# Learning 3D Part Assembly from a Single Image

Yichen Li[*1], Kaichun Mo[*1], Lin Shao[1], Minhyuk Sung[2], and Leonidas Guibas[1]

[1] Stanford University
[2] Adobe Research

**Abstract.** Autonomous assembly is a crucial capability for robots in many applications. For this task, several problems such as obstacle avoidance, motion planning, and actuator control have been extensively studied in robotics. However, when it comes to task specification, the space of possibilities remains underexplored. Towards this end, we introduce a novel problem, *single-image-guided 3D part assembly*, along with a learning-based solution. We study this problem in the setting of *furniture assembly* from a given complete set of parts and a single image depicting the entire assembled object. Multiple challenges exist in this setting, including handling ambiguity among parts (*e.g.*, slats in a chair back and leg stretchers) and 3D pose prediction for parts and part subassemblies, whether visible or occluded. We address these issues by proposing a two-module pipeline that leverages strong 2D-3D correspondences and assembly-oriented graph message-passing to infer part relationships. In experiments with a PartNet-based synthetic benchmark, we demonstrate the effectiveness of our framework as compared with three baseline approaches.

**Keywords:** single-image 3D part assembly, vision for robotic assembly.

## 1 Introduction

The important and seemingly straightforward task of furniture assembly presents serious difficulties for autonomous robots. A general robotic assembly task consists of action sequences incorporating the following stages: (1) picking up a particular part, (2) moving it to a desired 6D pose, (3) mating it precisely with the other parts, (4) returning the manipulator to a pose appropriate for the next pick-up movement. Solving such a complicated high-dimensional motion planning problem [25,21] requires considerable time and engineering effort. Current robotic assembly solutions first determine the desired 6D pose of parts [9] and then hard-code the motion trajectories for each specific object [54]. Such limited generalizability and painstaking process planning fail to meet demands for fast and flexible industrial manufacturing and household assembly tasks [31].

To generate smooth and collision-free motion planning and control solutions, it is required to accurately predict 6D poses of parts in 3D space [54,27]. We propose a *3D part assembly task* whose output can reduce the complexity of
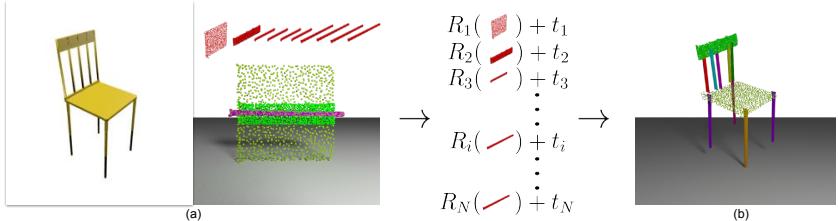
---

[*] :indicates equal contributions.

Fig. 1: **Single-Image-Based 3D Part Assembly Task.** Given as inputs an image and a set of part point clouds depicted in (a), the task is to predict 6D part poses in camera coordinates that assemble the parts to a 3D shape in the given image as shown in (b).

the high-dimensional motion planning problem. We aim to learn generalizable skills that allow robots to autonomously assemble unseen objects from parts [16]. Instead of hand-crafting a fixed set of rules to assemble one specific chair, for example, we explore category-wise structural priors that helps robots to assemble all kinds of chairs. The shared part relationships across instances in a category not only suggest potential pose estimation solutions for unseen objects but also lead to possible generalization ability for robotic control policies [64,53,42,60].

We introduce the task of *single-image-guided 3D part assembly*: inducing 6D poses of the parts in 3D space [30] from a set of 3D parts and an image depicting the complete object. Robots can acquire geometry information for each part using 3D sensing, but the only information provided for the entire object shape is the instruction image. Different from many structure-aware shape modeling works [40,71,17,62,70,32,52], we do not assume any specific granularity or semantics of the input parts, since the given furniture parts may not belong to any known part semantics and some of the parts may be provided pre-assembled into bigger units. We also step away from instruction manuals illustrating the step-by-step assembling process, as teaching machines to read sequential instructions depicted with natural languages and figures is still a hard problem.

At the core of the task lie several challenges. First, some parts may have similar geometry. For example, distinguishing the geometric subtlety of chair leg bars, stretcher bars, and back bars is a difficult problem. Second, 3D geometric reasoning is essential in finding a joint global solution, where every piece fits perfectly in the puzzle. Parts follow a more rigid relationship graph which determines a unique final solution that emerges from the interactions between the geometries of the parts. Third, the image grounds and selects one single solution from all possible part combinations that might all be valid for the generative task. Thus, the problem is at heart a reconstruction task where the final assembly needs to agree to the input image. Additionally, and different from object localization tasks, *the 3D Part Assembly Task* must locate all input parts, not only posing the parts visible in the image, but also hallucinating poses for the invisible ones by leveraging learned data priors. One can think of having multiple images to expose all parts to the robot, but this reduces the generalizability to real-world scenarios, and might not be easy to achieve. Thus, we focus on solving the task of single-image and category-prior-guided pose prediction.

In this paper, we introduce a learning-based method to tackle the proposed *single-image-guided 3D part assembly* problem. Given the input image and a set of 3D parts, we first focus on 2D structural guidance by predicting an part-instance image segmentation to serve as a 2D-3D grounding for the downstream pose prediction. To enforce reasoning involving fine geometric subtleties, we have designed a context-aware 3D geometric feature to help the network reason about each part pose, conditioned on the existence of other parts, which might be of similar geometry. Building on the 2D structural guidance, we can generate a pose proposal for each visible part and leverage these predictions to help hallucinate poses for invisible parts as well. Specifically, we use a part graph network, based on edges to encode different relationships among parts, and design a two-phase message-passing mechanism to take part relationship constraints into consideration in the assembly.

To best of our knowledge, we are the first to assemble *unlabeled* 3D parts with a *single image* input. We set up a testbed of the problem on the recently released PartNet [41] dataset. We pick three furniture categories with large shape variations that require part assembly: Chair, Table and Cabinet. We compare our method with several baseline methods to demonstrate the effectiveness of our approach. We follow the PartNet official train-test splits and evaluate all model performances on the unseen test shapes. Extensive ablation experiments also demonstrate the effectiveness and necessity of the proposed modules: 2D-mask-grounding component and the 3D-message-passing reasoning component.

In summary, our contributions are:

– we formulate the task of *single-image-guided 3D part assembly*;
– we propose a two-module method, consisting of a part-instance image segmentation network and an assembly-aware part graph convolution network;
– we compare with three baseline methods and conduct ablation studies demonstrating the effectiveness of our proposed method.

## 2   Related Work

We review previous works on 3D pose estimation, single-image 3D reconstruction, as well as part-based shape modeling, and discuss how they relate to our task.

**3D Pose Estimation.** Estimating the pose of objects or object parts is a long-standing problem with a rich literature. Early in 2001, Langley *et al.* [75] attempted to utilize visual sensors and neural networks to predict the pose for robotic assembly tasks. Andy *et al.* [77] built an robotic system taking multi-view RGB-D images as the input and predicting 6D pose of objects for Amazon Picking Challenge. Recently, Litvak *et al.* [37] proposed a two-stage pose estimation procedure taking depth images as input. In the vision community, there is also a line of works studying instance-level object pose estimation for known instances [1,48,59,28,72,58,2] and category-level pose estimation [19,44,3,63,7] that can possibly deal with unseen objects from known categories. There are also works on object re-localization from scenes [76,23,61]. Different from these works, our task takes as inputs unseen parts without any semantic labels at the

test time, and requires certain part relationships and constraints to be held in order to assemble a plausible and physically stable 3D shape.

**Single-Image 3D Reconstruction.** There are previous works of reconstructing 3D shape from a single image with the representations of voxel grids [10,57,67,49], point clouds [14,34,22], meshes [65,69], parametric surfaces [18], and implicit functions [8,39,45,51,74]. While one can consider employing such 2D-to-3D lifting techniques as a prior step in our assembly process so that the given parts can be matched to the predicted 3D shape, it can misguide the assembly in multiple ways. For instance, the 3D prediction can be inaccurate, and even some small geometric differences can be crucial for part pose prediction. Also, the occluded area can be hallucinated in different ways. In our case, the set of parts that should compose the object is given, and thus the poses of occluded parts can be more precisely specified. Given these, we do not leverage 3D shape generation techniques and directly predict the part poses from the input 2D image.

**Part-Based Shape Modeling.** 3D shapes have compositional part structures. Chaudhuri *et al.* [5], Kalogerakis *et al.* [26] and Jaiswal *et al.* [24] introduced frameworks learning probabilistic graphical models that describe pairwise relationships of parts. Chaudhuri and Koltun [6], Sung *et al.* [56] and Sung *et al.* [55] predict the compatibility between a part and a *partial* object for sequential shape synthesis by parts. Dubrovina *et al.* [13], PAGENet [32] and CompoNet [52] take the set of parts as the input and generates the shape of assembled parts. Different from these works that usually assume known part semantics or a part database, our task takes a set of unseen parts during the test time and we do not assume any provided part semantic labels.

GRASS [33], Im2Struct [43] and StructureNet [40] learns to generate box-abstracted shape hierarchical structures. SAGNet [71] and SDM-Net [17] learn the pairwise relationship among parts that are subsequently integrated into a latent representation of the global shape. G2LGAN [62] autoencodes the shape of an entire object with per-point part labels, and a subsequent network in the decoding refines the geometry of each part. PQ-Net [70] represents a shape as a sequence of parts and generates each part at every step of the iterative decoding process. All of these works are relevant but different from ours in that we obtain the final geometry of the object not by directly decoding the latent code into part geometry but by predicting the poses of the given parts and explicitly assembling them. There are also works studying partial-to-full shape matching [35,36,12]. Unlike these works, we use a single image as the guidance, instead of a 3D model.

## 3   Problem Formulation

We define the task of *single-image-guided 3D part assembly*: given a single RGB image $I$ of size $m \times m$ depicting a 3D object $S$ and a set of $N$ 3D part point clouds $\mathcal{P} = \{p_1, p_2, \cdots, p_N\}$ ($\forall i, p_i \in \mathbb{R}^{d_{pc} \times 3}$), we predict a set of part poses $\{(R_i, t_i) \mid R_i \in \mathbb{R}^{3 \times 3}, t_i \in \mathbb{R}^3, i = 1, 2, \cdots, N\}$ in $SE(3)$ space. After applying the predicted rigid transformation to all the input parts $p_i$'s, the union of them reconstructs the 3D object $S$. We predict output part poses $\{(R_i, t_i) \mid i = 1, 2, \cdots, N\}$
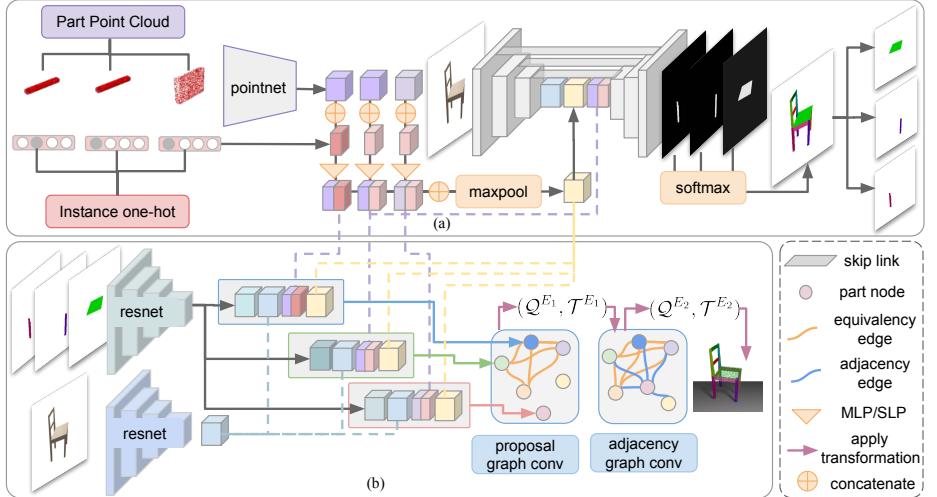
Fig. 2: **Network Architecture.** We propose a method that contains two network modules: (a) the part-instance image segmentation module, in which we predict 2D instance-level part masks on the image, and (b) the part pose prediction module, where we combine 2D mask information and 3D geometry feature for each part, push them through two phases of graph convolution, and finally predict 6D part poses.

in the camera space, following previous works [15,66]. In our paper, we use Quaternion to represent rotation and use $q_i$ and $R_i$ interchangeably.

We conduct a series of pose and scale normalization on the input part point clouds to ensure *synthetic-to-real* generalizability. We normalize each part point cloud pose $p_i \in \mathcal{P}$ to have a zero-mean center and use a local part coordinate system computed using PCA [46]. To normalize the global scale of all training and testing data, we compute Axis-Aligned-Bounding-Boxes (AABB) for all the parts and normalize them so that the longest box diagonal across all $p_i$'s of a shape has a unit length while preserving their relative scales. We cluster the normalized part point clouds $p_i$'s into sets of geometrically equivalent part classes $\mathcal{C} = \{C_1, C_2, \cdots, C_K\}$, where $C_1 = \{p_i\}_{i=1}^{N_1}$, $C_2 = \{p_i\}_{i=N_1+1}^{N_1+N_2}$, etc. For example, four legs of a chair are clustered together if their geometry is identical. This process of grouping indiscernible parts is essential to resolve the ambiguity among them in our framework. $\mathcal{C}$ is a disjoint complete set such that $C_k \cap C_l = \phi$ for every $C_k, C_l \in \mathcal{C}, k \neq l$ and $\cup_{k=1}^{K} C_k = \mathcal{P}$. We denote the representative point cloud $p_j$ for each class $C_j \in \mathcal{C}$.

## 4    Method

We propose a method for the task of *single-image-guided 3D part assembly*, which is composed of two network modules: the part-instance image segmentation module and the part pose prediction module; see Figure 2 for the overall architecture. We first extract a geometry feature of each part from the

input point cloud $p_j \in \mathcal{C}$ and generates $N$ instance-level 2D segmentation masks $\{M_i \in \{0,1\}^{m \times m} | i = 1, 2, \cdots, N\}$ on the input image ($m = 224$). Conditioned on the predicted segmentation masks, our model then leverages both the 2D mask features and the 3D geometry features to propose 6D part poses $\{(q_i, t_i) | i = 1, 2, \cdots, N\}$. We explain these two network modules in the following subsections. See supplementary for the implementation details.

## 4.1   Part-Instance Image Segmentation

To induce a faithful reconstruction of the object represented in the image, we need to learn a structural layout of the input parts from the 2D input. We predict a part instance mask $M_i \in \{0,1\}^{m \times m}$ for each part $p_i$. All part masks subject to the disjoint constraint, $i.e.$, $M_{bg} + \sum_{i=1}^{N} M_i = \mathbf{1}$, where $M_{bg}$ denotes a background mask. If a part is invisible, we simply predict an empty mask and let the second network to hallucinate a pose leveraging contextual information and learned data priors. The task difficulties are two folds. First, the network needs to distinguish between the geometric subtlety of the input part point clouds to establish a valid 2D-3D correspondence. Second, for the identical parts within each geometrically equivalent class, we need to identify separate 2D mask regions to pinpoint their exact locations. Below, we explain how our proposed method is designed to tackle the above challenges.

**Context-Aware 3D Part Features.** To enable the network to reason the delicate differences between parts, we construct the context-aware 3D conditional feature $f_{3d} \in \mathbb{R}^{2F_2}$ ($F_2 = 256$), which is computed from three components: part geometry feature $f_{geo} \in \mathbb{R}^{F_2}$, instance one-hot vector $f_{ins} \in \mathbb{R}^{P_{max}}$ ($P_{max} = 20$), and a global part contextual feature $f_{global} \in \mathbb{R}^{F_2}$. We use PointNet [47] to extract a global geometry feature $f_{geo}$ for each part point cloud $p_i$. If a part $p_j$ has multiple instances $k_j > 1$ within a geometrically equivalent class $\mathcal{C}_j$ ($e.g.$ four chair legs), we introduce an additional instance one-hot vector $f_{ins}$ to tell them apart. For part which has only one instance, we use an one-hot vector with the first element to be 1. For contextual awareness, we extract a global feature $f_{global}$ over all the input part point clouds, to facilitate the network to distinguish between similar but not equivalent part geometries ($e.g.$ a short bar or a long bar). Precisely, we first compute $f_{geo}$ and $f_{ins}$ for every part, then compute $f_{local} = SLP_1([f_{geo}; f_{ins}]) \in \mathbb{R}^{F_2}$ to obtain per-part local feature, where SLP is short for Single-Layer Perception. We aggregate over all part local features via a max-pooling symmetric function to compute the global contextual feature $f_{global} = SLP_2(MAX_{i=1,2,\cdots,N}(f_{i,local}))$. Finally, we define $f_{3d} = [f_{local}; f_{global}] \in \mathbb{R}^{2F_2}$ to be the context-aware 3D per-part feature.

**Conditional U-Net Segmentation.** We use a conditional U-Net [50] for the part-instance segmentation task. Preserving the standard U-Net CNN architecture, our encoder takes an 3-channel RGB image as input and produce a bottleneck feature map $f_{2d} \in \mathbb{R}^{F_1 \times 7 \times 7}$ ($F_1 = 512$). Concatenating the image feature $f_{2d}$ with our context-aware 3D part conditional feature $f_{3d}$, we obtain $f_{2d+3d} = [f_{2d}, f_{3d}] \in \mathbb{R}^{(F_1+2F_2) \times 7 \times 7}$, where we duplicate $f_{3d}$ along the spatial dimensions for $7 \times 7$ times. The decoder takes the conditional bottleneck feature

$f_{2d+3d}$ and decodes a part mask $M_i$ for evert input part $p_i$. We keep skip links as introduced in the original U-Net paper between encoder and decoder layers. To satisfy the non-overlapping constraint, we add a SoftMax layer across all predicted masks, augmented with a background mask $M_{bg} \in \{0, 1\}^{(m \times m)}$.

## 4.2   Part Pose Prediction

With the 2D grounding masks produced by the part-instance image segmentation module, we predict a 6D part pose $(R_i, t_i)$ for every input part $p_i \in \mathcal{P}$ using the part pose prediction module. We predict a unit Quaternion vector $q_i$ that corresponds to a 3D rotation $R_i$ and a translation vector $t_i$ denoting the part center position in the camera space.

Different from object pose estimation, the task of part assembly requires a joint prediction of all part poses. Part pose predictions should not be independent with each other, as part poses follow a set of more rigid relationships, such as symmetry and parallelism. For a valid assembly, parts must be in contact with adjacent parts. The rich part relationships restrict the solution space for each part pose. We leverage a two-phase graph convolutional neural network to address the joint communication of part poses for the task of part assembly.

**Mask-Conditioned Part Features.** We consider three sources of features for each part: 2D image feature $f_{img} \in \mathbb{R}^{F_3}$, 2D mask feature $f_{mask} \in \mathbb{R}^{F_3}$ ($F_3 = 512$), context-aware 3D part feature $f_{3d} \in \mathbb{R}^{2F_2}$. We use a ResNet-18 [20] pretrained on ImageNet [11] to extract 2D image feature $f_{img}$. We use a separate ResNet-18 that takes the 1-channel binary mask as input and extracts a 2D mask feature $f_{mask}$, where masks for invisible parts are predicted as empty. Then, finally, we propagate the 3D context-aware part feature $f_{3d}$ introduced in the Sec. 4.1 that encodes 3D part geometry information along with its global context.

**Two-Phase Graph Convolution.** We create a part graph $\mathcal{G} = (V, E)$, treating every part as a node and propose a two-phase of graph convolution to predict the pose of each part. We first describe how we construct the edges in each phase, and then introduce our assembly-oriented graph convolution operations.

During the first phase, we draw pairwise edges among all parts $p_i$ in every geometrically equivalent part classes $C_j$ and perform graph convolution over $\mathcal{G}^1 = (V, E^1)$, where

$$E^1 = \{(p_{i_1}, p_{i_2}) | \forall p_{i_1}, p_{i_2} \in C_j, i_1 \neq i_2, \forall C_j \in \mathcal{C}\} . \qquad (1)$$

Edges in $E^1$ allow message passing among geometrically identical parts that are likely to have certain spatial relationships or constraints (*e.g.* four legs of a chair have two orthogonal reflection planes). After the first-phase graph convolution, each node $p_i$ has an updated node feature. The updated node feature is then decoded as an 6D pose $(R_i, t_i)$ for each part. The predicted part poses produce an initial assembled shape.

We leverage a second phase of graph convolution to refine the predicted part poses. Besides the edges in $E^1$, we draw a new set of edges $E^2$ by finding top-5 nearest neighbors for each part based upon the initial assembly and define

$\mathcal{G}^2 = (V, E^1 \cup E^2)$. The intuition here is that once we have an initial part assembly, we are able to connect the adjacent parts so that they learn to attach to each other with certain joint constraints.

We implement the graph convolution as two iterations of message passing [73,68,40]. Given a part graph $\mathcal{G} = (V, E)$ with initial node features $f^0$ and edge features $e^0$, each iteration of message passing starts from computing edge features

$$e^{t+1}_{(p_{i_1}, p_{i_2})} = SLP_g \left( [f^t_{i_1}; f^t_{i_2}; e^t_{(p_{i_1}, p_{i_2})}] \right), t \in \{0, 1\}. \tag{2}$$

where we do not use $e^0$ during the first phase of graph convolution, and define $e^0_{(p_{i_1}, p_{i_2})} = 0$ if $(p_{i_1}, p_{i_2}) \in E^1$ and $e^0_{(p_{i_1}, p_{i_2})} = 1$ if $(p_{i_1}, p_{i_2}) \in E^2$ for the second phase. Then, we perform average-pooling over all edge features that are connected to a node and obtain the updated node feature

$$f^{t+1}_i = \frac{1}{|\{u \mid (p_i, p_u) \in E\}|} \sum_{(p_i, p_u) \in E} e^{t+1}_{(p_i, p_u)}, t \in \{0, 1\}. \tag{3}$$

We define $f^{t+1}_i = f^t_i$ if there is no edge drawn from node $i$. We define the final node features to be $f_i = [f^0_i; f^1_i; f^2_i]$ for each phase of graph convolution.

Respectively, we denote the final node feature of first phase and second phase graph convolution to be $^1f_i$ and $^2f_i$ for a part $p_i$.

**Part Pose Decoding.** After gathering the node features after conducting the two-phase graph convolution operations as $^1f_i$ and $^2f_i, i \in \{1, 2, \cdots, N\}$, we use a Multiple-Layer Perception (MLP) to decode part poses at each phase.

$$^sq_i, {}^st_i = MLP_{PoseDec} \left( {}^sf_i \right), s \in \{1, 2\}, i \in \{1, 2, \cdots, N\}. \tag{4}$$

To ensure the output of unit Quaternion prediction, we normalize the output vector length so that $\|^sq_i\|_2 = 1$.

## 4.3   Training and Losses

We first train the part-instance image segmentation module until its convergence and then train the part pose prediction module. Empirically, we find that having a good mask prediction is necessary before training for the part poses.

**Loss for Part-Instance Image Segmentation.** We adapt the negative *soft-iou* loss from [38] to supervise the training of the part-instance image segmentation module. We perform Hungarian matching [29] within each geometrically equivalent class to guarantee that the loss is invariant to the order of part poses in ground-truth and prediction. The loss is defined as

$$\mathcal{L}_{mask_i} = -\frac{\sum_{u,v \in [m,m]} \hat{M}^{(u,v)}_i \cdot M^{(u,v)}_{\mathcal{M}(i)}}{\sum_{u,v \in [m,m]} \left( \hat{M}^{(u,v)}_{\mathcal{M}(i)} + M^{(u,v)}_i - \hat{M}^{(u,v)}_{\mathcal{M}(i)} \cdot M^{(u,v)}_i \right)}. \tag{5}$$

where $M_i \in \{0, 1\}^{[m,m]}$ and $\hat{M}_{\mathcal{M}(i)} \in [0, 1]^{[m,m]}$ denote the ground truth and the matched predicted mask. $\mathcal{M}$ refers to the matching results that match ground-truth part indices to the predicted ones. $[m, m]$ includes all 2D index $(u, v)$'s on a $224 \times 224$ image plane.

**Losses for Part Pose Prediction.** For the pose prediction module, we design an order-invariant loss by conducting Hungarian matching within each geometry-equivalent classes $C_i \in \mathcal{C}$. Additionally, we observe that separating supervision loss for translation and rotation helps stabilize training. We use the following training loss for the pose prediction module.

$$\mathcal{L}_{pose} = \sum_{i=1}^{N}(\lambda_1 \times \mathcal{L}_T + \lambda_2 \times \mathcal{L}_C + \lambda_3 \times \mathcal{L}_E) + \lambda_4 \times \mathcal{L}_W \qquad (6)$$

We use the $L_2$ Euclidean distance to measure the difference between the 3D translation prediction and ground truth translation for each part. We denote $\mathcal{M}$ as the matching results.

$$\mathcal{L}_{T_i} = \|\hat{t}_{\mathcal{M}(i)} - t_i\|_2, \forall i \in \{1, 2, \cdots, N\}. \qquad (7)$$

where $\hat{t}_{\mathcal{M}(i)}$ and $t_i$ denote the matched predicted translation and the ground truth 3D translation. We use weight parameter of $\lambda_1 = 1$ in training.

We use two losses for rotation prediction: Chamfer distance [15] $\mathcal{L}_C$ and $L2$ distance $\mathcal{L}_E$. Because many parts have symmetric geometry (*e.g.* bars and boards) which results in multiple rotation solutions, we use Chamfer distance as the primary supervising loss to address such pose ambiguity. Given the point cloud of part $p_i$, the ground truth rotation $R_i$, and the matched predicted rotation $\hat{R}_{\mathcal{M}(i)}$, the Chamfer distance loss is defined as

$$\mathcal{L}_{C_i} = \frac{1}{d_{pc}} \sum_{x \in \hat{R}_{\mathcal{M}(i)}(p_i)} \min_{y \in R_i(p_i)} \|x - y\|_2^2 + \frac{1}{d_{pc}} \sum_{y \in R_i(p_i)} \min_{x \in \hat{R}_{\mathcal{M}(i)}(p_i)} \|x - y\|_2^2, \quad (8)$$

where $R_i(p_i)$ and $\hat{R}_{\mathcal{M}(i)}(p_i)$ denote the rotated part point clouds using $R_i$ and $\hat{R}_{\mathcal{M}(i)}$ respectively. We use $\lambda_2 = 20$ for the Chamfer loss. Some parts may be not *perfectly* symmetric (*e.g.* one bar that has small but noticeable different geometry at two ends), using Chamfer distance by itself in this case would make the network fall into local minima. We encourage the network to correct this situation by penalizing the $L_2$ distance between the matched predicted rotated point cloud and the ground truth rotated point cloud in Euclidean distance.

$$\mathcal{L}_{E_i} = \frac{1}{d_{pc}} \left\| \hat{R}_{\mathcal{M}(i)}(p_i) - R_i(p_i) \right\|_F^2, \qquad (9)$$

where $\|\cdot\|_F$ denotes the Frobenius norm, $d_{pc} = 1000$ is the number of points per part. Note that $\mathcal{L}_{E_i}$ on its own is not sufficient in cases when the parts are completely symmetric. Thus, we add the $\mathcal{L}_E$ loss as a regularizing term with a smaller weight of $\lambda_3 = 1$. We conducted an ablation experiment demonstrating the $\mathcal{L}_E$ loss contributes to correcting rotation for some parts.

Finally, we compute a shape holistic Chamfer distance as the predicted assembly should be close to the ground truth Chamfer distance.

$$\mathcal{L}_W = \frac{1}{N \cdot d_{pc}} \sum_{x \in \hat{S}} \min_{y \in S} \|x - y\|_2^2 + \frac{1}{N \cdot d_{pc}} \sum_{y \in S} \min_{x \in \hat{S}} \|x - y\|_2^2, \qquad (10)$$
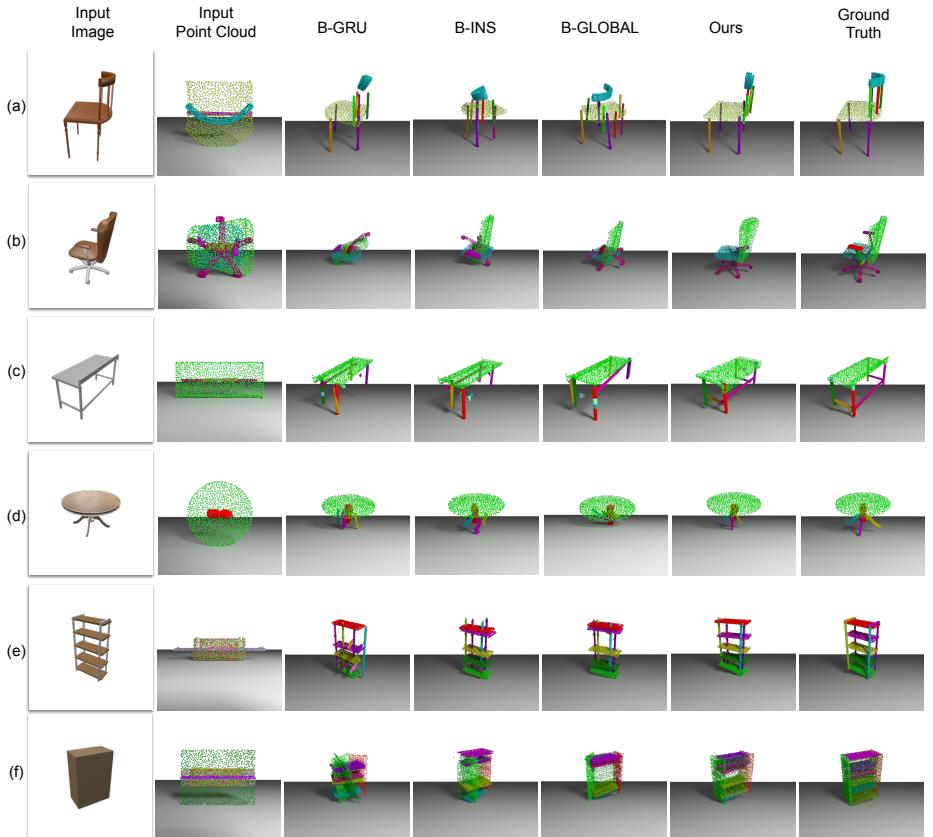
Fig. 3: **Qualitative Results.** We show six examples of the three furniture categories, each in two different modalities. The upper and lower rows correspond to modality Level-3 and Level mixed respectively.

where $\hat{S} = \cup_{i=1}^{N}(\hat{R}_{\mathcal{M}(i)}(p_i) + \hat{t}_i)$ denotes the predicted assembled shape point cloud and $S = \cup_{i=1}^{N}(R_i(p_i) + t_i)$ denotes the ground truth shape point cloud. This loss encourages the holistic shape appearance and the part relationships to be close to the ground-truth. We use $\lambda_4 = 1$.

## 5    Experiments

In this section, we set up the testbed for the proposed *single-image-guided 3D part assembly* problem on the PartNet [41] dataset. To validate the proposed approach, we compare against three baseline methods. Both qualitative and quantitative results demonstrate the effectiveness of our method.

### 5.1    Dataset

Recently, Mo et. al. [41] proposed the PartNet dataset, which is the largest 3D object dataset with fine-grained and hierarchical part annotation. Every

PartNet object is provided with a ground-truth hierarchical part instance-level semantic segmentation, from coarse to fine-grained levels , which provides a good complexity of parts. In our work, we use the three largest furniture categories that the requires real-world assembly: Chair, Table and Cabinet. We follow the official PartNet train/validation/test split (roughly $70\% : 10\% : 20\%$) and filter out the shapes with more than 20 parts.

For each object category, we create two data modalities: *Level-3* and *Level-mixed*. The *Level-3* corresponds to the most fine-grained PartNet segmentation. While we do not assume known part semantics, an algorithm can implicitly learn the semantic priors dealing with the *Level-3* data, which is undesired in our goal of generalizing to real-life assembly settings, as it is unrealistic to assume taht IKEA furnitures also follow the PartNet same semantics. To enforce the network to reason with part geometries, we created an additional category modality, *Level-mixed*, which contains part segmentation at all levels in the PartNet hierarchy. Specifically, for each shape, we traverse every path of the ground-truth part hierarchy and stop at any level randomly. We have 3736 chairs, 2431 tables, 704 cabinets in *Level-3* and 4664 chairs, 5987 tables, 888 cabinets in *Level-mixed*.

For the input image, we render a set of $224 \times 224$ images the PartNet models with ShapeNet textures [4]. We then compute the world-to-camera matrix accordingly and obtain the ground-truth 3D object position in the camera space, which is used for supervising part-instance segmentation supervision. For the input point cloud, we use Furthest Point Sampling (FPS) to sample $d_{pc} = 1000$ points over the each part mesh. We then normalize them following the descriptions in Sec. 3. After parts are normalized, we detect geometrically equivalent classes of parts by first filtering out parts comparing dimensions of AABB under a threshold of 0.1. We further process the remaining parts computing all possible pairwise part Chamfer distance normalized by their average diagonal length under a hand-picked threshold of 0.02.

## 5.2   Evaluation Metric

To evaluate the part assembly performance, we use two metrics: *part accuracy* and *shape Chamfer distance*. The community of object pose estimation usually uses metrics such as 5-degree-5-cm. However, fine-grained part segments usually show abundant pose ambiguity. For example, a chair leg may be simply a cylinder which has a full rotational and reflective symmetry. Thus, we introduce the *part accuracy* metric that leverages Chamfer distance between the part point clouds after applying the predicted part pose and the ground truth pose to address such ambiguity. Following previously defined notation in Section 4.3, we define the Part Accuracy Score (PA) as follows and set a threshold of $\tau = 0.1$.

$$PA = \frac{1}{N} \sum_{i=1}^{N} \mathbb{1} \left( \left\| (\hat{R}_{\mathcal{M}(i)}(p_i) + \hat{t}_i) - (R_i(p_i) + t_i) \right\|_{chamfer} < \tau \right) \qquad (11)$$

Borrowing the evaluation metric heavily used in the community of 3D object reconstruction, we also measure the *shape Chamfer distance* from the predicted

assembled shape to the ground-truth assembly. Formally, we define the *shape Chamfer distance* metric $SC$ borrowing notations defined in Section 4.3 as follows.

$$SC(S, \hat{S}) = \frac{1}{N \cdot d_{pc}} \sum_{x \in \hat{S}} \min_{y \in S} \|x - y\|_2 + \frac{1}{N \cdot d_{pc}} \sum_{y \in S} \min_{x \in \hat{S}} \|x - y\|_2 \qquad (12)$$

## 5.3  Baseline Methods

We compare our approach to three baseline methods. Since there is no direct comparison from previous works that address the exactly same task, we try to adapt previous works on part-based shape generative modeling [70,56,40,43] to our setting and compare with them. Most of these works require known part semantics and thus perform part-aware shape generation without the input part conditions. However, in our task, there is no assumption for part semantics or part priors, and thus all methods must explicitly take the part input point clouds as input conditions. We train all three baselines with the same pose loss used in our method defined in Section 4.3.

**Sequential Pose Proposal (B-GRU)** The first baseline is a sequential model, similar to the method proposed by [70,56], instead of sequentially generating parts, we sequentially decode $k$ candidate possible poses for a given part geometry, conditioned on an image. For each input part, if there is $n$ geometrically equivalent parts , where $n \leq k$, we take the first n pose proposal generated using GRU, and conduct Hungarian matching to match with the $n$ ground truth part poses.

**Instance One-hot Pose Proposal (B-InsOneHot)** The second baseline uses MLP to directly infer pose for a given part from its geometry and the input image, similar to previous works [40,43] that output box abstraction for shapes. Here, instead of predicting a box for each part, we predict a 6D part pose $(R_j, t_j)$. We use instance one-hot features to differentiate between the equivalent part point clouds, and conduct Hungarian matching to match with the ground truth part poses regardless of the onehot encoding.

**Global Feature Model (B-Global)** The third baseline is proposed by improving upon the second baseline by adding our the context-aware 3D part feature defined in Section 4.1. Each part pose proposal not only considers the part-specific 3D feature and the 2D image feature, but also a 3D global feature obtained by

Table 1: Part Accuracy and Assembly Chamfer Distance(CD)

| Modality | Method | Part Accuracy ↑ | | | Assembly CD ↓ | | |
|---|---|---|---|---|---|---|---|
| | | Chair | Table | Cabinet | Chair | Table | Cabinet |
| Level-3 | B-GRU | 0.310 | 0.574 | 0.334 | 0.107 | 0.057 | 0.062 |
| | B-InsOnehot | 0.173 | 0.507 | 0.295 | 0.130 | 0.064 | 0.065 |
| | B-Global | 0.170 | 0.530 | 0.339 | 0.125 | 0.061 | 0.065 |
| | Ours | **0.454** | **0.716** | **0.402** | **0.067** | **0.037** | **0.050** |
| Mixed | B-GRU | 0.326 | 0.567 | 0.283 | 0.101 | 0.070 | 0.066 |
| | B-InsOnehot | 0.286 | 0.572 | 0.320 | 0.108 | 0.067 | 0.061 |
| | B-Global | 0.337 | 0.619 | 0.290 | 0.093 | 0.062 | 0.0677 |
| | Ours | **0.491** | **0.778** | **0.483** | **0.065** | **0.037** | **0.043** |

Table 2: Visible and Invisible Part Accuracies

| Modality | Method | Part Accuracy (Visible) ↑ | | | Part Accuracy (Invisible) ↑ | | |
|---|---|---|---|---|---|---|---|
| | | Chair | Table | Cabinet | Chair | Table | Cabinet |
| Level-3 | B-GRU | 0.3182 | 0.598 | 0.353 | 0.206 | 0.481 | 0.304 |
| | B-InsOnehot | 0.178 | 0.572 | 0.291 | 0.104 | 0.369 | 0.289 |
| | B-Global | 0.174 | 0.563 | 0.354 | 0.120 | 0.427 | 0.269 |
| | Ours | **0.471** | **0.753** | **0.455** | **0.270** | **0.557** | **0.358** |
| Mixed | B-GRU | 0.335 | 0.593 | 0.302 | 0.180 | 0.267 | 0.258 |
| | B-InsOnehot | 0.295 | 0.592 | 0.346 | 0.133 | 0.275 | 0.279 |
| | B-Global | 0.334 | 0.638 | 0.320 | 0.184 | 0.349 | 0.227 |
| | Ours | **0.505** | **0.803** | **0.537** | **0.262** | **0.515** | **0.360** |

aggregating the all 3D part feature then max-pool to a global 3D feature containing information of all parts. This baseline shares similar ideas to PAGENet [32] and CompoNet [52] that also compute global features to assemble each of the generated parts.

## 5.4   Results and Analysis

We compare with the three baselines and observe that our method outperforms the baseline methods both qualitatively and quantitatively using the two evaluation metrics, PA and SC. We show significant improvement for occluded part pose hallucination as Table 2 demonstrates. Qualitatively, we observe that our method can learn to infer part poses for invisible parts by (1) learning a category prior and (2) leveraging visible parts of the same geometric equivalent class. Our network can reason the stacked placement structure of cabinets as shown in the last row in Fig 3. The input image does not reveal the inner structure of the cabinet and our proposed approach learns to vertically distribute the geometrically equivalent boards to fit inside the cabinet walls, similar to the ground truth shape instance. The top row of Fig 3 demonstrates how our network learns to place the occluded back bar along the visible ones. This could be contributed to our first stage of graph convolution where we leverage visible parts to infer the pose for occluded parts in the same geometrically equivalent class.

Our method demonstrates the most faithful part pose prediction for the shape instance depicted by the input image. As shown in Fig 3 row (e), our method equally spaces the board parts vertically, which is consistent with the shape structure revealed by the input image. This is likely resulted from our part-instance image segmentation module where we explicitly predict a 2D-3D grounding, whereas the baseline methods lack such components, and we further demonstrate its effectiveness with an ablation experiments.

However, our proposed method has its limitations in dealing with unusual image views, exotic shape instance, and shapes composed of only one type of part geometry, which result in noisy mask prediction. The 2D-3D grounding error cascades to later network modules resulting in poor pose predictions. As shown in Fig 4 row (a), the image view is not very informative of the shape structure, making it difficult to leverage 3D geometric cues to find 2D-3D grounding. Additionally, this chair instance itself is foreign to Chair category. We avoided
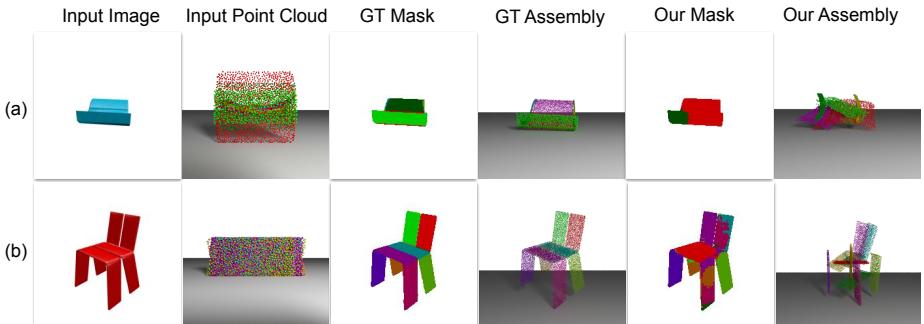
Fig. 4: **Failure Case** . We show two examples of failure cases. Case (a), the input image is not geometrically informative. Case (b), the chair has only one type of part geometry.

employing differentiable rendering because it does not help address such failure cases. Fig 4 row (b) reflects a case where a shape instance is composed of a single modality of part geometry. Geometric affinity of the board parts makes it difficult for the network to come to a determinant answer for the segmentation prediction, resulting in a sub-optimal part pose prediction. These obstacles arise from the task itself that all baselines also suffer from the same difficulties.

**Ablation Experiments** We conduct several ablation experiments on our proposed method and losses trained on PartNet Chair Level-3. Table 3 in Appendix demonstrates the effectiveness of each ablated component. The part-instance image segmentation module plays the most important role in our pipeline. Removing it results in the most significant performance decrease.

## 6    Conclusion and Future Works

We formulated a novel problem of *single-image-guided 3D part assembly* and proposed a neural-net-based pipeline for the task that leverages information from both 2D grounding and 3D geometric reasoning. We established a test bed on the PartNet dataset. Quantitative evaluation demonstrates that the proposed method achieves a significant improvement upon three baseline methods. For the future works, one can study how to leverage multiple images or 3D partial scans as inputs to achieve better results. We also do not explicitly consider the connecting junctions between parts (*e.g.* pegs and holes) in our framework, which are strong constraints for real-world robotic assembly.

## 7    Acknowledgements

# 8   Appendix

This document provides supplementary materials accompanying the main paper, including

- Ablation Experiments
- Discussion of failure cases and future works;
- More Architecture Details;
- More Qualitative Examples.

## A. Ablation Experiments

Table 3: Ablation Experiment Results

| Ablated Module | Part Accuracies ↑ | | | Assembly CD ↓ |
|---|---|---|---|---|
| | Total | Visible | Invisible | |
| w/o L2 Rotation loss | 0.426 | 0.445 | 0.207 | 0.070 |
| w/o Segmentation | 0.363 | 0.378 | 0.164 | 0.084 |
| w/o Graph Conv 1, 2 | 0.403 | 0.423 | 0.178 | 0.073 |
| w/o Graph Conv 2 | 0.434 | 0.456 | 0.239 | 0.073 |
| w/o Image Feature | 0.403 | 0.419 | 0.208 | 0.077 |
| w/o Global Feature | 0.418 | 0.437 | 0.202 | 0.072 |
| Ours - Full | **0.454** | **0.470** | **0.270** | **0.067** |

## B. More Failure Cases and Discussion

**Disconnected Parts** We notice that our prediction on very fine-grained instances sometimes results in unconnected parts. The assembly setting requires the physical constraint that each part must be in contact with another part. However, the implicit soft constraint enforced using the second stage graph graph convolution is not sufficient enough for this task. Ideally, the translation and rotation predicted for each part is only valid if they can transform the part to be in contact at the joints between relevant parts. For example, in Figure 5 we can see that the back of the chair base bars does not connect. We plan to address this problem in future works by explicitly enforcing contact between parts in a range of contact neighborhood.

**Geometric Reasoning** Additionally, though our current proposed method makes many design choices geared for geometric reasoning between fitting of parts, however, we still see some cases that the fitting between parts is not yet perfect. For example, in Figure 5, We can see that the back pad does not fit perfectly into the back frame bar. This problem need to be addressed in future work where the method design should discover some pairwise or triplet-level geometric properties that allow fitting between parts.
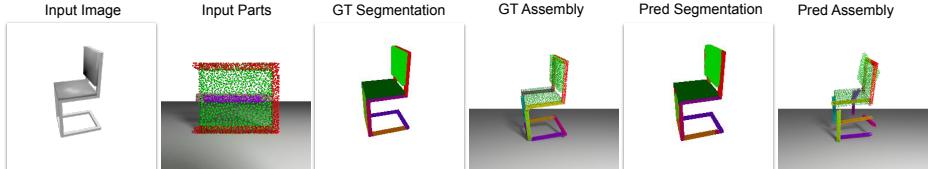
Fig. 5: **Failure Case** This figure shows that our proposed method does not well handle disconnected parts, and needs to leverage more geometric reasoning.

## C. Architecture Details

Table 4: Part-instance Segmentation Architecture.

| layer | configuration |
|---|---|
| | UNet Encoding |
| 1 | Conv2D (3, 32, 3, 1, 1), ReLU, BN, Conv2D (32, 32, 3, 1, 1), ReLU, BN, |
| 2 | Conv2D (32, 64, 3, 1, 1), ReLU, BN, Conv2D (64, 64, 3, 1, 1), ReLU, BN, |
| 3 | Conv2D (64, 128, 3, 1, 1), ReLU, BN, Conv2D (128, 128, 3, 1, 1), ReLU, BN, |
| 4 | Conv2D (128, 256, 3, 1, 1), ReLU, BN, Conv2D (256, 256, 3, 1, 1), ReLU, BN, |
| 5 | Conv2D (256, 512, 3, 1, 1), ReLU, BN, Conv2D (512, 512, 3, 1, 1), ReLU, BN, |
| | UNet Decoding |
| 1 | ConvTranspose2D(1301, 256, 2, 2) |
| 2 | ConvTranspose2D(256, 128, 2, 2) |
| 3 | ConvTranspose2D(128, 64 , 2, 2) |
| 4 | ConvTranspose2D(64, 32, 2, 2) |
| 5 | ConvTranspose2D(32, 1, 1, 1) |
| | PointNet |
| 1 | Conv1D (3, 64, 1, 1), BN, ReLU |
| 2 | Conv1D (64, 64, 1, 1), BN, ReLU |
| 3 | Conv1D (64, 64, 1, 1), BN, ReLU |
| 4 | Conv1D (64, 128, 1, 1), BN, ReLU |
| 5 | Conv1D (128, 512, 1, 1), BN, ReLU |
| | SLP1 |
| 1 | FC (512, 256), ReLU, MaxPool1D |
| | SLP2 |
| 1 | FC (256, 256), ReLU |

Table 5: Pose Prediction Architecture.

| layer | configuration |
|-------|---------------|
| SLP 3 | |
| 1 | FC(1301, 256), ReLU |
| Pose Decoder 2 | |
| 1 | FC (1301, 256), ReLU |
| 2 | FC(256, 3) |
| 3 | FC(256, 4) |
| SLP 4 | |
| 1 | FC(1031, 256), ReLU |
| Pose Decoder 2 | |
| 1 | FC (1031, 256), ReLU |
| 2 | FC(256, 3) |
| 3 | FC(256, 4) |

## D. More Qualitative Results

Fig. 6: **Qualitative Results for the Chair Category.** The top 5 rows show the results of Chair Level-3, and the bottom 5 rows contains the results of Chair Level-mixed.
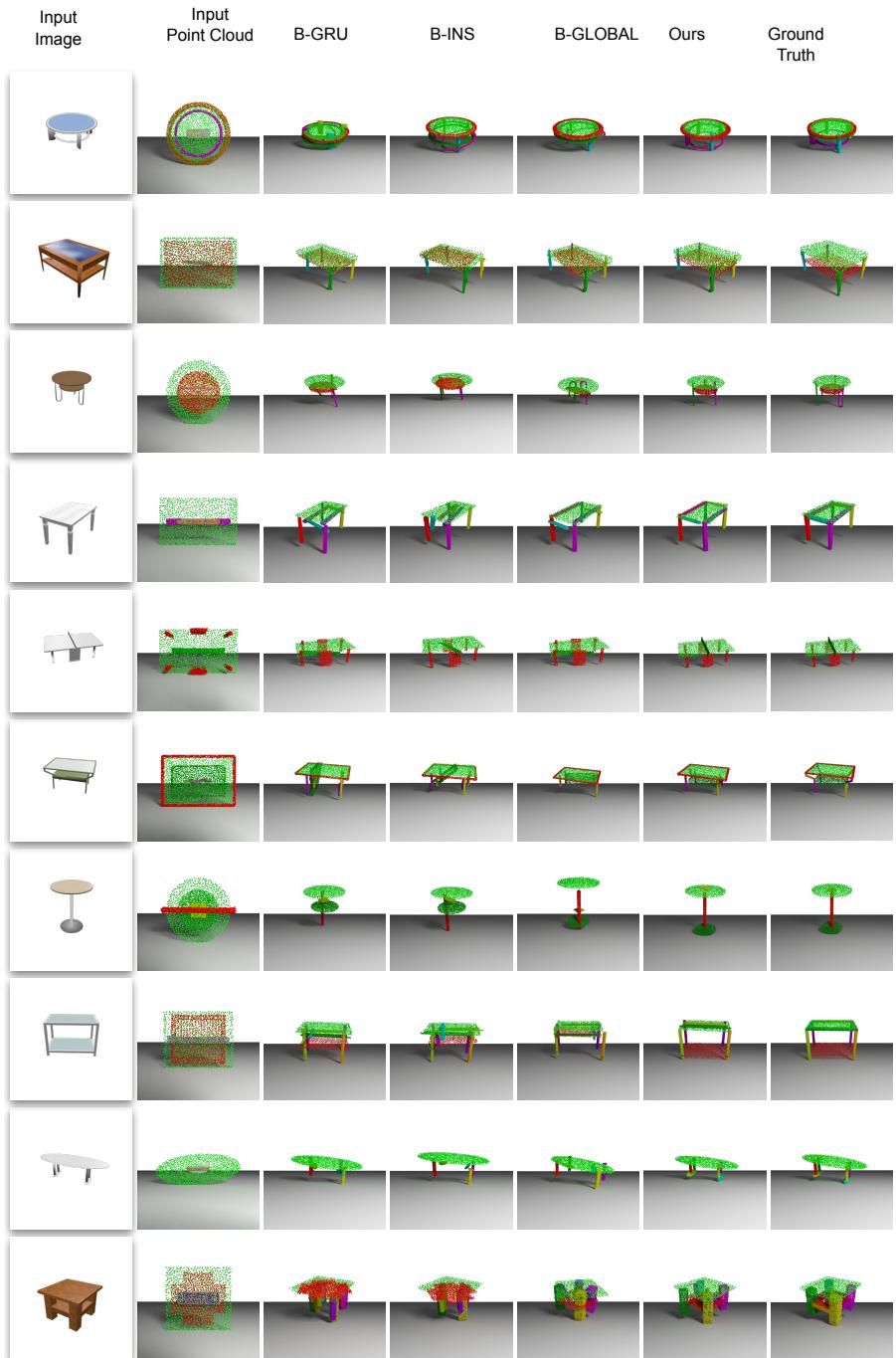
Fig. 7: **Qualitative Results for the Table Category.** The top 5 rows show the results of Table Level-3, and the bottom 5 rows contains the results of Table Level-mixed.
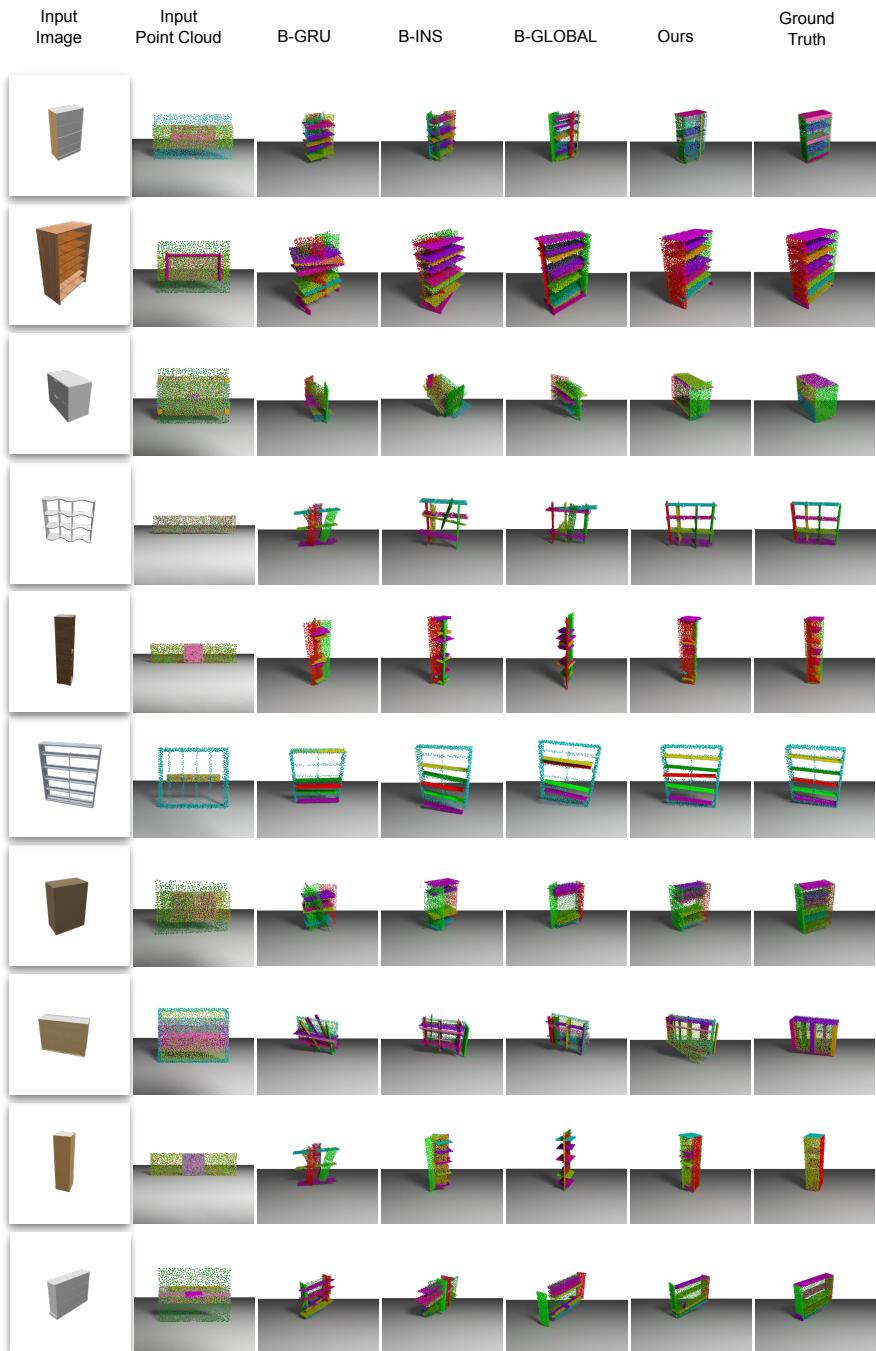
Fig. 8: **Qualitative Results for the Cabinet Category.** The top 5 rows show the results of Cabinet Level-3, and the bottom 5 rows contains the results of Cabinet Level-mixed.

# References

1. Brachmann, E., Krull, A., Michel, F., Gumhold, S., Shotton, J., Rother, C.: Learning 6d object pose estimation using 3d object coordinates. In: European conference on computer vision. pp. 536–551. Springer (2014)
2. Brachmann, E., Michel, F., Krull, A., Ying Yang, M., Gumhold, S., et al.: Uncertainty-driven 6d pose estimation of objects and scenes from a single rgb image. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3364–3372 (2016)
3. Braun, M., Rao, Q., Wang, Y., Flohr, F.: Pose-rcnn: Joint object detection and pose estimation using 3d object proposals. In: 2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC). pp. 1546–1551. IEEE (2016)
4. Chang, A.X., Funkhouser, T.A., Guibas, L.J., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., Xiao, J., Yi, L., Yu, F.: Shapenet: An information-rich 3d model repository (2015)
5. Chaudhuri, S., Kalogerakis, E., Guibas, L.J., Koltun, V.: Probabilistic reasoning for assembly-based 3d modeling. In: ACM SIGGRAPH (2011)
6. Chaudhuri, S., Koltun, V.: Data-driven suggestions for creativity support in 3d modeling. In: ACM SIGGRAPH Asia (2010)
7. Chen, D., Li, J., Xu, K.: Learning canonical shape space for category-level 6d object pose and size estimation. arXiv preprint arXiv:2001.09322 (2020)
8. Chen, Z., Zhang, H.: Learning implicit fields for generative shape modeling. In: CVPR (2019)
9. Choi, C., Taguchi, Y., Tuzel, O., Liu, M., Ramalingam, S.: Voting-based pose estimation for robotic assembly using a 3d sensor. In: 2012 IEEE International Conference on Robotics and Automation. pp. 1724–1731 (May 2012). https://doi.org/10.1109/ICRA.2012.6225371
10. Choy, C.B., Xu, D., Gwak, J., Chen, K., Savarese, S.: 3D-R2N2: A unified approach for single and multi-view 3d object reconstruction. In: ECCV (2016)
11. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
12. Domokos, C., Kato, Z.: Realigning 2d and 3d object fragments without correspondences. IEEE transactions on pattern analysis and machine intelligence **38**(1), 195–202 (2015)
13. Dubrovina, A., Xia, F., Achlioptas, P., Shalah, M., Groscot, R., Guibas, L.J.: Composite shape modeling via latent space factorization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 8140–8149 (2019)
14. Fan, H., Su, H., Guibas, L.: A point set generation network for 3d object reconstruction from a single image. In: CVPR (2017)
15. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 605–613 (2017)
16. Feng, C., Xiao, Y., Willette, A., McGee, W., Kamat, V.R.: Vision guided autonomous robotic assembly and as-built scanning on unstructured construction sites. Automation in Construction **59**, 128–138 (2015)
17. Gao, L., Yang, J., Wu, T., Yuan, Y.J., Fu, H., Lai, Y.K., Zhang, H.: SDM-NET: Deep generative network for structured deformable mesh. In: ACM SIGGRAPH Asia (2019)

18. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: AtlasNet: A papier-mch approach to learning 3d surface generation. In: CVPR (2019)
19. Gupta, S., Arbeláez, P., Girshick, R., Malik, J.: Aligning 3d models to rgb-d images of cluttered scenes. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4731–4740 (2015)
20. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2015)
21. Hutchinson, S.A., Kak, A.C.: Extending the classical ai planning paradigm to robotic assembly planning. In: Proceedings., IEEE International Conference on Robotics and Automation. pp. 182–189. IEEE (1990)
22. Insafutdinov, E., Dosovitskiy, A.: Unsupervised learning of shape and pose with differentiable point clouds. In: NeurIPS (2018)
23. Izadinia, H., Shan, Q., Seitz, S.M.: Im2cad. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5134–5143 (2017)
24. Jaiswal, P., Huang, J., Rai, R.: Assembly-based conceptual 3d modeling with unlabeled components using probabilistic factor graph. Computer-Aided Design (2016)
25. Jiménez, P.: Survey on assembly sequencing: a combinatorial and geometrical perspective. Journal of Intelligent Manufacturing **24**(2), 235–250 (2013)
26. Kalogerakis, E., Chaudhuri, S., Koller, D., Koltun, V.: A Probabilistic Model of Component-Based Shape Synthesis. In: ACM SIGGRAPH (2012)
27. Kaufman, S.G., Wilson, R.H., Jones, R.E., Calton, T.L., Ames, A.L.: The archimedes 2 mechanical assembly planning system. In: Proceedings of IEEE international conference on Robotics and Automation. vol. 4, pp. 3361–3368. IEEE (1996)
28. Kehl, W., Manhardt, F., Tombari, F., Ilic, S., Navab, N.: Ssd-6d: Making rgb-based 3d detection and 6d pose estimation great again. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1521–1529 (2017)
29. Kuhn, H.W.: The hungarian method for the assignment problem. Naval Research Logistics Quarterly (1955)
30. Langley, C.S., D'Eleuterio, G.M.T.: Neural network-based pose estimation for fixtureless assembly. In: Proceedings 2001 IEEE International Symposium on Computational Intelligence in Robotics and Automation (Cat. No.01EX515). pp. 248–253 (July 2001). https://doi.org/10.1109/CIRA.2001.1013205
31. Levitin, G., Rubinovitz, J., Shnits, B.: A genetic algorithm for robotic assembly line balancing. European Journal of Operational Research **168**(3), 811–825 (2006)
32. Li, J., Niu, C., Xu, K.: Learning part generation and assembly for structure-aware shape synthesis. In: AAAI (2020)
33. Li, J., Xu, K., Chaudhuri, S., Yumer, E., Zhang, H., Guibas, L.: GRASS: Generative recursive autoencoders for shape structures. In: ACM SIGGRAPH (2019)
34. Lin, C.H., Kong, C., Lucey, S.: Learning efficient point cloud generation for dense 3d object reconstruction. In: AAAI (2018)
35. Litany, O., Bronstein, A.M., Bronstein, M.M.: Putting the pieces together: Regularized multi-part shape matching. In: European Conference on Computer Vision. pp. 1–11. Springer (2012)
36. Litany, O., Rodolà, E., Bronstein, A.M., Bronstein, M.M., Cremers, D.: Non-rigid puzzles. In: Computer Graphics Forum. vol. 35, pp. 135–143. Wiley Online Library (2016)
37. Litvak, Y., Biess, A., Bar-Hillel, A.: Learning pose estimation for high-precision robotic assembly using simulated depth images. In: 2019 International Conference on Robotics and Automation (ICRA). pp. 3521–3527 (May 2019). https://doi.org/10.1109/ICRA.2019.8794226

38. Máttyus, G., Luo, W., Urtasun, R.: Deeproadmapper: Extracting road topology from aerial images. 2017 IEEE International Conference on Computer Vision (ICCV) pp. 3458–3466 (2017)

39. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: CVPR (2019)

40. Mo, K., Guerrero, P., Yi, L., Su, H., Wonka, P., Mitra, N., Guibas, L.: StructureNet: Hierarchical graph networks for 3d shape generation. In: ACM SIGGRAPH Asia (2019)

41. Mo, K., Zhu, S., Chang, A.X., Yi, L., Tripathi, S., Guibas, L.J., Su, H.: PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding. In: CVPR (2019)

42. Nelson, B., Papanikolopoulos, N., Khosla, P.: Visual servoing for robotic assembly. In: Visual Servoing: Real-Time Control of Robot Manipulators Based on Visual Sensory Feedback, pp. 139–164. World Scientific (1993)

43. Niu, C., Li, J., Xu, K.: Im2Struct: Recovering 3d shape structure from a single rgb image. In: CVPR (2018)

44. Papon, J., Schoeler, M.: Semantic pose using deep networks trained on synthetic rgb-d. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 774–782 (2015)

45. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: DeepSDF: Learning continuous signed distance functions for shape representation. In: CVPR (2019)

46. Pearson, K.: On lines and planes of closest fit to systems of points in space. Philosophical Magazine **2**, 559–572 (1901)

47. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: PointNet: Deep learning on point sets for 3d classification and segmentation. In: CVPR (2017)

48. Rad, M., Lepetit, V.: Bb8: A scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3828–3836 (2017)

49. Richter, S.R., Roth, S.: Matryoshka networks: Predicting 3D geometry via nested shape layers. In: CVPR (2018)

50. Ronneberger, O., Fischer, P., Brox, T.: U-Net: Convolutional networks for biomedical image segmentation. In: MICCAI (2015)

51. Saito, S., Huang, Z., Natsume, R., Morishima, S., Kanazawa, A., Li, H.: PIFu: Pixel-aligned implicit function for high-resolution clothed human digitization. In: ICCV (2019)

52. Schor, N., Katzir, O., Zhang, H., Cohen-Or, D.: CompoNet: Learning to generate the unseen by part synthesis and composition. In: ICCV (2019)

53. Shao, L., Migimatsu, T., Bohg, J.: Learning to scaffold the development of robotic manipulation skills. arXiv preprint arXiv:1911.00969 (2019)

54. Suárez-Ruiz, F., Zhou, X., Pham, Q.C.: Can robots assemble an ikea chair? Science Robotics **3**(17) (2018). https://doi.org/10.1126/scirobotics.aat6385, https://robotics.sciencemag.org/content/3/17/eaat6385

55. Sung, M., Dubrovina, A., Kim, V.G., Guibas, L.: Learning fuzzy set representations of partial shapes on dual embedding spaces. In: Eurographics Symposium on Geometry Processing (2018)

56. Sung, M., Su, H., Kim, V.G., Chaudhuri, S., Guibas, L.: ComplementMe: Weakly-supervised component suggestions for 3d modeling. In: ACM SIGGRAPH Asia (2017)

57. Tatarchenko, M., Dosovitskiy, A., Brox, T.: Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In: ICCV (2017)

58. Tejani, A., Kouskouridas, R., Doumanoglou, A., Tang, D., Kim, T.K.: Latent-class hough forests for 6 dof object pose estimation. IEEE transactions on pattern analysis and machine intelligence **40**(1), 119–132 (2017)

59. Tekin, B., Sinha, S.N., Fua, P.: Real-time seamless single shot 6d object pose prediction. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 292–301 (2018)

60. Thorsley, M., Okouneva, G., Karpynczyk, J.: Stereo vision algorithm for robotic assembly operations. In: First Canadian Conference on Computer and Robot Vision, 2004. Proceedings. pp. 361–366. IEEE (2004)

61. Wald, J., Avetisyan, A., Navab, N., Tombari, F., Nießner, M.: Rio: 3d object instance re-localization in changing indoor environments. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 7658–7667 (2019)

62. Wang, H., Schor, N., Hu, R., Huang, H., Cohen-Or, D., Huang, H.: Global-to-local generative model for 3d shapes. In: ACM SIGGRAPH Asia (2018)

63. Wang, H., Sridhar, S., Huang, J., Valentin, J., Song, S., Guibas, L.J.: Normalized object coordinate space for category-level 6d object pose and size estimation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2642–2651 (2019)

64. Wang, L., Schmidt, B., Givehchi, M., Adamson, G.: Robotic assembly planning and control with enhanced adaptability through function blocks. The International Journal of Advanced Manufacturing Technology **77**(1), 705–715 (2015)

65. Wang, N., Zhang, Y., Li, Z., Fu, Y., Liu, W., Jiang, Y.G.: Pixel2mesh: Generating 3d mesh models from single rgb images. In: ECCV (2018)

66. Wang, N., Zhang, Y., Li, Z., Fu, Y., Liu, W., Jiang, Y.G.: Pixel2mesh: Generating 3d mesh models from single rgb images. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 52–67 (2018)

67. Wang, P.S., Liu, Y., Guo, Y.X., Sun, C.Y., Tong, X.: Adaptive O-CNN: A Patch-based Deep Representation of 3D Shapes. In: ACM SIGGRAPH Asia (2018)

68. Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M.: Dynamic graph cnn for learning on point clouds. ACM Transactions on Graphics (TOG) **38**(5), 1–12 (2019)

69. Wen, C., Zhang, Y., Li, Z., Fu, Y.: Pixel2Mesh++: Multi-view 3d mesh generation via deformation. In: ICCV (2019)

70. Wu, R., Zhuang, Y., Xu, K., Zhang, H., Chen, B.: PQ-NET: A generative part seq2seq network for 3d shapes (2019)

71. Wu, Z., Wang, X., Lin, D., Lischinski, D., Cohen-Or, D., Huang, H.: SAGNet: Structure-aware generative network for 3d-shape modeling. In: ACM SIGGRAPH Asia (2019)

72. Xiang, Y., Schmidt, T., Narayanan, V., Fox, D.: Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. arXiv preprint arXiv:1711.00199 (2017)

73. Xu, K., Hu, W., Leskovec, J., Jegelka, S.: How powerful are graph neural networks? arXiv preprint arXiv:1810.00826 (2018)

74. Xu, Q., Wang, W., Ceylan, D., Mech, R., Neumann, U.: DISN: Deep implicit surface network for high-quality single-view 3D reconstruction. In: NeurIPS (2019)

75. Yoon, Y., DeSouza, G.N., Kak, A.C.: Real-time tracking and pose estimation for industrial objects using geometric features. In: 2003 IEEE International Conference on Robotics and Automation (Cat. No. 03CH37422). vol. 3, pp. 3473–3478. IEEE (2003)

76. Zeng, A., Song, S., Nießner, M., Fisher, M., Xiao, J., Funkhouser, T.: 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1802–1811 (2017)
77. Zeng, A., Yu, K.T., Song, S., Suo, D., Walker, E., Rodriguez, A., Xiao, J.: Multi-view self-supervised deep learning for 6d pose estimation in the amazon picking challenge. In: 2017 IEEE international conference on robotics and automation (ICRA). pp. 1386–1383. IEEE (2017)