this paper, we have developed a primitive-based graph neural network for 3D object estimation from single images. Our method first uses a proposal network to lift the image features into a pool of primitive proposals, and then builds a fully-connected graph network on these proposals to refine primitive features. Besides, we optionally take into account semantic part cues to provide better guidance for object shape. As a result, our method is able to conduct joint reasoning on the primitives for coherent 3D recovery. Moreover, we adopt a stage-wise strategy in our model learning, in which it first learns a recurrent network for primitive generation and then trains a graph network to propagate messages between the primitives. Evaluations on three benchmarks show that our approach consistently outperforms prior approaches with a sizable margin. Particularly, our method can robustly handle (self-)occlusion and challenging viewpoints in real-world scenarios.

# REFERENCES

[1] Irving Biederman. 1987. Recognition-by-components: a theory of human image understanding. *Psychological review* 94, 2 (1987), 115.

[2] Zhiqin Chen, Andrea Tagliasacchi, and Hao Zhang. 2020. Bsp-net: Generating compact meshes via binary space partitioning. In *CVPR*.

[3] Zhiqin Chen and Hao Zhang. 2019. Learning implicit fields for generative shape modeling. In *CVPR*.

[4] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 2016. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *ECCV*.

[5] Boyang Deng, Kyle Genova, Soroosh Yazdani, Sofien Bouaziz, Geoffrey Hinton, and Andrea Tagliasacchi. 2020. Cvxnet: Learnable convex decomposition. In *CVPR*.

[6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*.

[7] Haoqiang Fan, Hao Su, and Leonidas J Guibas. 2017. A point set generation network for 3d object reconstruction from a single image. In *CVPR*.

[8] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. 2020. Local Deep Implicit Functions for 3D Shape. In *CVPR*.

[9] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. 2019. Learning shape templates with structured implicit functions. In *ICCV*.

[10] Rohit Girdhar, David F Fouhey, Mikel Rodriguez, and Abhinav Gupta. 2016. Learning a predictable and generative vector representation for objects. In *ECCV*.

[11] Ruiqi Guo and Derek Hoiem. 2013. Support surface prediction in indoor scenes. In *ICCV*.

[12] Zhizhong Han, Chao Chen, Yu-Shen Liu, and Matthias Zwicker. 2020. ShapeCaptioner: Generative caption network for 3D shapes by learning a mapping from parts detected in multiple views to sentences. In *Proceedings of the 28th ACM International Conference on Multimedia*. 1018–1027.

[13] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *ICCV*.

[14] Qian He, Desen Zhou, Xuming He, et al. 2018. 3D Object Structure Recovery via Semi-supervised Learning on Videos. In *BMVC*.

[15] Lin Huang, Jianchao Tan, Jingjing Meng, Ji Liu, and Junsong Yuan. 2020. HOT-Net: Non-Autoregressive Transformer for 3D Hand-Object Pose Estimation. In *Proceedings of the 28th ACM International Conference on Multimedia*. 3136–3145.

[16] Qixing Huang, Hai Wang, and Vladlen Koltun. 2015. Single-view reconstruction via joint analysis of image and shape collections. *ACM TOG* (2015).

[17] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. 2018. Learning category-specific mesh reconstruction from image collections. In *ECCV*.

[18] Abhishek Kar, Shubham Tulsiani, Joao Carreira, and Jitendra Malik. 2015. Category-specific object reconstruction from a single image. In *CVPR*.

[19] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. 2018. Neural 3d mesh renderer. In *CVPR*.

[20] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[21] Chen Kong, Chen-Hsuan Lin, and Simon Lucey. 2017. Using locally corresponding CAD models for dense 3D reconstructions from a single image. In *CVPR*.

[22] Chi Li, M Zeeshan Zia, Quoc-Huy Tran, Xiang Yu, Gregory D Hager, and Manmohan Chandraker. 2017. Deep supervision with shape concepts for occlusion-aware 3d object parsing. In *CVPR*.

[23] Jun Li, Kai Xu, Siddhartha Chaudhuri, Ersin Yumer, Hao Zhang, and Leonidas Guibas. 2017. Grass: Generative recursive autoencoders for shape structures. *ACM TOG* (2017).

[24] Yangyan Li, Hao Su, Charles Ruizhongtai Qi, Noa Fish, Daniel Cohen-Or, and Leonidas J Guibas. 2015. Joint embeddings of shapes and images via cnn image purification. *ACM TOG* (2015).

[25] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *CVPR*.

[26] Yawen Lu, Yuxing Wang, and Guoyu Lu. 2020. Single Image Shape-from-Silhouettes. In *Proceedings of the 28th ACM International Conference on Multimedia*. 3604–3613.

[27] Quan Meng, Jiakai Zhang, Qiang Hu, Xuming He, and Jingyi Yu. 2020. LGNN: A Context-aware Line Segment Detector. In *Proceedings of the 28th ACM International Conference on Multimedia*. 4364–4372.

[28] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. 2019. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*.

[29] Niloy J Mitra, Michael Wand, Hao Zhang, Daniel Cohen-Or, Vladimir Kim, and Qi-Xing Huang. 2014. Structure-aware shape processing. In *ACM SIGGRAPH 2014 Courses*.

[30] Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas Guibas. 2019. StructureNet: Hierarchical Graph Networks for 3D Shape Generation. *ACM TOG, Siggraph Asia 2019* (2019).

[31] Chengjie Niu, Jun Li, and Kai Xu. 2018. Im2struct: Recovering 3d shape structure from a single rgb image. In *CVPR*.

[32] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. 2019. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*.

[33] Despoina Paschalidou, Luc Van Gool, and Andreas Geiger. 2020. Learning Unsupervised Hierarchical Part Decomposition of 3D Objects from a Single RGB Image. In *CVPR*.

[34] Despoina Paschalidou, Ali Osman Ulusoy, and Andreas Geiger. 2019. Superquadrics Revisited: Learning 3D Shape Parsing beyond Cuboids. In *CVPR*.

[35] Gopal Sharma, Rishabh Goyal, Difan Liu, Evangelos Kalogerakis, and Subhransu Maji. 2018. CSGNet: Neural Shape Parser for Constructive Solid Geometry. In *CVPR*.

[36] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. 2012. Indoor segmentation and support inference from rgbd images. In *ECCV*.

[37] Xingyuan Sun, Jiajun Wu, Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Tianfan Xue, Joshua B Tenenbaum, and William T Freeman. 2018. Pix3d: Dataset and methods for single-image 3d shape modeling. In *CVPR*.

[38] Supasorn Suwajanakorn, Noah Snavely, Jonathan J Tompson, and Mohammad Norouzi. 2018. Discovery of latent 3d keypoints via end-to-end geometric reasoning. *NeurIPS* (2018).

[39] Yonglong Tian, Andrew Luo, Xingyuan Sun, Kevin Ellis, William T. Freeman, Joshua B. Tenenbaum, and Jiajun Wu. 2019. Learning to Infer and Execute 3D Shape Programs. In *ICLR*.

[40] Shubham Tulsiani, Hao Su, Leonidas J Guibas, Alexei A Efros, and Jitendra Malik. 2017. Learning shape abstractions by assembling volumetric primitives. In *CVPR*.

[41] Meng Wang, Lingjing Wang, and Yi Fang. 2017. 3densinet: A robust neural network architecture towards 3d volumetric object prediction from 2d image. In *Proceedings of the 25th ACM international conference on Multimedia*. 961–969.

[42] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. 2018. Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images. In *ECCV*.

[43] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. 2018. Non-local neural networks. In *CVPR*.

[44] Yanzhen Wang, Kai Xu, Jun Li, Hao Zhang, Ariel Shamir, Ligang Liu, Zhiquan Cheng, and Yueshan Xiong. 2011. Symmetry hierarchy of man-made objects. In *CGF*.

[45] Jiajun Wu, Yifan Wang, Tianfan Xue, Xingyuan Sun, Bill Freeman, and Josh Tenenbaum. 2017. Marrnet: 3d shape reconstruction via 2.5 d sketches. In *NeurIPS*.

[46] Jiajun Wu, Tianfan Xue, Joseph J Lim, Yuandong Tian, Joshua B Tenenbaum, Antonio Torralba, and William T Freeman. 2016. Single image 3d interpreter network. In *ECCV*.

[47] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. 2016. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *NeurIPS*.

[48] Jiajun Wu, Chengkai Zhang, Xiuming Zhang, Zhoutong Zhang, William T Freeman, and Joshua B Tenenbaum. 2018. Learning shape priors for single-view 3d completion and reconstruction. In *ECCV*.

[49] Rundi Wu, Yixin Zhuang, Kai Xu, Hao Zhang, and Baoquan Chen. 2020. PQ-NET: A generative part seq2seq network for 3D shapes. In *CVPR*.

[50] Zhenyu Wu, Duc Hoang, Shih-Yao Lin, Yusheng Xie, Liangjian Chen, Yen-Yu Lin, Zhangyang Wang, and Wei Fan. 2020. Mm-hand: 3d-aware multi-modal guided hand generative network for 3d hand pose synthesis. *arXiv preprint arXiv:2010.01158* (2020).

[51] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 2015. 3d shapenets: A deep representation for volumetric shapes. In *CVPR*.

[52] Hongwen Zhang, Jie Cao, Guo Lu, Wanli Ouyang, and Zhenan Sun. 2019. Danet: Decompose-and-aggregate network for 3d human shape and pose estimation. In *Proceedings of the 27th ACM International Conference on Multimedia*. 935–944.

[53] Yumeng Zhang, Li Chen, Yufeng Liu, Wen Zheng, and Junhai Yong. 2020. Adaptive Wasserstein Hourglass for Weakly Supervised RGB 3D Hand Pose Estimation. In *Proceedings of the 28th ACM International Conference on Multimedia*. 2076–2084.

[54] Youyi Zheng, Daniel Cohen-Or, Melinos Averkiou, and Niloy J Mitra. 2014. Recurring part arrangements in shape collections. In *CGF*.

[55] Fan Zhu, Li Liu, Jin Xie, Fumin Shen, Ling Shao, and Yi Fang. 2018. Learning to synthesize 3d indoor scenes from monocular images. In *Proceedings of the 26th ACM international conference on Multimedia*. 501–509.

[56] Chuhang Zou, Ersin Yumer, Jimei Yang, Duygu Ceylan, and Derek Hoiem. 2017. 3d-prnn: Generating shape primitives with recurrent neural networks. In *ICCV*.
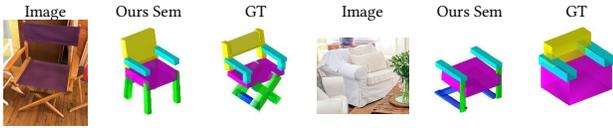
**Figure 4: Two failure cases of our method on Pix3D chair. The chair on the left has bilaterally crossed legs which are unique in the dataset. The chair on the right has few texture and is occluded by objects with very similar appearance.**

**Table 5: Performance on different subsets of Pix3D [37] chair. Set A, B, C, and D are truncated, occluded, slightly occluded, and complete (neither truncated nor occluded), respectively.**

| Methods | HErr ↓ | | | | $IoU_p$ ↑ | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Set A | Set B | Set C | Set D | Set A | Set B | Set C | Set D |
| 3D-PRNN | 0.261 | 0.249 | 0.311 | 0.233 | 21.0 | 23.3 | 17.0 | 27.8 |
| Ours Sem | **0.176** | **0.186** | **0.249** | **0.184** | **40.8** | **38.1** | **30.7** | **38.2** |
| PQ-Net | *0.221* | *0.220* | *0.265* | *0.210* | *35.4* | *36.5* | *31.7* | *37.3* |
| Ours Sem* | *0.175* | *0.185* | *0.247* | *0.183* | *42.1* | *40.1* | *32.6* | *40.6* |

## A  IMPLEMENTATION DETAILS

We choose ResNet18 as backbone for its strong performance and extensive usage in the prior work [28, 33]. Empirically, we also found that deeper ResNets provide little improvement. To handle different inputs, we adopt the same network architecture, ResNet18+FPN, as our backbone, which takes three-channel images as input. As a result, we process the RGB and depth inputs in slightly different manners. For experiments on Pix3D, we directly send the input RGB images into the network, while for the ModelNet and NYUv2 dataset, we duplicate each depth image three times to build a 3-channel input for the network. Our ResNet18 module is pretrained on ImageNet [6].

We first train the proposal network for 20 epochs and then freeze the parameters during the training of primitive reasoning, to prevent it from overfitting on training data and providing no training signals for the primitive reasoning network. We train our model with batch size 16 for a maximum of 400 epochs and with an Adam [20] optimizer, whose learning rate is 1e-4 and betas are (0.95, 0.999). We explore batch size from 4 to 64, epochs from 200 to 500, and learning rate from 1e-2 to 1e-5. Final decisions on hyper-parameters are based on the performance on validation split. Then

we fix all hyper-parameters to train on training and validation splits together to obtain our final models.

The implementation of our method is based on PyTorch 1.0.0. We run experiments on single TITAN Xp GPU cards, using less than 2GB GPU memory, in Ubuntu 16.04.4. It takes about four days to train a class-agnostic model on all three categories of Pix3D or ModelNet. The seeds to randomness are all fixed to 0.

## B  MORE QUALITATIVE RESULTS

We also visually compare the results of Tulsiani, PQ-Net, 3D-PRNN, **Ours Agn** and **Ours Sem**, as shown in Fig. 5 for Pix3D and Fig. 6 for ModelNet. For Pix3D real-world images, our models can more robustly handle occlusion, truncation, unusual viewpoints and challenging novel instances, compared to Tulsiani, PQ-Net and 3D-PRNN. For ModelNet synthetic data, our recovered shapes are also better than Tulsiani, PQ-Net and 3D-PRNN, in terms of both detailed shape consistency with images and global shape regularity.

## C  FURTHER ANALYSIS

### C.1  Ablation on Proposal Network

We also evaluate our proposal network design comparing to other variations in Table 6. Similar to $TAcc^\delta$, we also use *thresholded recall* ($TRec^\delta$ [%]) to evaluate the quality of proposals, which is the percentage of ground truth primitives that $H(\mathcal{V}_i, \mathcal{V}_i^*)/L(\mathcal{V}_i^*) < \delta$, where for each ground truth primitive $\mathcal{V}_i^*$, $\mathcal{V}_i$ is its nearest match in predicted primitives. Note that each predicted primitive can only match one ground truth primitive. Our proposal module achieves the best $TRec^\delta$ with all thresholds, which verifies the capability of our proposal network to capture rich semantic dependency and to attend to local image features.

### C.2  Robustness to Occlusion

In Table 5, we show results on different subsets of Pix3D chair according to truncation and occlusion. **Ours Sem** consistently outperforms 3D-PRNN and PQ-Net on all different subsets. In addition, our performance stays comparable in truncated and occluded subsets, which further demonstrates the robustness of our method.

### C.3  Failure Cases

Our results include two typical failure cases: 1) novel object instances with uncommon 3D shapes; 2) heavily occluded objects. Two examples are shown in Fig. 4. Both scenarios can potentially be improved with better modeling of object priors or utilizing larger datasets.
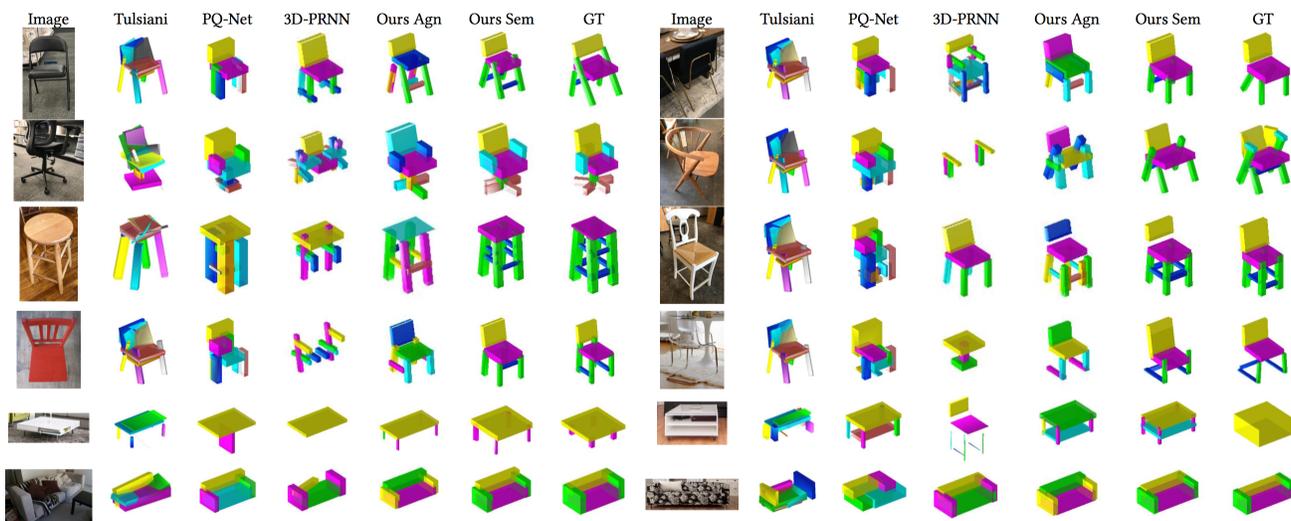
**Figure 5: Qualitative results on Pix3D [37]. Note that the color codings of Tulsiani, PQ-Net, 3D-PRNN, and Ours Agn are based on the order of drawing primitives and have no semantic meaning. In contrast, the color codings of Ours Sem and GT are based on the semantics of primitives.**
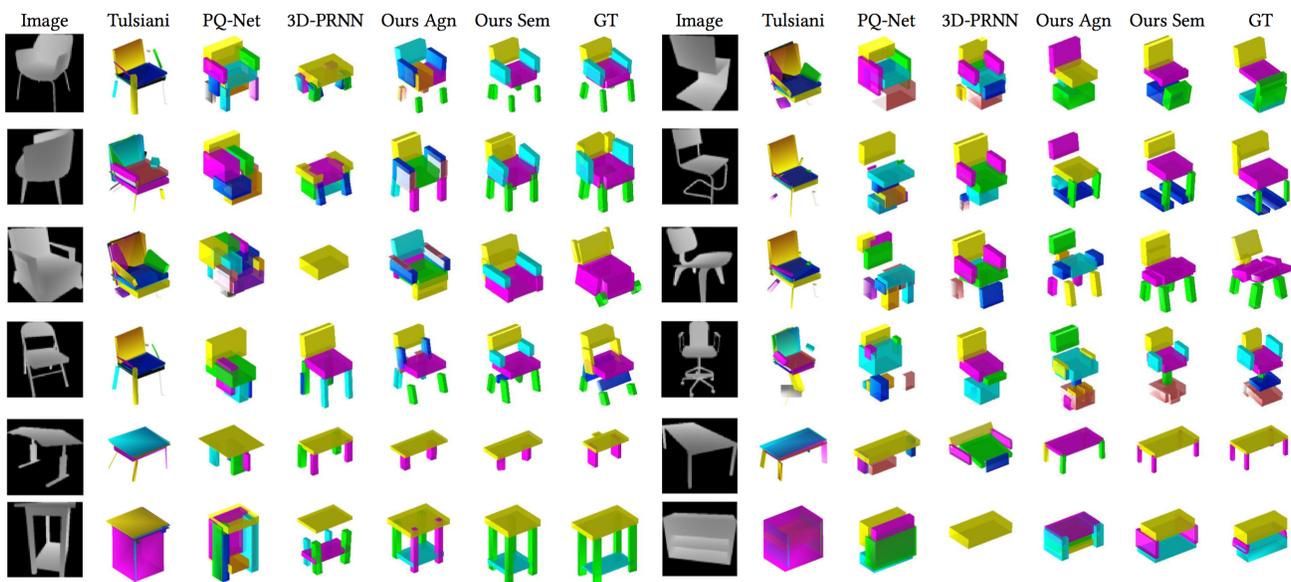


**Figure 6: Qualitative results on ModelNet [51]. Note that the color codings of Tulsiani, PQ-Net, 3D-PRNN, and Ours Agn are based on the order of drawing primitives and have no semantic meaning. In contrast, the color codings of Ours Sem and GT are based on the semantics of primitives.**

Table 6: Ablation of proposal networks on Pix3D [37] chair. 1 or 2 under LSTM denotes the number of LSTM units the model uses to capture primitive sequence in different order. We conduct the experiments in drop-one-out manner.

| Proposal Network | Module | | | Thresholded Recall ($\text{TRec}^\delta$) ↑ | | | | | |
| | LSTM | Sem | BConv | $\text{TRec}^{0.1}$ | $\text{TRec}^{0.2}$ | $\text{TRec}^{0.3}$ | $\text{TRec}^{0.4}$ | $\text{TRec}^{0.5}$ | $\text{TRec}^{0.6}$ |
|---|---|---|---|---|---|---|---|---|---|
| Ours Sem | 2 | ✓ | ✓ | **14.2** | **43.2** | **63.5** | **77.1** | **87.1** | **92.1** |
| | 2 | ✓ | - | 13.8 | 41.3 | 61.4 | 74.9 | 86.7 | 92.0 |
| | 2 | - | ✓ | 7.9 | 30.4 | 48.9 | 67.2 | 82.1 | 90.4 |
| | 1 | ✓ | ✓ | 9.5 | 33.2 | 50.6 | 65.5 | 77.3 | 82.7 |
| | Faseter R-CNN | ✓ | ✓ | 8.7 | 30.8 | 47.3 | 58.3 | 69.4 | 77.6 |