# Towards Part-Based Understanding of RGB-D Scans

Alexey Bokhovkin[1,2]   Vladislav Ishimtsev[2]   Emil Bogomolov[2]   Denis Zorin[3,2]

Alexey Artemov[2]   Evgeny Burnaev[2]   Angela Dai[1]

[1]Technical University of Munich
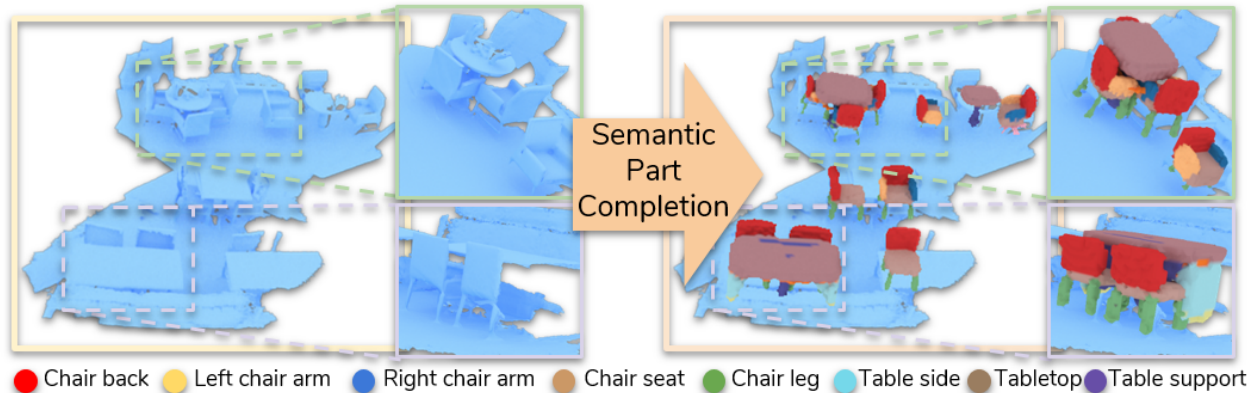[2]Skolkovo Institute of Science and Technology
[3]New York University

Figure 1: From an input RGB-D scan (left), we propose to detect objects in the scan and predict their complete part decompositions as *semantic part completion*; that is, we predict the part masks for the complete object, inferring the part geometry of any missing or unobserved regions in the scan. To achieve this, we predict the part structure of each detected object to drive a geometric prior-driven prediction of the complete part masks.

## Abstract

*Recent advances in 3D semantic scene understanding have shown impressive progress in 3D instance segmentation, enabling object-level reasoning about 3D scenes; however, a finer-grained understanding is required to enable interactions with objects and their functional understanding. Thus, we propose the task of part-based scene understanding of real-world 3D environments: from an RGB-D scan of a scene, we detect objects, and for each object predict its decomposition into geometric part masks, which composed together form the complete geometry of the observed object. We leverage an intermediary part graph representation to enable robust completion as well as building of part priors, which we use to construct the final part mask predictions. Our experiments demonstrate that guiding part understanding through part graph to part prior-based predictions significantly outperforms alternative approaches to the task of semantic part completion.*

## 1. Introduction

Recently, we have seen remarkable advances in 3D semantic scene understanding, driven by efforts in large-scale data collection and annotation of 3D reconstructions of RGB-D scanned environments [5, 2], coupled with exploration of 3D deep learning approaches across 3D representations such as sparse or dense volumetric grids [55, 39, 5, 15, 4], point clouds [38, 40], meshes [13, 23], and multi-view [7, 50]. This has led to significant progress in both 3D semantic segmentation as well as 3D semantic instance segmentation [16, 15, 4, 25]. These have enabled a basis for 3D perception at the level of objects, which is essential for semantic understanding, but lacks finer-grained understanding often critical for enabling interactions with objects and reasoning about functionality (e.g., the seat part of a chair is for sitting on, a knob or handle enables opening doors or drawers).

At the same time, notable progress has been made in part segmentation for shapes [33, 32, 18]. However, these methods have been developed on synthetic datasets such as ShapeNet [3], of objects in isolation; this scenario is

much less complex than the objects observed in real-world environments. Thus, we aim to bring these two directions together and propose the task of *semantic part completion*, predicting the part decomposition of objects in real-world 3D environments, where observations are often cluttered and geometrically incomplete (e.g., due to occlusions, sensor limitations, etc). That is, from an RGB-D scan of a scene, we detect objects characterized by 3D bounding boxes and class labels, and for each object, we predict its complete part decomposition into binary part masks, with each part mask reflecting the part geometry of the complete object, including unobserved missing regions, to achieve a holistic understanding of the objects in an observed scene.

To achieve this part-based understanding of a scene, we propose to predict the full part graph for each detected object, and based on the predicted part graph, the geometric masks for each complete part. Predicting the part graph structure enables capturing the complete semantic structure of the object in a low-dimensional representation, allowing reliable prediction of missing and unobserved parts (e.g., for a four-legged table with one leg unobserved, the missing leg is easy to predict based on commonly observed table part patterns). Furthermore, this enables us to build and exploit strong part geometry priors for each predicted part in the part graph. We can then predict the part masks by finding similar part priors and refining them to produce final part mask predictions. This enables a robust decomposition of an RGB-D scan of a scene into its component objects and their constituent parts, including regions of objects that have been unobserved. We believe that this takes an important step towards enabling local interactions with objects and functionality analysis in real-world 3D scenes.

We formulate the task of semantic part completion for 3D scene understanding, informing comprehensive part-based object understanding of real-world scans. To address this part understanding, we propose an approach to decompose a 3D scan of a scene into its complete object parts, outperforming state-of-the-art alternative approaches for the task:

- We propose to predict part graph information for objects in real-world scan scenes as an intermediary representation that enables robust, part-based completion of objects.

- We leverage the predicted part graphs to guide prior-based prediction for effective inference of geometric part mask decomposition for the objects of a scanned scene.

## 2. Related Work

**3D Object Detection and Instance Segmentation.** Following the success of convolutional neural networks for object detection and instance segmentation in 2D images [12, 42, 41, 19], we are now seeing notable advances in 3D object localization and segmentation. Earlier approaches leveraging 3D convolutional neural networks developed methods operating on dense voxel grids using 3D region proposal techniques for detection and segmentation [47, 20]. Sparse volumetric backbones have also been leveraged to enable effective feature extraction on high-resolution inputs for improved 3D detection and segmentation performance [10, 16]. Recently, VoteNet [37] introduced a Hough Voting-inspired scheme for 3D object detection on point clouds. This was extended by MLCVNet [56] to incorporate multi-scale contextual information for improved detection performance. These approaches have now shown impressive performance for instance-level scene understanding; we aim to build upon this and propose to infer finer-grained part decomposition for each object in a 3D scan.

**3D Scan Completion.** Repairing and completing holes or broken meshes has been well-studied for 3D shapes. Traditional methods have mainly focused on repairing small holes by fitting geometric primitives, continuous energy minimization, or leveraging surface reconstruction for interpolation of missing regions [34, 59, 49, 27, 28]. Structural or symmetry priors have also been leveraged for shape completion [52, 31, 36, 46, 49]. Recently, generative deep learning approaches have been developed, with significant progress in 3D shape reconstruction and completion [55, 9, 17, 35].

In addition to operating on the limited spatial context of shapes, generative deep learning approaches have also been developed for completion of 3D scenes. Song et al. [48] developed a voxel-based approach to predict geometric occupancy of a single depth frame, leveraging a large-scale synthetic 3D dataset of scenes. Dai et al. [8] proposed an autoregressive approach for scan completion, enabling very large scale completion. SG-NN [6] presented a self-supervised approach towards 3D scan completion, enabling training only on real scan data. These approaches operate on geometric completion but without knowledge of individual object instances, which is fundamental to many perception-based tasks. RevealNet [21] introduced an approach to detect objects in a 3D scan and infer each object's complete geometry, joining together geometric reconstruction with object-based understanding. We similarly aim to infer each object's complete geometry from a partial scan observation, but infer a part decomposition of the object structure, enabling both finer-grained understanding as well as more effective object completion through its part structure.

**Part Segmentation of 3D Shapes.** Understanding the structure of a 3D shape by identifying shape parts has been long-studied in shape analysis. Various approaches have been developed for finding a consistent segmentation across a set of shapes without supervision of part labels [14, 24, 45, 22]. Recently, deep learning based approaches
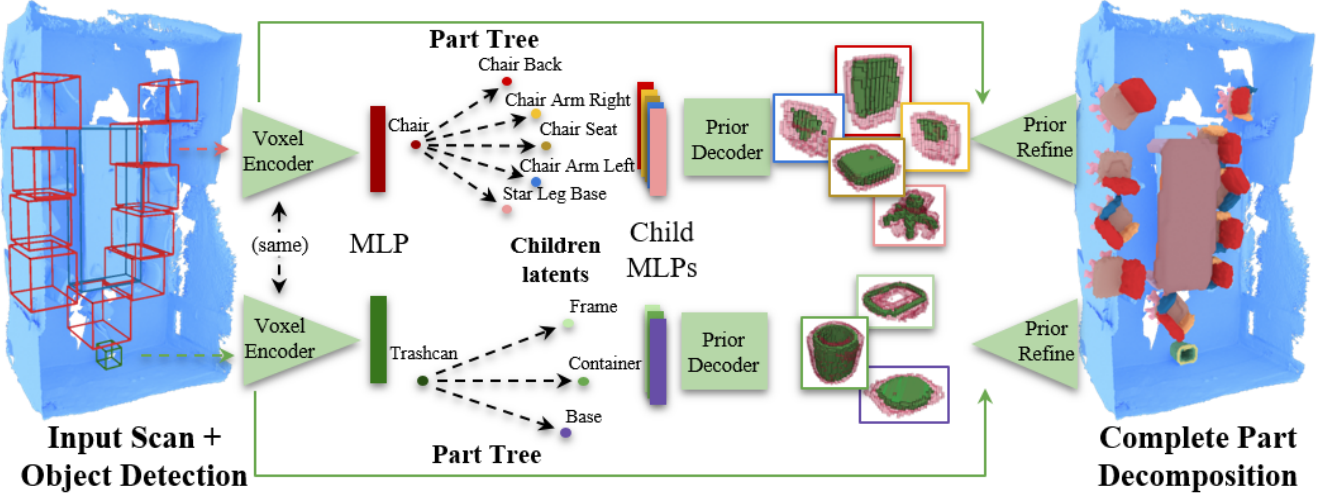
Figure 2: Overview of our approach. From an input scan, we detect objects as their 3D bounding boxes, and for each object (a chair and a trash can visualized top and bottom, respectively), we predict their part tree structure, which is then used to guide a geometric prior-based part mask prediction. This results in a part decomposition of the scene where each object is decomposed into its complete part geometry, including any missing or unobserved regions.

have been developed to find part segmentation of shapes in a data-driven fashion [26, 58, 18]. To better capture more complex structures in the part layout of shapes, several methods propose to parse object parts as hierarchies [54, 53, 57, 33, 32]. Such hierarchically structured representations have also been adopted for 3D scene synthesis, leveraging a scene graph [11, 60, 30], where object instances rather than parts form the node primitives. We also adopt a hierarchically structured approach towards part decomposition, but aim to operate on noisy, incomplete real-world scans of scenes with multiple objects, and so propose to combine our hierarchical part decomposition with strong geometric part priors.

## 3. Method

### 3.1. Overview

We address the problem of simultaneous part segmentation and completion of objects of real-world RGB-D scans, which are often noisy and incomplete. An overview of our approach is illustrated on Fig. 2. Given an input 3D scan $\mathbb{S}$, we aim to predict a set of parts for each object in the scan, with each part representing the complete geometry of the part, including any missing or unobserved regions. From $\mathbb{S}$, we first detect a set of object instances $\mathbb{O} = \{o_i\}$ in the scene, as 3D bounding box locations and class category predictions. For each detected object in $\mathbb{O}$, we then convert it into a $32^3$ occupancy grid representation, to inform our part segmentation and completion.

We then predict the part segmentation and completion for each detected object $o_i \in \mathbb{O}$. First, for a detected object $o_i$,

we predict its full part decomposition as a tree graph $T_i$ corresponding to the part hierarchy structure (nodes representing the part class types), with part hierarchies derived from those of PartNet [33]. This enables encoding the high-level, semantic part structure of the shape, which both facilitates completion of the shape structure, as missing parts are easy to identify in their part tree structure, as well as guides the prediction of the geometry of each part. In particular, this allows us to leverage geometric part priors built for each part category. We construct the part priors based on clustering of train part masks for each part category, and learn to predict similar priors for each leaf in our predicted $T_i$, followed by a refinement of these priors to predict the final part mask geometry. This produces a semantic part decomposition of objects in a 3D scan while simultaneously inferring their complete part geometry.

### 3.2. Object Detection

From an input 3D scan, we first detect objects in the scene. We leverage a state-of-the-art 3D object detection approach, MLCVNet [56], as our object detection backbone. The input scan sampled to a point cloud, and object proposals are produced by voting [37], leveraging global contextual information at various scales. As output, we obtain 3D bounding box locations for each detected object. We then resample the input scan geometry within each detected box into $32^3$ occupancy grids $o_i \in \mathbb{O}$ to inform our part decomposition.

For a detected object $o_i$ from the scan, represented as a $32^3$ occupancy grid of the scan geometry within its predicted bounding box, we encode the occupancy grid with four 3D convolutional blocks (consisting of convolution,

group normalization and ReLU activation) and extract a feature encoding $z_i$ of dimension 128, which is used to inform the part decomposition.

**Object Orientation Prediction** Since our object detection backbone predicts axis-aligned bounding boxes for each object, we additionally predict the orientation $r_i$ of each object $o_i$ from its feature $z_i$ using an MLP. We assume that the up (gravity) vector is known in the scene, and thus predict the angle around the up vector by classifying the angle in $n_\alpha = 8$ bins of discretized angles ($\{0°, 45°, \ldots, 315°\}$) with a cross entropy loss. The predicted object orientation helps to guide our prior-based part decomposition as described in Section 3.4.

### 3.3. Part Tree Prediction

For a detected object $o_i$ from the scan, represented as a $32^3$ occupancy grid of the scan geometry within its predicted bounding box, we aim to capture its high-level part structure from its cluttered and partial observation. We predict the part tree structure of the object; this facilitates completion of the object by predicting its high-level structure, as well as enables our prior-guided part geometry prediction.

We first encode the occupancy grid of $o_i$ with four 3D convolutional blocks (consisting of convolution, group normalization and ReLU activation), and extract a feature encoding $z_i$ of dimension 128. We then decode $z_i$ into a part tree prediction, constructing a part tree $T_i$ with each node represented by its predicted part category and a 128-dimensional feature encoding. Inspired by StructureNet [32], we leverage a message-passing graph neural network for our part tree prediction. From $z_i$, we predict tree children nodes using an MLP to predict $n_{children} = 10$ latent vectors $\{z'_k\}$ that correspond to potential parts of $o$. We additionally predict a tuple $t_k = (e_k, s_k)$ for every child $z'_k$, where $e_k$ is the probability of child existence, $s_k$ is the one-hot representation of the part category label. For each pair $(z'_i, z'_j)$ of nodes, we predict if they are adjacent or not, enforcing structural features to be learned by the message-passing network. We employ a cross entropy loss for the part category label, and binary cross entropy losses for node existence and adjacency relationships. This produces a high-level part summary of $o_i$, where nodes $\{z'_k\}$ represent part semantic information of the complete structure of $o_i$, even if $o_i$ has been partially observed. We leverage this part semantic information to guide our final part decomposition as geometric part masks.

### 3.4. Prior-guided Part Decomposition

We then predict the final part decomposition by generating part masks for each node in the predicted part tree $T_i$, where each mask represents the complete geometry associated with the part, including regions that were unobserved
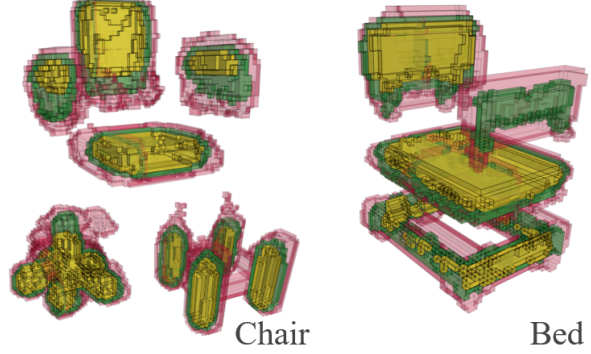


Chair          Bed

Figure 3: Several geometric part priors for part types belonging to the 'chair' and 'bed' class categories. Each part prior represents a cluster of train parts, visualized at three different isolevels.

in the initial scan observation. Rather than directly reconstructing the part geometry of each predicted part in the part tree, we observe that object parts often maintain very similar geometry structures, which we leverage to obtain our final part decomposition. That is, we construct geometric part priors to aid in generating our complete part mask predictions, and learn to find similar geometric part priors which we then refine for a final prediction.

We construct our geometric part priors by $k$-means clustering of the binary part masks in the train set, inspired by the ShapeMask [29] approach to building priors for novel 2D object segmentation. For each part type, we find $K = 10$ centroids of the part masks of that type, and perform the clustering on the part masks in $32^3$ grids of the canonical object space. This produces a set of part priors $\{P_1, \ldots, P_M\}$ with $M = n_{\text{classes}}K$. Various resulting part priors are visualized in Figure 3. Since objects in the real-world scan inputs may not be oriented in the canonical orientation of the object, we use the predicted orientation $r_i$ to transform the priors to $\{P_1^r, \ldots, P_M^r\}$.

Thus, to predict the part geometry associated for a node in the predicted part tree $T_i$ with feature encoding $z'_k$ and predicted part type $t$, we use a one-layer MLP which takes as input $z'_k$ and predicts a set of weights $w_m$ used to construct an initial part reconstruction as:

$$P_k^{\text{coarse}} = \sum_{m=1}^{M_t} w_m P_m^r,$$

where $w = softmax(\phi(z'_k))$, and $\phi$ is a linear layer. We employ a proxy loss on this initial part reconstruction, using a mean squared error with a target part mask.

Such prior-guided part decomposition helps to reconstruct global structures in part masks such as symmetry and geometry in missing regions in the input observation. We then refine the predicted $P_k^{\text{coarse}}$ using four 3D convolutional

| Method | Chamfer Distance (↓) | | | | | | | | IoU (↑) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg |
| SG-NN + MLCVNet + UNet | 0.050 | 0.118 | 0.080 | 0.053 | 0.083 | 0.108 | 0.082 | 0.073 | 17.5 | 6.4 | 7.6 | 12.4 | 13.3 | 13.9 | 11.9 | 13.3 |
| SG-NN + MLCVNet + PointGroup | 0.074 | 0.102 | 0.100 | 0.063 | 0.091 | 0.140 | 0.095 | 0.093 | 5.1 | 1.5 | 1.0 | 4.5 | 4.5 | 0.9 | 2.9 | 2.9 |
| MLCVNet + StructureNet | **0.029** | 0.095 | **0.065** | 0.037 | 0.076 | 0.106 | 0.068 | 0.057 | 13.8 | 0.5 | 3.8 | 9.0 | 3.9 | 9.3 | 6.8 | 8.9 |
| **Ours** | 0.033 | **0.089** | 0.069 | **0.033** | **0.054** | **0.096** | **0.062** | **0.053** | **22.1** | **7.7** | **13.0** | **18.1** | **17.3** | **22.0** | **16.7** | **18.3** |

Table 1: Evaluation on semantic part completion on Scan2CAD [1]. We compare with state-of-the-art approaches for scan completion [6], followed by object detection [56], and then part segmentation [25, 32]. By leveraging part structures to guide our prior-based approach, we obtain more accurate part decompositions.

blocks (consisting of convolution, batch normalization and ReLU activation) taking as input the concatenation of the geometry of $o_i$ and $P_k^{\text{coarse}}$ to produce $P_k^{\text{refine}}$; we then obtain the final part mask prediction,

$$P_k = P_k^{\text{coarse}} + P_k^{\text{refine}}.$$

Empirically, we found that predicting the refinement as residuals to modify the initial $P_k^{\text{coarse}}$ to perform better than a direct refinement (c.f. Section 4). We then employ a binary cross entropy loss on $P_k$ with a target part mask. This encourages an improved local fit to the observed geometry that may not have been captured in the global structure of the geometric priors.

### 3.5. Training Details

**Data generation.** In order to train our approach, we leverage the Scan2CAD dataset [1] in combination with PartNet [33]. Scan2CAD contains annotations of CAD models from ShapeNet [3] aligned to the 3D scans of ScanNet [5], and we use the part annotations of PartNet for these ShapeNet CAD models to obtain our ground truth part decompositions of the 3D scans. We leverage the ground truth CAD alignments to compute our geometric part priors in the canonically-oriented space of the objects, and use our rotation prediction during training and inference to orient them to the scan observations. In all our experiments we use original ScanNet geometry with typical number of points as 200k per scene, for MLCVNet method 40k points are randomly sampled from each scene to train object detection.

**Training.** We train our part tree prediction and geometric decomposition model with an Adam optimizer, using a batch size of 24, learning rate of 0.001, and weight decay of 0.01. The learning rate is decayed every 8 epochs by a factor of 0.8. We first pre-train for 20 epochs using ground truth 3D bounding boxes, and then fine-tune for 10 epochs with geometry from MLCVNet detections. MLCVNet is trained using the original proposed parameters: using an Adam optimizer with batch size 8, learning rate 0.01, for 250 epochs.

## 4. Results

We evaluate our proposed approach in comparison to alternative approaches for semantic part completion on real-world

RGB-D scans. We use scans from the ScanNet dataset [5], containing 1513 reconstructed RGB-D scans, and evaluate with their train/val/test split of 1045/156/312 scenes, respectively. To train and evaluate the complete part decomposition for each object, we use the Scan2CAD [1] annotations of CAD model alignments from ShapeNet [3] to the ScanNet scans, coupled with the PartNet [33] annotations for the part decomposition of the ShapeNet CAD models. We train and evaluate on 6 object class categories representing the majority of parts (45 part types in total that we train and evaluate on) for these annotations. For a detailed specification of the part types used, we refer to the appendix.

To evaluate our part decompositions of the objects in a scan, we use a Chamfer Distance metric as well as an intersection over union (IoU) metric. For IoU, we evaluate $32^3$ voxelizations of each predicted part in object space, compared to the Scan2CAD ground truth part. For Chamfer Distance, we use the predicted voxel centers as points, normalized to the unit box of the object. For both Chamfer Distance and IoU, we compute the metrics for each part type and average over all part types corresponding to an object class category. The class average is computed by averaging all resulting category numbers, and instance average computed by averaging the metrics of all part instances regardless of their object category. Note that to evaluate part segmentation without completion, we consider only predictions which overlap with the original scan geometry.

**Comparison to alternative approaches.** In Table 1, we compare to several state-of-the-art approaches for part segmentation and scan completion, coupled together to provide a complete part decomposition of the objects in a scan. As an alternative approach for this task, we consider scan completion followed by object detection and part instance segmentation. We employ the state-of-the-art scan completion approach SG-NN [6] to generate a prediction for the complete geometry of a partial scan observation, then detect object instances with MLCVNet [56], obtain a final complete part decomposition by the state-of-the-art instance segmentation of PointGroup [25]. We also compare to StructureNet [32] on MLCVNet detections, following their approach of using a pretraining a decoder for complete part decompositions and then learning an encoder to map this space. We additionally

5

| | Chamfer Distance (↓) | | | | | | | | IoU (↑) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg |
| MLCVNet + UNet | 0.052 | 0.082 | **0.062** | 0.034 | 0.093 | 0.068 | 0.065 | 0.060 | 24.1 | 13.4 | 9.3 | **31.8** | 14.3 | 14.6 | 17.9 | 18.9 |
| MLCVNet + PointGroup | 0.054 | **0.057** | 0.077 | 0.045 | **0.072** | 0.086 | 0.065 | 0.061 | 28.4 | 14.9 | 9.6 | 27.5 | 18.8 | 11.9 | 18.5 | 19.6 |
| MLCVNet + StructureNet | **0.039** | 0.084 | **0.062** | 0.034 | 0.075 | 0.083 | 0.063 | 0.056 | **32.6** | 2.1 | 9.4 | 23.1 | 16.1 | 15.4 | 16.5 | 15.4 |
| **Ours** | 0.044 | 0.072 | 0.063 | **0.031** | 0.092 | **0.063** | **0.061** | **0.054** | 30.9 | **16.5** | **10.9** | 31.8 | **20.6** | **20.9** | **21.9** | **24.7** |

Table 2: Evaluation of part segmentation on Scan2CAD [1]. We evaluate part segmentation of visible geometry only, in comparison with state-of-the-art part segmentation [25, 32].
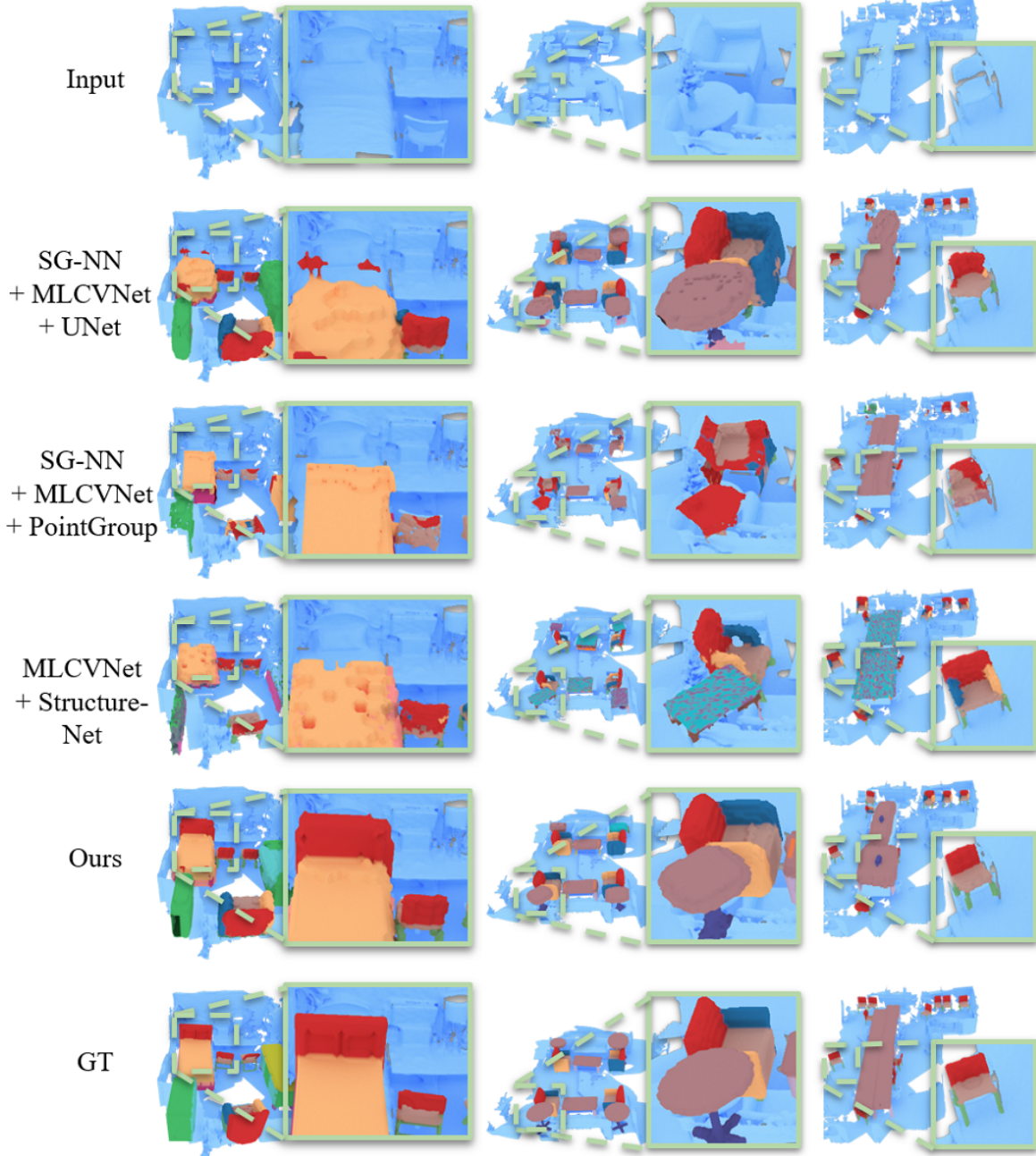


Figure 4: Qualitative evaluation on semantic part completion in comparison with state-of-the-art approaches for part decomposition, including scan completion followed by part segmentation. Our approach produces more consistent, accurate part decompositions.

| | Chamfer Distance (↓) | | | | | | | IoU (↑) | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | chair | table | cab. | bkshlf | bed | bin | **avg** | chair | table | cab. | bkshlf | bed | bin | **avg** |
| MLCVNet + StructureNet | **0.0032** | 0.0106 | 0.0074 | 0.0046 | 0.0194 | **0.0025** | 0.0079 | 29.5 | 21.9 | 22.4 | 24.1 | 23.5 | 32.4 | 25.6 |
| RevealNet | 0.0035 | **0.0070** | 0.0043 | **0.0020** | 0.0076 | 0.0078 | 0.0053 | 35.1 | 26.2 | 46.1 | **38.3** | 19.6 | 24.7 | 31.7 |
| MLCVNet + UNet | 0.0038 | 0.0103 | **0.0011** | 0.0050 | 0.0119 | 0.0028 | 0.0059 | **39.7** | 30.0 | **62.6** | 28.2 | 17.4 | 37.7 | 35.9 |
| **Ours** | 0.0038 | 0.0075 | 0.0022 | 0.0053 | **0.0061** | 0.0045 | **0.0049** | 38.6 | **30.7** | 57.4 | 33.6 | **37.6** | **37.8** | **39.3** |

Table 3: Evaluation of instance completion on Scan2CAD [1]. We evaluate object completion as a union of predicted part decompositions, in comparison with state-of-the-art instance completion [21] and the union of StructureNet [32] parts as instances.

consider a UNet [44] composed of 3D volumetric convolutions as a baseline for the final part segmentation. We train these alternative approaches on our part decomposition data for ScanNet. These approaches do not consider explicit part structure reasoning, whereas our part tree prediction guiding our part decomposition with geometric priors enables a more effective complete part decomposition.

In Figure 4, we show a qualitative comparison: without part structure reasoning, the PointGroup approach can often mix up geometrically similar parts such as the left and right chair arms, and the UNet baseline suffers in generating complete part structures. StructureNet provides part structure reasoning, but their approach to train an encoder into a pretrained decoder can tend to predict only the dominant part decompositions for a class category (e.g., an office-type chair instead of an armchair in the third row of Figure 4). Our part structure guided priors enable more effective and accurate part decompositions of the objects in the scenes.

**Part segmentation on 3D scans.** In addition to our task of semantic part completion, we evaluate our approach in comparison to state of the art on part segmentation in Table 2. To evaluate part segmentation, we consider only the part predictions that intersect with the original scan geometry, and compare to PointGroup [25], StructureNet [32], and a UNet baseline, using the object detection of by MLCVNet [56]. For part segmentation, we see that our part structure reasoning coupled with geometric priors also produces more consistent part segmentations of the objects in a scan.

**Object completion on 3D scans.** In Table 3, we additionally evaluate our approach on object instance completion by taking the union of our part mask predictions as a com-

plete object mask prediction. We compare to RevealNet [21], which established this task, as well as a state-of-the-art object detection using MLCVNet [56] followed by a UNet for completion or by StructureNet [32]. Our part reasoning enables more effective instance completion by explicitly leveraging shared structural knowledge of objects.

**Ablations.** In Table 4, we analyze the effect of our design decisions for part tree and prior-guided part mask prediction. We evaluate our approach without message-passing in our part tree prediction (*w/o Part Msg Pass*), without using priors and directly decoding with convolutions to a part mask prediction (*w/o Priors*), without refinement of priors (*No Prior Refine*), and prior refinement with absolute predictions instead of our relative offsets that are added to the raw prior prediction (*Prior Refine (Abs)*). Our prior-guided predictions, with refinement learned as a residual offset, helps to produce more accurate results.

We additionally consider the effect of varying voxel resolutions in Table 5. All resolutions produce meaningful results, although a (twice) higher resolution can result in somewhat noisier results, and a (half) lower resolution tends to lack detail. Thus we choose to employ a $32^3$ resolution for each object.

**Limitations.** While our approach for semantic part completion shows promise towards a finer-grained, semantically part-based understanding of 3D environments, we believe there are many avenues for further development. For instance, a dense volumetric representation of parts may suffice for functionality analysis of furniture-type objects, but can struggle to generate very high resolution parts for small objects; we believe sparse [15, 4] or hierarchical [43, 51]

| | Chamfer Distance (↓) | | | | | | | | IoU (↑) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg |
| w/o Part Msg Pass | 0.037 | 0.094 | 0.069 | 0.039 | 0.077 | 0.096 | 0.069 | 0.057 | 20.0 | 6.6 | 9.8 | 14.2 | 14.4 | 20.6 | 14.3 | 16.3 |
| w/o Priors | 0.036 | 0.093 | 0.067 | 0.044 | 0.058 | 0.101 | 0.067 | 0.056 | 21.8 | 7.3 | 11.0 | 13.8 | 16.4 | 21.9 | 15.4 | 17.7 |
| No Prior Refine | 0.034 | 0.093 | 0.069 | 0.034 | 0.057 | 0.096 | 0.064 | 0.055 | **22.5** | 7.6 | 12.2 | 17.9 | 16.6 | **22.0** | 16.4 | 18.2 |
| Prior Refine (Abs) | 0.036 | **0.089** | **0.065** | 0.034 | 0.067 | 0.105 | 0.066 | 0.055 | 21.4 | 7.5 | 11.5 | 17.4 | 16.5 | 20.7 | 15.8 | 17.6 |
| **Ours** | **0.033** | **0.089** | 0.069 | **0.033** | **0.054** | **0.096** | **0.062** | **0.053** | 22.1 | **7.7** | **13.0** | **18.1** | **17.3** | **22.0** | **16.7** | **18.3** |

Table 4: Ablation study for our design decisions, evaluated for semantic part completion on Scan2CAD [1].
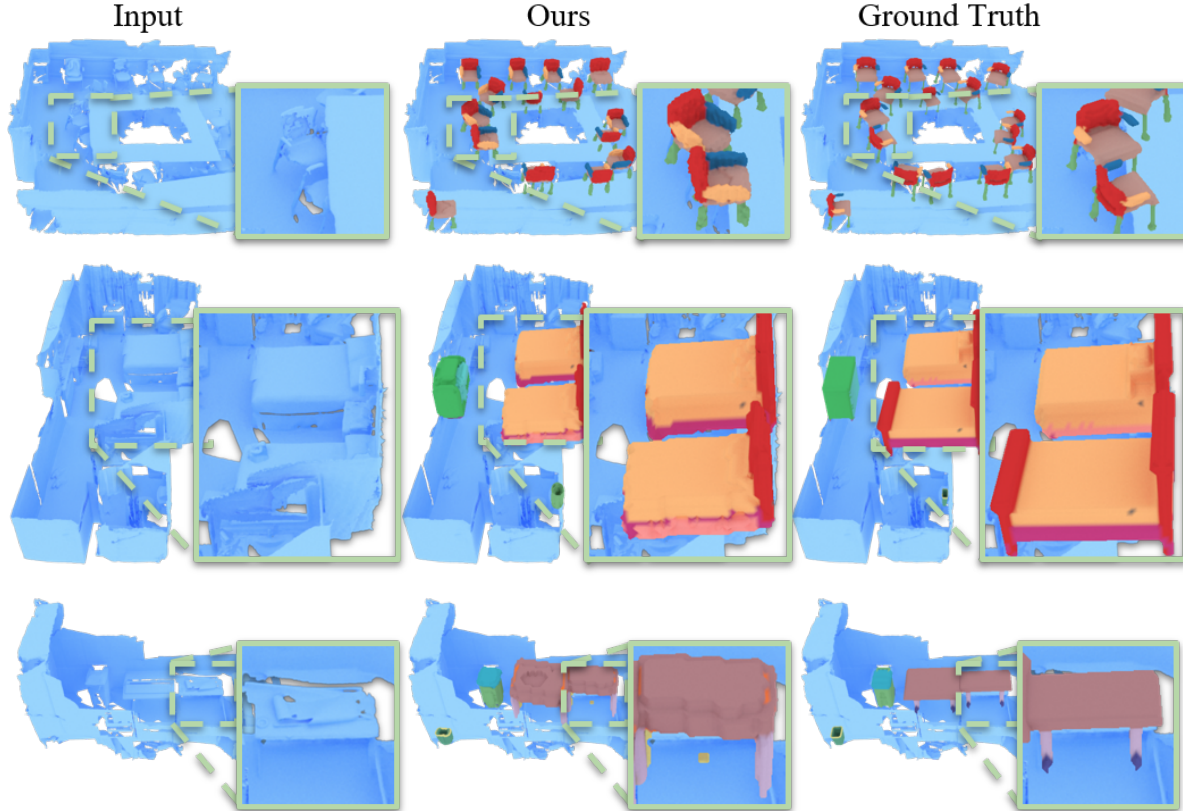
Figure 5: Qualitative results on real-world ScanNet [5] scenes using Scan2CAD [1] and PartNet [33] targets. Our approach effectively predicts each object's complete geometry as a decomposition into semantic parts.

## 5. Conclusion

In this paper, we have presented a new approach for the semantic part completion task of predicting a geometrically complete part decomposition for each object in a 3D scan. For each detected object in a scene, we exploit explicit part structure prediction in order to guide a geometric part prior prediction, which is then refined to a final part decomposi-

approaches would complement our prior-based approach. Furthermore, objects are currently considered independently for each part decomposition, where relational inference between objects in a scene would help to explain noisy or unobserved part regions (e.g., multiple chairs or tables in a scene are often repeated instances of the same geometry).

tion, where each part is represented by its semantic part type as well as the geometry corresponding to the part, including any missing or unobserved regions in the scan. We show that our structural and prior-guided reasoning about object parts notably outperforms alternative approaches on this task. We believe that our approach makes an important step towards part-based understanding of 3D environments, and opens up new possibilities for part-level functionality analysis, autonomous agent interactions with an environment, and more.

## 6. Acknowledgements

| Method | Chamfer Distance (↓) | | | | | | | | IoU (↑) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg |
| Res. 16 | 0.034 | 0.088 | 0.072 | 0.054 | 0.061 | 0.109 | 0.070 | 0.055 | 28.4 | 10.5 | 13.5 | 20.9 | 18.5 | 21.2 | **18.8** | **22.8** |
| Res. 32 | 0.033 | 0.089 | 0.069 | 0.033 | 0.054 | 0.096 | **0.062** | **0.053** | 22.1 | 7.7 | 13.0 | 18.1 | 17.3 | 22.0 | 16.7 | 18.3 |
| Res. 64 | 0.045 | 0.098 | 0.058 | 0.044 | 0.067 | 0.100 | 0.069 | 0.060 | 18.8 | 5.6 | 9.9 | 10.5 | 14.7 | 19.3 | 13.1 | 15.4 |

Table 5: Evaluation of various object resolutions during training for semantic part completion on Scan2CAD [1].

# References

[1] Armen Avetisyan, Manuel Dahnert, Angela Dai, Manolis Savva, Angel X. Chang, and Matthias Nießner. Scan2cad: Learning CAD model alignment in RGB-D scans. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 2614–2623, 2019. 5, 6, 7, 8, 12, 16

[2] Angel X. Chang, Angela Dai, Thomas A. Funkhouser, Maciej Halber, Matthias Nießner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from RGB-D data in indoor environments. In *2017 International Conference on 3D Vision, 3DV 2017, Qingdao, China, October 10-12, 2017*, pages 667–676, 2017. 1

[3] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 1, 5

[4] Christopher B. Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 3075–3084, 2019. 1, 7

[5] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Niessner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 1, 5, 8, 16

[6] Angela Dai, Christian Diller, and Matthias Nießner. Sg-nn: Sparse generative neural networks for self-supervised scene completion of rgb-d scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 849–858, 2020. 2, 5, 12

[7] Angela Dai and Matthias Nießner. 3dmv: Joint 3d-multi-view prediction for 3d semantic scene segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 452–468, 2018. 1

[8] Angela Dai, Daniel Ritchie, Martin Bokeloh, Scott Reed, Jürgen Sturm, and Matthias Nießner. Scancomplete: Large-scale scene completion and semantic segmentation for 3d scans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2018. 2

[9] Angela Dai, Charles Ruizhongtai Qi, and Matthias Nießner. Shape completion using 3d-encoder-predictor cnns and shape synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5868–5877, 2017. 2

[10] Francis Engelmann, Martin Bokeloh, Alireza Fathi, Bastian Leibe, and Matthias Nießner. 3d-mpa: Multi-proposal aggregation for 3d semantic instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9031–9040, 2020. 2

[11] Matthew Fisher, Daniel Ritchie, Manolis Savva, Thomas Funkhouser, and Pat Hanrahan. Example-based synthesis

[12] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 2

[13] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9785–9795, 2019. 1

[14] Aleksey Golovinskiy and Thomas Funkhouser. Consistent segmentation of 3d models. *Computers & Graphics*, 33(3):262–269, 2009. 2

[15] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 9224–9232, 2018. 1, 7

[16] Lei Han, Tian Zheng, Lan Xu, and Lu Fang. Occuseg: Occupancy-aware 3d instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2940–2949, 2020. 1, 2

[17] Christian Häne, Shubham Tulsiani, and Jitendra Malik. Hierarchical surface prediction for 3d object reconstruction. In *2017 International Conference on 3D Vision (3DV)*, pages 412–420. IEEE, 2017. 2

[18] Rana Hanocka, Amir Hertz, Noa Fish, Raja Giryes, Shachar Fleishman, and Daniel Cohen-Or. Meshcnn: a network with an edge. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 1, 3

[19] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 2

[20] Ji Hou, Angela Dai, and Matthias Nießner. 3d-sis: 3d semantic instance segmentation of rgb-d scans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4421–4430, 2019. 2

[21] Ji Hou, Angela Dai, and Matthias Nießner. Revealnet: Seeing behind objects in rgb-d scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2098–2107, 2020. 2, 7, 12

[22] Ruizhen Hu, Lubin Fan, and Ligang Liu. Co-segmentation of 3d shapes via subspace clustering. In *Computer graphics forum*, volume 31, pages 1703–1713. Wiley Online Library, 2012. 2

[23] Jingwei Huang, Haotian Zhang, Li Yi, Thomas Funkhouser, Matthias Nießner, and Leonidas J Guibas. Texturenet: Consistent local parametrizations for learning from high-resolution signals on meshes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4440–4449, 2019. 1

[24] Qixing Huang, Vladlen Koltun, and Leonidas Guibas. Joint shape segmentation with linear programming. In *Proceedings of the 2011 SIGGRAPH Asia Conference*, pages 1–12, 2011. 2

[25] Li Jiang, Hengshuang Zhao, Shaoshuai Shi, Shu Liu, Chi-Wing Fu, and Jiaya Jia. Pointgroup: Dual-set point grouping for 3d instance segmentation. In *Proceedings of the*

[11] (continued) of 3d object arrangements. *ACM Transactions on Graphics (TOG)*, 31(6):1–11, 2012. 3

*IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4867–4876, 2020. 1, 5, 6, 7

[26] Evangelos Kalogerakis, Melinos Averkiou, Subhransu Maji, and Siddhartha Chaudhuri. 3d shape segmentation with projective convolutional networks. In *proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3779–3788, 2017. 3

[27] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, volume 7, 2006. 2

[28] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013. 2

[29] Weicheng Kuo, Anelia Angelova, Jitendra Malik, and Tsung-Yi Lin. Shapemask: Learning to segment novel objects by refining shape priors. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9207–9216, 2019. 4

[30] Manyi Li, Akshay Gadi Patil, Kai Xu, Siddhartha Chaudhuri, Owais Khan, Ariel Shamir, Changhe Tu, Baoquan Chen, Daniel Cohen-Or, and Hao Zhang. Grains: Generative recursive autoencoders for indoor scenes. *ACM Transactions on Graphics (TOG)*, 38(2):1–16, 2019. 3

[31] Niloy J Mitra, Leonidas J Guibas, and Mark Pauly. Partial and approximate symmetry detection for 3d geometry. *ACM Transactions on Graphics (TOG)*, 25(3):560–568, 2006. 2

[32] Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas J Guibas. Structurenet: Hierarchical graph networks for 3d shape generation. *arXiv preprint arXiv:1908.00575*, 2019. 1, 3, 4, 5, 6, 7, 12

[33] Kaichun Mo, Shilin Zhu, Angel X Chang, Li Yi, Subarna Tripathi, Leonidas J Guibas, and Hao Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 909–918, 2019. 1, 3, 5, 8, 16

[34] Andrew Nealen, Takeo Igarashi, Olga Sorkine, and Marc Alexa. Laplacian mesh optimization. In *Proceedings of the 4th international conference on Computer graphics and interactive techniques in Australasia and Southeast Asia*, pages 381–389, 2006. 2

[35] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 2

[36] Mark Pauly, Niloy J Mitra, Johannes Wallner, Helmut Pottmann, and Leonidas J Guibas. Discovering structural regularity in 3d geometry. In *ACM SIGGRAPH 2008 papers*, pages 1–11. 2008. 2

[37] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9277–9286, 2019. 2, 3

[38] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 77–85, 2017. 1

[39] Charles R Qi, Hao Su, Matthias Nießner, Angela Dai, Mengyuan Yan, and Leonidas J Guibas. Volumetric and multi-view cnns for object classification on 3d data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5648–5656, 2016. 1

[40] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 5099–5108, 2017. 1

[41] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. 2

[42] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 2

[43] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3577–3586, 2017. 7

[44] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 7

[45] Oana Sidi, Oliver van Kaick, Yanir Kleiman, Hao Zhang, and Daniel Cohen-Or. Unsupervised co-segmentation of a set of shapes via descriptor-space spectral clustering. In *Proceedings of the 2011 SIGGRAPH Asia Conference*, pages 1–10, 2011. 2

[46] Ivan Sipiran, Robert Gregor, and Tobias Schreck. Approximate symmetry detection in partial 3d meshes. In *Computer Graphics Forum*, volume 33, pages 131–140. Wiley Online Library, 2014. 2

[47] Shuran Song and Jianxiong Xiao. Deep sliding shapes for amodal 3d object detection in rgb-d images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 808–816, 2016. 2

[48] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1746–1754, 2017. 2

[49] Pablo Speciale, Martin R Oswald, Andrea Cohen, and Marc Pollefeys. A symmetry prior for convex variational 3d reconstruction. In *European Conference on Computer Vision*, pages 313–328. Springer, 2016. 2

[50] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015. 1

[51] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2088–2096, 2017. 7

[52] Sebastian Thrun and Ben Wegbreit. Shape from symmetry. In *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, volume 2, pages 1824–1831. IEEE, 2005. 2

[53] Oliver Van Kaick, Kai Xu, Hao Zhang, Yanzhen Wang, Shuyang Sun, Ariel Shamir, and Daniel Cohen-Or. Co-hierarchical analysis of shape structures. *ACM Transactions on Graphics (TOG)*, 32(4):1–10, 2013. 3

[54] Yanzhen Wang, Kai Xu, Jun Li, Hao Zhang, Ariel Shamir, Ligang Liu, Zhiquan Cheng, and Yueshan Xiong. Symmetry hierarchy of man-made objects. In *Computer graphics forum*, volume 30, pages 287–296. Wiley Online Library, 2011. 3

[55] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1912–1920, 2015. 1, 2

[56] Qian Xie, Yu-Kun Lai, Jing Wu, Zhoutao Wang, Yiming Zhang, Kai Xu, and Jun Wang. Mlcvnet: Multi-level context votenet for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10447–10456, 2020. 2, 3, 5, 7, 12

[57] Li Yi, Leonidas Guibas, Aaron Hertzmann, Vladimir G Kim, Hao Su, and Ersin Yumer. Learning hierarchical shape segmentation and labeling from online repositories. *arXiv preprint arXiv:1705.01661*, 2017. 3

[58] Li Yi, Vladimir G Kim, Duygu Ceylan, I-Chao Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, and Leonidas Guibas. A scalable active framework for region annotation in 3d shape collections. *ACM Transactions on Graphics (ToG)*, 35(6):1–12, 2016. 3

[59] Wei Zhao, Shuming Gao, and Hongwei Lin. A robust hole-filling algorithm for triangular mesh. *The Visual Computer*, 23(12):987–997, 2007. 2

[60] Xi Zhao, Ruizhen Hu, Paul Guerrero, Niloy Mitra, and Taku Komura. Relationship templates for creating scene variations. *ACM Transactions on Graphics (TOG)*, 35(6):1–13, 2016. 3

## A. Appendix

In this appendix, we detail our network architecture in Section B; in Section C, we provide details of our baselines designs; in Section D, we provide specifications of parts that we used in our experiments; in Section E, we additionally provide more quantitative results, visualize examples of part priors combinations for each main category and examples of our predictions compared to ground-truth.

## B. Network Architecture Details

We detail our network architecture specification in Tables 8-9. Table 8 describes the layers for encoding the detected objects to a feature code. The feature code is then input to a decoder which predicts the part tree, as detailed in Table 10; here, the output of the last layer, lin3, represents a tuple of children latent codes, which predict part prior weights, as specified in Section 3.4 of the main paper. The final part refinement is then described in Table 9. Our volumetric object encoder and part refinement are fully convolutional, while the part tree prediction operates on the latent feature representations of shapes and parts with MLP structure.

## C. Additional Baseline Training Details

In all our experiments in comparison with state of the art, we leveraged a combination of various approaches. For the task of Semantic Part Completion, we performed scan completion with SG-NN [6] and object detection with ML-CVNet [56]. Our UNet baseline is developed as a baseline without any part tree or geometric part prior inference; it consists of only a 3D voxel encoder (four convolutional blocks consisting of 3D convolution, Group Normalization, ReLU activation) and 3D voxel decoder (five convolutional blocks consisting of 3D transposed convolution, 3D convolution, Group Normalization, ReLU activation) with 45 output feature channels, corresponding to binary masks for each part type, and trained with a binary cross entropy loss. Without the explicit part structure representations, this UNet baseline tends to predict noisy part masks, or part types from incorrect classes which remain functionally different.

Note that for experiments with StructureNet [32], we used the same experimental setup as described in their original paper, training different models for each class category. Since StructureNet operates in the canonical space of the objects, we provided our predicted object orientations from our approach to guide the StructureNet predictions.

## D. Part Types

In Figure 6, we visualize all part types which we trained on. Note that the classes 'cabinet' and 'bookshelf' share the same set of parts, so we use the same part types and priors.

## E. Additional Results

**Additional Quantitative Results**   In Table 6 we additionally evaluate object instance completion using an mAP@25 metric, in comparison to state-of-the-art RevealNet [21] and a combination of MLCVNet [56] with StructureNet [32]. Additionally, in Table 7, we evaluate our approach with ground truth 3D detection, i.e., ground truth oriented 3D bounding boxes for each object in the scene. Under ground truth detection, our structural part priors enable more robust part decomposition than StructureNet [32].

**Additional Part Prior Visualizations**   We show additional examples of computed part priors for each object class category in Figure 7. All priors are visualized with three level-sets.

**Additional Qualitative Semantic Part Completion Results**   Figure 8 shows additional examples of our predictions compared with ground-truth. Our method predicts meaningful part completion across a variety of object categories.

| Method | mAP@25 (↑) | | | | | | |
|---|---|---|---|---|---|---|---|
| | chair | table | cab. | bkshlf | bed | bin | **avg** |
| MLCVNet + StructureNet | 45.7 | 25.7 | 19.8 | 50.0 | 36.4 | 53.0 | 38.4 |
| RevealNet | 70.3 | 40.6 | 90.5 | 87.2 | 22.7 | 20.6 | 55.3 |
| Ours | 78.4 | 47.2 | 90.5 | 77.8 | 22.7 | 72.4 | 64.8 |

Table 6: Evaluation of instance completion on Scan2CAD [1]. We evaluate object completion as a union of predicted part decompositions, in comparison with state-of-the-art instance completion [21] and the union of StructureNet [32] parts as instances.

| Method | Chamfer Distance (↓) | | | | | | | | IoU (↑) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg | chair | table | cab. | bkshlf | bed | bin | class avg | inst avg |
| StructureNet [32] | 0.019 | 0.089 | 0.048 | 0.032 | 0.069 | 0.105 | 0.061 | 0.049 | 18.5 | 1.0 | 10.1 | 16.8 | 6.8 | 12.1 | 10.9 | 12.8 |
| **Ours** | 0.029 | 0.089 | 0.055 | 0.037 | 0.058 | 0.081 | 0.058 | 0.048 | 27.6 | 8.0 | 17.3 | 20.9 | 19.8 | 28.7 | 20.4 | 22.6 |

Table 7: Evaluation on semantic part completion on Scan2CAD [1] with ground truth 3D object detection (oriented 3D bounding boxes) as input.

| Encoder | Input Layer | Type | Input Size | Output Size | Kernel Size | Stride | Padding |
|---|---|---|---|---|---|---|---|
| conv0 | scan occ. grid | Conv3D | (1, 32, 32, 32) | (16, 16, 16, 16) | (5, 5, 5) | (2, 2, 2) | (2, 2, 2) |
| gnorm0 | conv0 | GroupNorm | (16, 16, 16, 16) | (16, 16, 16, 16) | - | - | - |
| relu0 | gnorm0 | ReLU | (16, 16, 16, 16) | (16, 16, 16, 16) | - | - | - |
| pool1 | relu0 | MaxPooling | (16, 16, 16, 16) | (16, 8, 8, 8) | (2, 2, 2) | (2, 2, 2) | (0, 0, 0) |
| conv1 | pool1 | Conv3D | (16, 8, 8, 8) | (32, 8, 8, 8) | (3, 3, 3) | (1, 1, 1) | (1, 1, 1) |
| gnorm1 | conv1 | GroupNorm | (32, 8, 8, 8) | (32, 8, 8, 8) | - | - | - |
| relu1 | gnorm1 | ReLU | (32, 8, 8, 8) | (32, 8, 8, 8) | - | - | - |
| pool2 | relu1 | MaxPooling | (32, 8, 8, 8) | (32, 4, 4, 4) | (2, 2, 2) | (2, 2, 2) | (0, 0, 0) |
| conv2 | pool2 | Conv3D | (32, 4, 4, 4) | (64, 2, 2, 2) | (5, 5, 5) | (2, 2, 2) | (2, 2, 2) |
| gnorm2 | conv2 | GroupNorm | (64, 2, 2, 2) | (64, 2, 2, 2) | - | - | - |
| relu2 | gnorm2 | ReLU | (64, 2, 2, 2) | (64, 2, 2, 2) | - | - | - |
| pool3 | relu2 | MaxPooling | (64, 2, 2, 2) | (64, 1, 1, 1) | (2, 2, 2) | (2, 2, 2) | (0, 0, 0) |
| conv3 | pool3 | Conv3D | (64, 1, 1, 1) | (128, 1, 1, 1) | (1, 1, 1) | (1, 1, 1) | (0, 0, 0) |
| gnorm3 | conv3 | GroupNorm | (128, 1, 1, 1) | (128, 1, 1, 1) | - | - | - |
| relu3 | gnorm3 | ReLU | (128, 1, 1, 1) | (128, 1, 1, 1) | - | - | - |
| flat0 | node feature | Flatten | (128, 1, 1, 1) | (128) | - | - | - |

Table 8: Layer specification for detected object encoder.

| Child decoder | Input Layer | Type | Input Size | Output Size |
|---|---|---|---|---|
| lin0 | node feature | Linear | 128 | 1280 |
| relu0 | lin0 | ReLU | 1280 | 1280 |
| reshape0 | relu0 | Reshape | 1280 | (10, 128) |
| node_exist | reshape0 | Linear | (10, 128) | (10, 1) |
| concat0 | (reshape0, reshape0) | Concat. | (10, 128), (10, 128) | (10, 10, 256) |
| lin1 | concat0 | Linear | (10, 10, 256) | (10, 10, 128) |
| relu1 | lin1 | ReLU | (10, 10, 128) | (10, 10, 128) |
| edge_exist | relu1 | Linear | (10, 10, 128) | (10, 10, 1) |
| mp | (relu1, edge_exist, reshape0) | Mes. Passing | (10, 10, 128), (10, 10, 1), (10, 128) | (10, 384) |
| lin2 | mp | Linear | (10, 384) | (10, 128) |
| relu2 | lin2 | ReLU | (10, 128) | (10, 128) |
| node_sem | relu2 | Linear | (10, 128) | (10, #classes) |
| lin3 | relu2 | Linear | (10, 128) | (10, 128) |
| relu3 | lin3 | ReLU | (10, 128) | (10, 128) |

Table 9: Layer specification for decoding an object into a part tree.

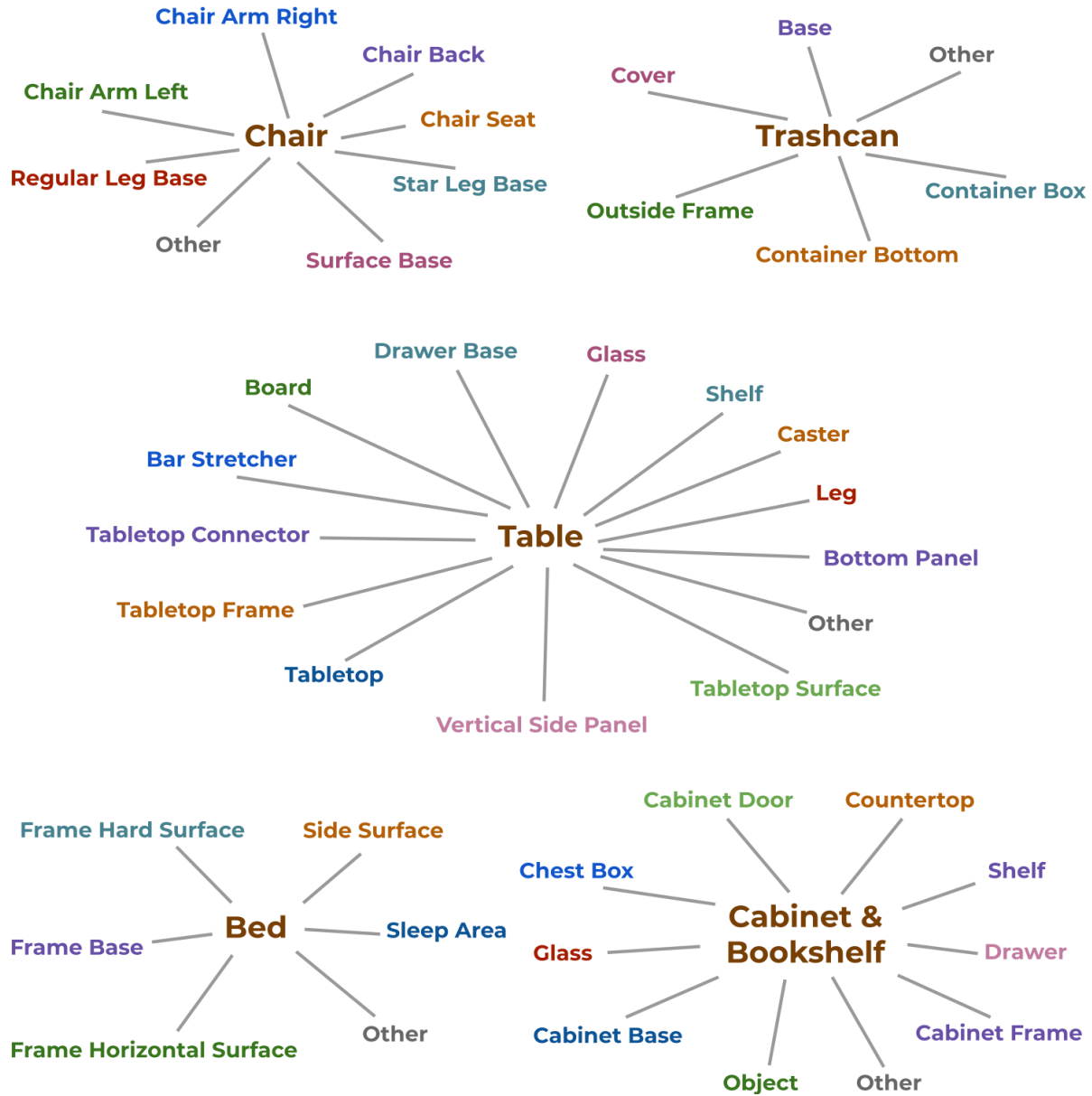| Prior refiner | Input Layer | Type | Input Size | Output Size | Kernel Size | Stride | Padding |
|---|---|---|---|---|---|---|---|
| concat0 | (prior, scan occ. grid) | Concat. | (1, 32, 32, 32), (1, 32, 32, 32) | (2, 32, 32, 32) | - | - | - |
| conv0 | concat0 | Conv3D | (2, 32, 32, 32) | (8, 32, 32, 32) | (3, 3, 3) | (1, 1, 1) | (1, 1, 1) |
| bnorm0 | conv0 | BatchNorm | (8, 32, 32, 32) | (8, 32, 32, 32) | - | - | - |
| relu0 | bnorm0 | ReLU | (8, 32, 32, 32) | (8, 32, 32, 32) | - | - | - |
| conv1 | relu0 | Conv3D | (8, 32, 32, 32) | (16, 32, 32, 32) | (3, 3, 3) | (1, 1, 1) | (1, 1, 1) |
| bnorm1 | conv1 | BatchNorm | (16, 32, 32, 32) | (16, 32, 32, 32) | - | - | - |
| relu1 | bnorm1 | ReLU | (16, 32, 32, 32) | (16, 32, 32, 32) | - | - | - |
| conv2 | relu1 | Conv3D | (16, 32, 32, 32) | (8, 32, 32, 32) | (3, 3, 3) | (1, 1, 1) | (1, 1, 1) |
| bnorm2 | conv2 | BatchNorm | (8, 32, 32, 32) | (8, 32, 32, 32) | - | - | - |
| relu2 | bnorm2 | ReLU | (8, 32, 32, 32) | (8, 32, 32, 32) | - | - | - |
| conv3 | relu2 | Conv3D | (8, 32, 32, 32) | (1, 32, 32, 32) | (1, 1, 1) | (1, 1, 1) | (0, 0, 0) |
| add3 | (prior, conv3) | Add | (1, 32, 32, 32), (1, 32, 32, 32) | (1, 32, 32, 32) | - | - | - |
| sigmoid3 | add3 | Sigmoid | (1, 32, 32, 32) | (1, 32, 32, 32) | - | - | - |

Table 10: Layer specification for final part mask refinement.

Figure 6: Part specification for the parts used in our approach. Note that 'cabinet' and 'bookshelf' classes have the same set of parts.

Chairs
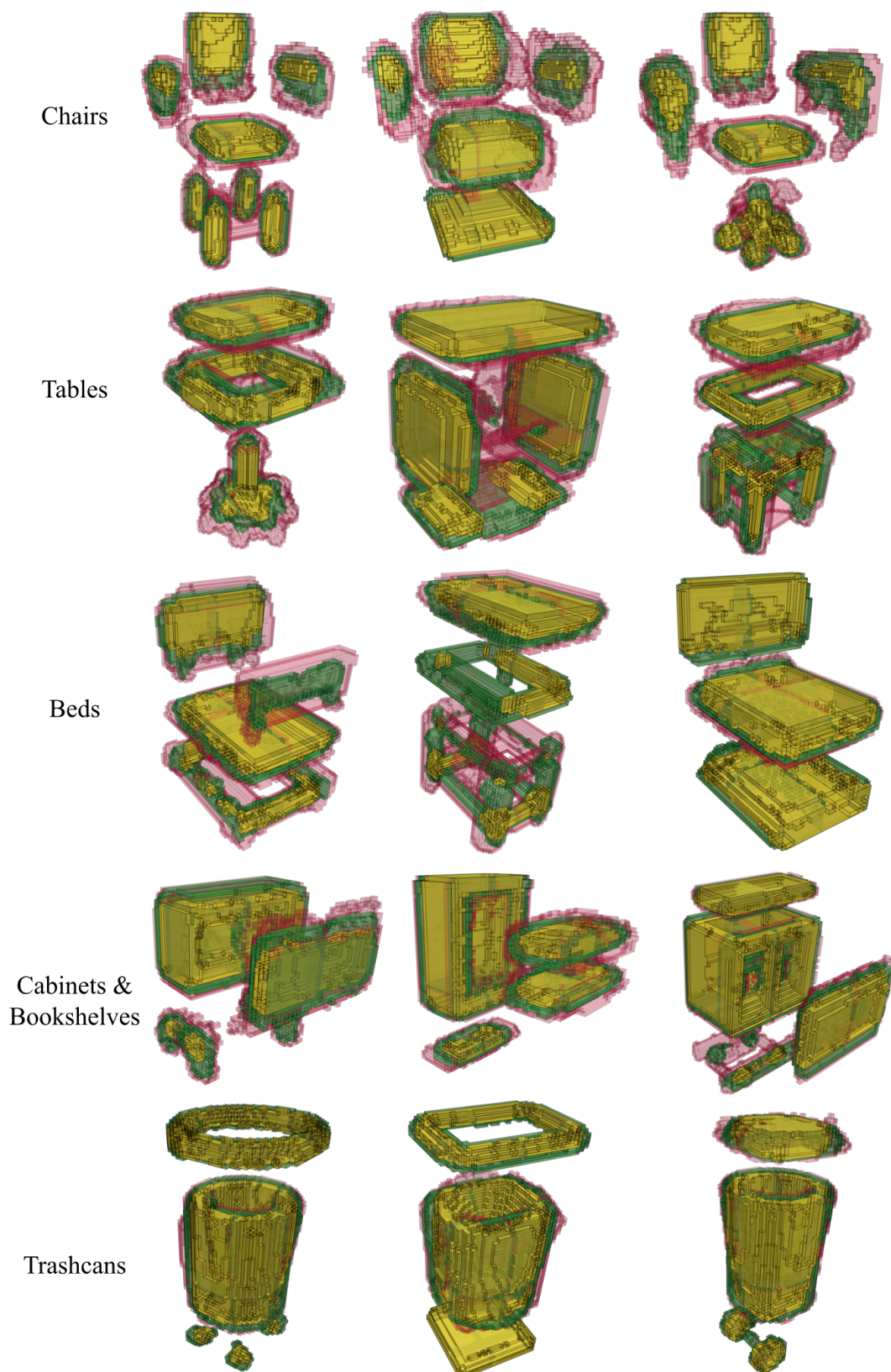
Tables

Beds

Cabinets &
Bookshelves

Trashcans

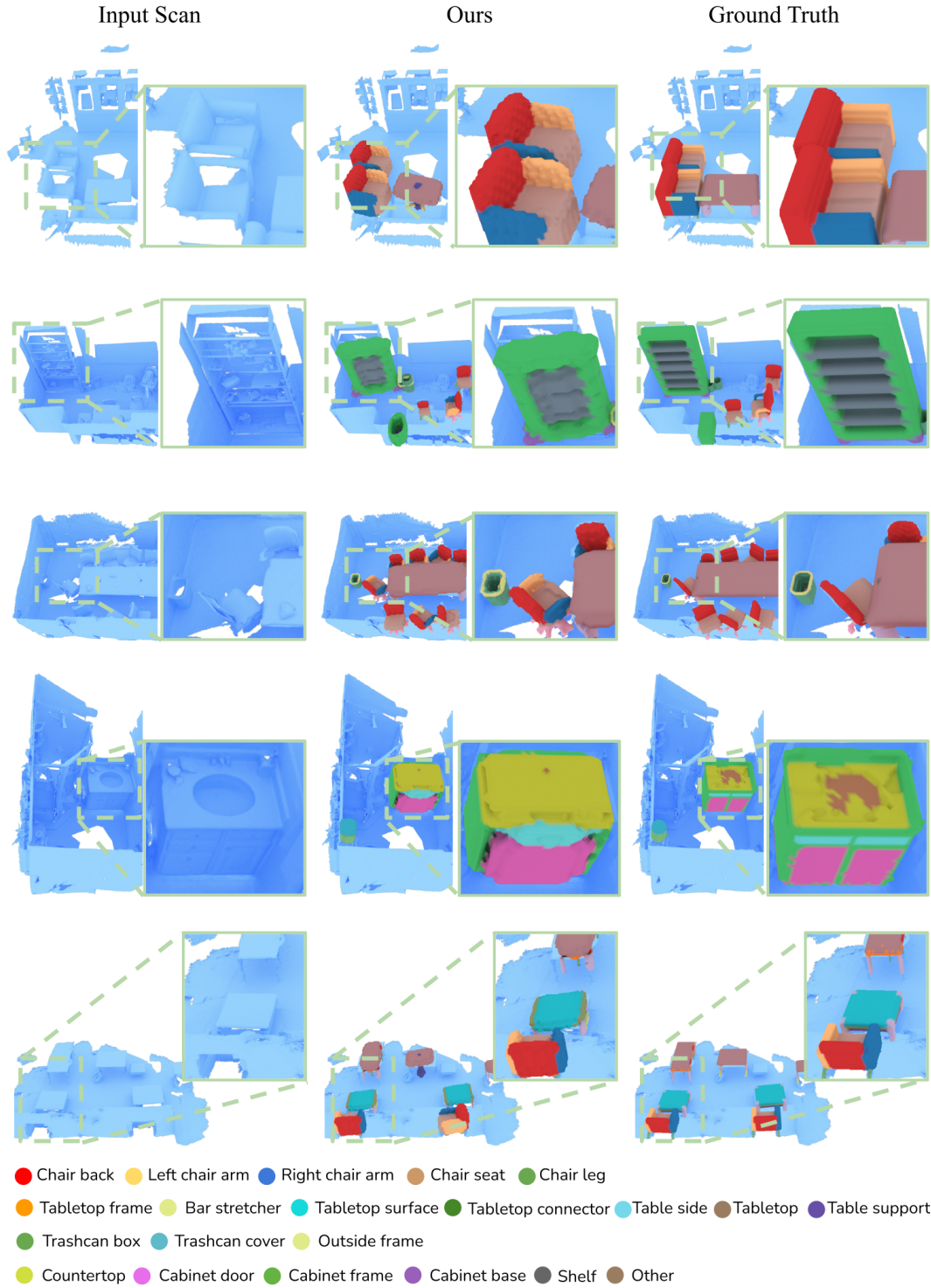Figure 7: Visualization of various part priors.

Figure 8: Additional qualitative results for our method on ScanNet [5] scenes and ground truth from Scan2CAD [1] and PartNet [33].