

# PartNet: A Large-scale Benchmark for Fine-grained and Hierarchical Part-level 3D Object Understanding

Kaichun Mo<sup>1</sup> Shilin Zhu<sup>2</sup> Angel X. Chang<sup>3</sup> Li Yi<sup>1</sup> Subarna Tripathi<sup>4</sup> Leonidas J. Guibas<sup>1</sup> Hao Su<sup>2</sup>

<sup>1</sup>Stanford University <sup>2</sup>University of California San Diego <sup>3</sup>Simon Fraser University <sup>4</sup>Intel AI Lab

<https://cs.stanford.edu/~kaichun/partnet/>

## Abstract

We present *PartNet*: a consistent, large-scale dataset of 3D objects annotated with fine-grained, instance-level, and hierarchical 3D part information. Our dataset consists of 573,585 part instances over 26,671 3D models covering 24 object categories. This dataset enables and serves as a catalyst for many tasks such as shape analysis, dynamic 3D scene modeling and simulation, affordance analysis, and others. Using our dataset, we establish three benchmarking tasks for evaluating 3D part recognition: fine-grained semantic segmentation, hierarchical semantic segmentation, and instance segmentation. We benchmark four state-of-the-art 3D deep learning algorithms for fine-grained semantic segmentation and three baseline methods for hierarchical semantic segmentation. We also propose a novel method for part instance segmentation and demonstrate its superior performance over existing methods.

## 1. Introduction

Being able to parse objects into parts is critical for humans to understand and interact with the world. People recognize, categorize, and organize objects based on the knowledge of their parts [8]. Many actions that people take in the real world require detection of parts and reasoning over parts. For instance, we open doors using doorknobs and pull out drawers by grasping their handles. Teaching machines to analyze parts is thus essential for many vision, graphics, and robotics applications, such as predicting object functionality [11, 12], human-object interactions [16], simulation [18], shape editing [27, 14], and shape generation [23, 39].

To enable part-level object understanding by learning approaches, 3D data with part annotations are in high demand. Many cutting-edge learning algorithms, especially for 3D understanding [43, 42, 29], intuitive physics [25], and reinforcement learning [45, 28], require such data to train the

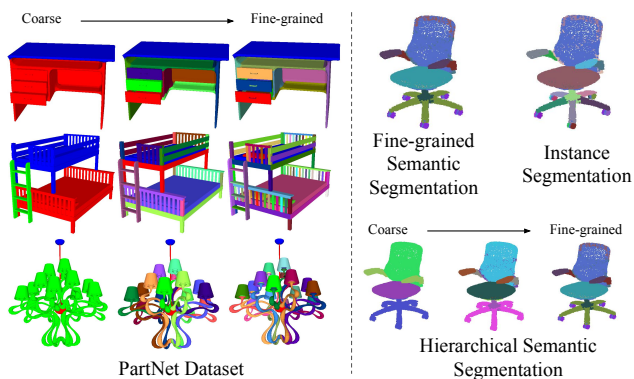


Figure 1. **PartNet dataset and three benchmarking tasks.** Left: we show example annotations at three levels of segmentation in the hierarchy. Right: we propose three fundamental and challenging segmentation tasks and establish benchmarks using PartNet.

networks and benchmark the performances. Researchers are also increasingly interested in synthesizing dynamic data through physical simulation engines [18, 41, 28]. Creation of large-scale animatable scenes will require a large amount of 3D data with affordances and mobility information. Object parts serve as a critical stepping stone to access this information. Thus it is necessary to have a big 3D object dataset with part annotation.

With the availability of the existing 3D shape datasets with part annotations [5, 3, 43], we witness increasing research interests and advances in 3D part-level object understanding. Recently, a variety of learning methods have been proposed to push the state-of-the-art for 3D shape segmentation [29, 30, 44, 17, 33, 22, 7, 37, 38, 40, 31, 6, 24, 21]. However, existing datasets only provide part annotations on relatively small numbers of object instances [5], or on coarse yet non-hierarchical part annotations [43], restricting the applications that involves understanding fine-grained and hierarchical shape segmentation.

In this paper, we introduce PartNet: a *consistent, large-scale* dataset on top of ShapeNet [3] with *fine-grained, hierarchical, instance-level* 3D part information. Collecting

Dataset	#Shape	#Part	#Category	Granularity	Semantics	Hierarchical	Instance-level	Consistent
Chen et al. [5]	380	4,300	19	<b>Fine-grained</b>	No	No	<b>Yes</b>	No
MCL [35]	1,016	7,537	10	<b>Fine-grained</b>	<b>Yes</b>	No	No	<b>Yes</b>
Chang et al. [4]	2,278	27,477	<b>90</b>	<b>Fine-grained</b>	<b>Yes</b>	No	<b>Yes</b>	No
Yi et al. [43]	<b>31,963</b>	80,323	16	Coarse	<b>Yes</b>	No	No	<b>Yes</b>
<b>PartNet (ours)</b>	26,671	<b>573,585</b>	24	<b>Fine-grained</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>

Table 1. Comparison to the other shape part datasets.

such fine-grained and hierarchical segmentation is challenging. The boundary between fine-grained part concepts are more obscure than defining coarse parts. Thus, we define a common set of part concepts by carefully examining the 3D objects to annotate, balancing over several criteria: well-defined, consistent, compact, hierarchical, atomic and complete. Shape segmentation involves multiple levels of granularity. Coarse parts describe more global semantics and fine-grained parts convey richer geometric and semantic details. We organize expert-defined part concepts in hierarchical segmentation templates to guide annotation.

PartNet provides a large-scale benchmark for many part-level object understanding tasks. In this paper, we focus on three fundamental and challenging shape segmentation tasks: fine-grained semantic segmentation, hierarchical semantic segmentation, and instance segmentation. We benchmark four state-of-the-art algorithms on fine-grained semantic segmentation and propose three baseline methods for hierarchical semantic segmentation. We propose the task of part instance segmentation using PartNet. By taking advantages of rich shape structures, we propose a novel method that outperforms the existing baseline algorithm by a clear margin.

PartNet contains highly structured, fine-grained and heterogeneous parts. Our experiments reveals that existing algorithms developed for coarse and homogeneous part understanding cannot work well on PartNet. First, small and fine-grained parts, *e.g.* door handles and keyboard buttons, are abundant and present new challenges for part recognition. Second, many geometrically similar but semantically different parts requires more global shape context to distinguish. Third, understanding the heterogeneous variation of shapes and parts necessitate hierarchical understanding. We expect that PartNet could serve as a better platform for part-level object understanding in the next few years.

In summary, we make the following contributions:

- We introduce PartNet, consisting of 573,585 fine-grained part annotations for 26,671 shapes across 24 object categories. To the best of our knowledge, it is the first *large-scale* dataset with *fine-grained, hierarchical, instance-level* part annotations;
- We propose three part-level object understanding tasks to demonstrate the usefulness of this data: fine-grained semantic segmentation, hierarchical semantic segmentation, and instance segmentation.
- We benchmark four state-of-the-art algorithms for semantic segmentation and three baseline methods for hierarchical segmentation using PartNet;
- We propose the task of part instance segmentation on PartNet and describe a novel method that outperforms the existing baseline method by a large margin.

## 2. Related Work

Understanding shape parts is a long-standing problem in computer vision and graphics. Lacking large-scale annotated datasets, early research efforts evaluated algorithm results qualitatively and conducted quantitative comparison on small sets of 3D models. Attene *et al.* [1] compared 5 mesh segmentation algorithms using 11 3D surface meshes and presented side-by-side qualitative comparison. Chen *et al.* [5] collected 380 surface meshes from 19 object categories with instance-level part decomposition for each shape and proposed quantitative evaluation metrics for shape segmentation. Concurrently, Benhabiles *et al.* [2] proposed similar evaluation criteria and methodology. Kalogerakis *et al.* [15] further assigned semantic labels to the segmented components. Shape co-segmentation benchmarks [36, 9] were proposed to study co-segmentation among similar shapes.

Recent advances in deep learning have demonstrated the power and efficiency of data-driven methods on 3D shape understanding tasks such as classification, segmentation and generation. ShapeNet [3] collected a large-scale synthetic 3D CAD models from online open-sourced 3D repositories, including more than 3,000,000 models and 3,135 object categories. Yi *et al.* [43] took an active learning approach to annotate the ShapeNet models with semantic segmentation for 31,963 shapes covering 16 object categories. In their dataset, each object is usually decomposed into 2~5 coarse semantic parts. PartNet provides more fine-grained part annotations that contains 18 parts per shape on average.

Many recent works studied fine-grained and hierarchical shape segmentation. Yi *et al.* [42] leveraged the noisy part decomposition inputs in the CAD model designs and trained per-category models to learn consistent shape hierarchy. Chang *et al.* [4] collected 27,477 part instances from 2,278 models covering 90 object categories and studied the part properties related to language. Wang *et al.* [35] proposed multi-component labeling benchmark containing

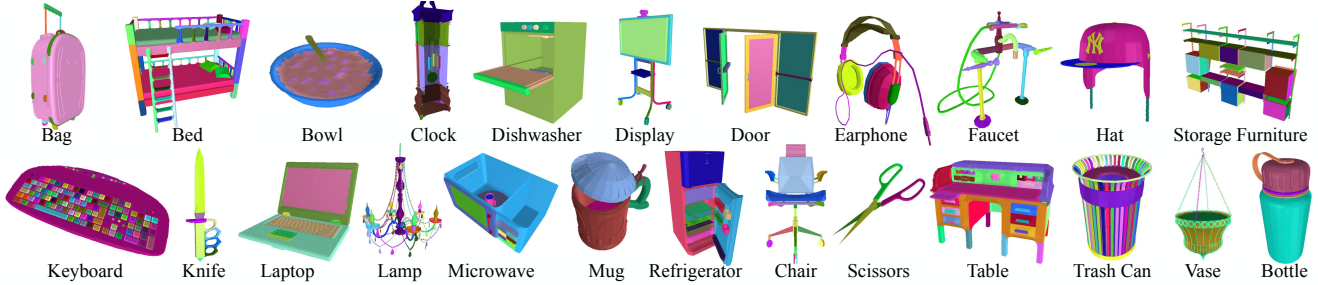


Figure 2. **PartNet dataset.** We visualize example shapes with fine-grained part annotations for the 24 object categories in PartNet.

	All	Bag	Bed	Bott	Bowl	Chair	Clock	Dish	Disp	Door	Ear	Fauc	Hat	Key	Knife	Lamp	Lap	Micro	Mug	Frid	Scis	Stora	Table	Trash	Vase
#A	32537	186	248	519	247	8176	624	241	1005	285	285	840	287	210	514	3408	485	268	252	247	127	2639	9906	378	1160
#S	26671	146	212	464	208	6400	579	201	954	245	247	708	250	174	384	2271	453	212	212	207	88	2303	8309	340	1104
#M	771	20	18	28	20	77	25	20	26	20	19	60	19	18	57	64	20	28	20	20	34	91	19	28	
#PS	480	4	24	12	4	57	23	12	8	8	15	18	8	3	16	83	8	12	4	13	5	36	82	15	10
#PI	573K	664	9K	2K	615	176K	4K	2K	7K	2K	3K	8K	1K	20K	3K	50K	3K	2K	839	2K	981	77K	177K	8K	5K
$P_{med}$	14	4	33	5	2	19	5	9	8	7	12	9	4	106	7	12	8	7	3	9	8	24	15	9	4
$P_{max}$	230	7	169	7	4	153	32	16	12	20	14	34	5	127	10	230	8	17	6	33	9	220	214	143	200
$D_{med}$	3	1	5	2	1	3	3	3	3	3	3	3	2	1	3	5	2	3	1	3	2	4	4	2	2
$D_{max}$	7	1	5	2	1	5	4	3	3	3	3	3	2	1	3	7	2	3	1	3	2	5	6	2	3

Table 2. **PartNet statistics.** Row #A, #S, #M respectively show the number of shape annotations, the number of distinct shape instances and the number of shapes that we collect multiple annotations. Row #PS and #PI show the number of different part semantics and part instances that we finally collect. Row  $P_{med}$  and  $P_{max}$  respectively indicate the median and maximum number of part instances per shape. Row  $D_{med}$  and  $D_{max}$  respectively indicate the median and maximum hierarchy depth per shape, with root node as depth 0.

1,016 3D models from ShapeNet [3] from 10 object categories with manually annotated fine-grained level part semantics and studied to learn neural networks for grouping and labeling fine-grained part components. PartNet proposes a large-scale dataset with 573,585 fine-grained and hierarchical shape part annotations covering 26,671 models from 24 object categories.

There are also many previous works that attempted to understand parts by their functionality and articulation. Hu *et al.* [11] constructed a dataset of 608 objects from 15 object categories annotated with the object functionality and introduced a co-analysis method to learn category-wise object functionality. Hu *et al.* [10] proposed a dataset of 368 mobility units with diverse types of articulation and learned to predict part mobility information from a single static segmented 3D mesh. In PartNet, we assign consistent semantic labels that entail such functionality and articulation information for part components within each object category, which potentially makes PartNet support such research.

### 3. Data Annotation

The data annotation is performed in a hierarchical manner. Expert-defined hierarchical part templates are provided to guarantee labeling consistency among multiple annotators. We design a single-thread question-answering 3D GUI to guide the annotation. We hire 66 professional annotators and train them for the annotation. The average annotation time per shape is 8 minutes, and at least one pass of verification is performed for each annotation to ensure accuracy.

### 3.1. Expert-Defined Part Hierarchy

Shape segmentation naturally involves hierarchical understanding. People understand shapes at different segmentation granularity. Coarse parts convey global semantics while fine-grained parts provide more detailed understanding. Moreover, fine-grained part concepts are more obscure to define than coarse parts. Different annotators have different knowledge and background so that they may name parts differently when using free-form annotation [4]. To address the issues, we introduce And-Or-Graph-style hierarchical templates and collect part annotations according to the pre-defined templates.

Due to the lack of well-acknowledged rules of thumb to define good templates, the task of designing hierarchical part templates for a category becomes a non-trivial task. Furthermore, the requirement for the designed template to cover all variations of shapes and parts, makes the problem more challenging. Below we summarize the criteria that we use to guide our template design:

- **Well-defined:** Part concepts are well-delineated such that parts are identifiable by multiple annotators;
- **Consistent:** Part concepts are shared and reused across different parts, shapes and object categories;
- **Compact:** There is no unnecessary part concept and part concepts are reused when it is possible;
- **Hierarchical:** Part concepts are organized in a taxonomy to cover both coarse and fine-grained parts;
- **Atomic:** Leaf nodes in the part taxonomy consist of primitive, non-decomposable shapes;

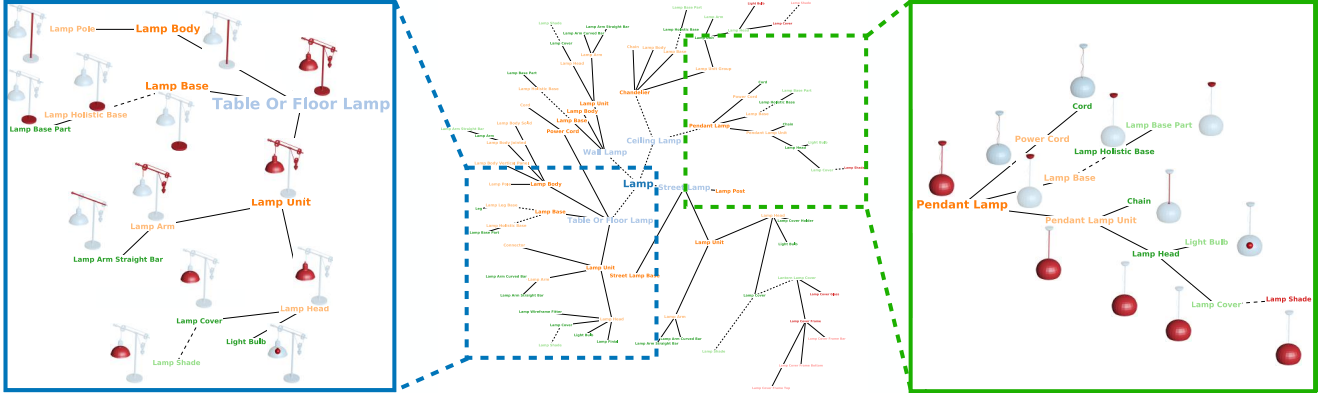


Figure 3. We show the expert-defined hierarchical template for lamp (middle) and the instantiations for a table lamp (left) and a ceiling lamp (right). The And-nodes are drawn in solid lines and Or-nodes in dash lines. The template is deep and comprehensive to cover structurally different types of lamps. In the meantime, the same part concepts, such as light bulb and lamp shade, are shared across the different types.

- **Complete:** The part taxonomy covers a heterogeneous variety of shapes as completely as possible.

Guided by these general principles, we build an And-Or-Graph-style part template for each object category. The templates are defined by experts after examining a broad variety of objects in the category. Each template is designed in a hierarchical manner from the coarse semantic parts to the fine-grained primitive-level components. Figure 3 (middle) shows the lamp template. And-nodes segment a part into small subcomponents. Or-nodes indicate subcategorization for the current part. The combination of And-nodes and Or-nodes allows us to cover structurally different shapes using the same template while sharing as much common part labels as possible. As in Figure 3 (left) and (right), both table lamps and ceiling lamps are explained by the same template through the first-level Or-node for lamp types.

Despite the depth and comprehensiveness of these templates, it is still impossible to cover all cases. Thus, we allow our annotators to improve upon the structure of the template and to annotate parts that are out of the scope of our definition. We also conduct template refinements to resolve part ambiguity after we obtain the data annotation according to the original templates. To systematically identify ambiguities, we reserve a subset of shapes from each class and collect multiple human annotations for each shape. We compute the confusion matrix among different annotators and address data inconsistencies. For example, we merge two concepts with high confusion scores or remove a part if it is frequently segmented in the wrong way. We provide more details about this in the supplementary material.

### 3.2. Annotation Interface

Figure 4 (a) shows our web-based annotation interface. Based on the template hierarchy, the annotation process is designed to be a single-thread question-answering workflow, traversing the template graph in a depth-first manner, as shown in Figure 4 (b). Starting from the root node, the

annotator is asked a sequence of questions. The answers automatically construct the final hierarchical segmentation for the current shape instance. For each question, the annotator is asked to mark the number of subparts (And-node) or pick one among all subtypes (Or-node) for a given part. For each leaf node part, the annotator annotates the part geometry in the 3D interface. To help them understand the part definition and specification, we provide rich textual definitions and visual examples for each part. In addition, our interface supports cross-section and visibility control to annotate the interior structure of a 3D model.

The collected 3D CAD models often include original mesh subgroups and part information. Some of the grouping information is detailed enough to determine the final segmentation we need. Considering this, we provide the annotators with the original groupings at the beginning of the annotation, to speed up annotation. The annotators can simply select multiple predefined pieces to form a part of the final segmentation. We also provide mesh cutting tools to split large pieces into smaller ones following [5], when the original groupings are coarser than the desired segmentation, as shown in Figure 4 (c). The annotators draw boundary lines on the remeshed watertight surface [13] and the mesh cutting algorithm automatically splits the mesh into multiple smaller subcomponents.

In contrast to prior work, our UI is designed for operating directly on 3D models and collecting fine-grained and hierarchical part instances. Compared to Yi *et al.* [43] where the annotation is performed in 2D, our approach allows the annotators to directly annotate on the 3D shapes and thus be able to pick up more subtle part details that are hidden from 2D renderings. Chang *et al.* [4] proposes a 3D UI that paints regions on mesh surfaces for part labeling. However, their interface is limited to existing over-segmentations on part components and does not support hierarchical annotations.

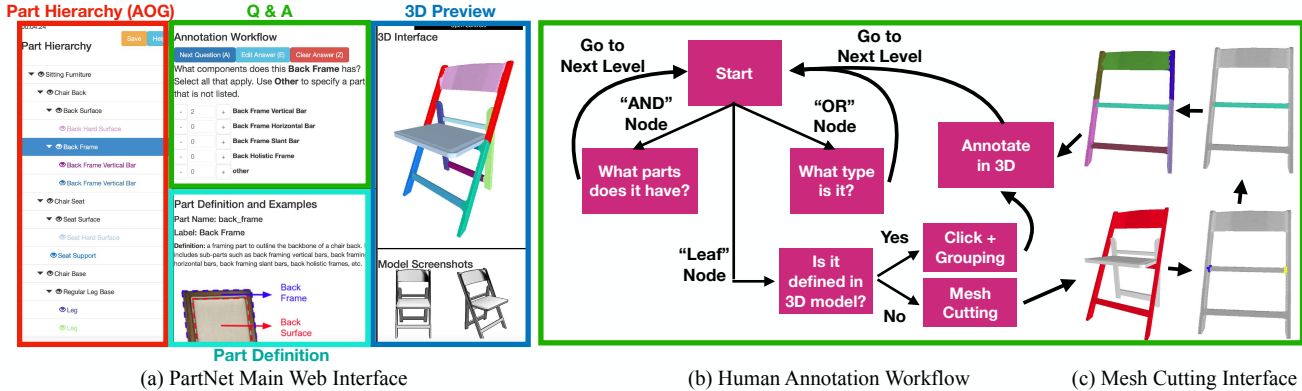


Figure 4. We show our annotation interface with its components, the proposed question-answering workflow and the mesh cutting interface.

## 4. PartNet Dataset

The final PartNet dataset provides fine-grained and hierarchical instance-level part segmentation annotation for 26,671 shapes with 573,585 part instances from 24 object categories. Most of the shapes and object categories are from ShapeNetCore [3]. We supplement 3 object categories that are commonly present in indoor scenes (*i.e.* scissors, refrigerators, and doors) and augment 7 of the existing categories with more 3D models from 3D Warehouse<sup>1</sup>.

Figure 2 and Table 2 show the PartNet data and statistics. More visualization and statistics are included in supplemental material. Our templates define hierarchical segmentation with 3 depth in median and 7 maximum. In total, we annotate 573,585 fine-grained part instances, with a median of 14 parts per shape and a maximum of 230. To study annotation consistency, we also collect a subset of 771 shapes and ask for multiple annotations per shape.

## 5. Tasks and Benchmarks

We benchmark three part-level object understanding tasks using PartNet: fine-grained semantic segmentation, hierarchical semantic segmentation and instance segmentation. Four state-of-the-art algorithms for semantic segmentation are evaluated and three baseline methods are proposed for hierarchical segmentation. Moreover, we propose a novel method for instance segmentation that outperforms the existing baseline method.

**Data Preparation.** In this section, we only consider parts that can be fully determined by their shape geometry<sup>2</sup>. We ignore the parts in evaluation that require additional information to identify, such as the glass parts on the cabinet doors which opacity is needed to identify, and the buttons on microwaves that texture information is desired to distinguish it from the main frame. We also remove rarely ap-

<sup>1</sup><https://3dwarehouse.sketchup.com>

<sup>2</sup>Although 3D models in ShapeNet [3] come with face normal, textures, material and other information, there is no guarantee for the quality of such information. Thus, we leave this as a future work.

peared parts from the evaluation, as the lacking of samples is insufficient for training and evaluating networks.

We sample 10,000 points from each CAD model with furthest point sampling and use the 3D coordinates as the neural network inputs for all the experiments in the paper. The proposed dataset is split into train, validation and test sets with the ratio 70%: 10%: 20%. The shapes with multiple human annotations are not used in the experiments.

### 5.1. Fine-grained Semantic Segmentation

Recent advances of 3D semantic segmentation [29, 30, 44, 17, 33, 22, 7, 37, 38, 40, 31, 6, 24, 21] have accomplished promising achievement in coarse-level segmentation on the ShapeNet Part dataset [3, 43]. However, few work focus on the fine-grained 3D semantic segmentation, due to the lack of large-scale fine-grained dataset. With the help of the proposed PartNet dataset, researchers can now work on this more challenging task with little overhead.

Fine-grained 3D semantic segmentation requires recognizing and distinguishing small and similar semantic parts. For example, door handles are usually small, 77 out of 10,000 points on average in PartNet, but semantically important on doors. Beds have several geometrically similar parts such as side vertical bars, post bars and base legs. To recognize the subtle part details, segmentation systems need to understand them locally, through discriminative features, and globally, in the context of the whole shape.

**Benchmark Algorithms.** We benchmark four state-of-the-art semantic segmentation algorithms on the fine-grained PartNet segmentation: PointNet [29], PointNet++ [30], SpiderCNN [40] and PointCNN [24]<sup>3</sup>. PointNet [29] takes unordered point sets as inputs and extracts features for shape classification and segmentation. To better learn local geometric features, the follow-up work PointNet++ [30] proposes a hierarchical feature extraction scheme. SpiderCNN [40] extends traditional convolution operations on 2D

<sup>3</sup>There are many other algorithm candidates: [44, 17, 33, 22, 7, 37, 38, 31, 6, 21]. We will host an online leaderboard to report the performances.

	Avg	Bag	Bed	Bott	Bowl	Chair	Clock	Dish	Disp	Door	Ear	Fauc	Hat	Key	Knife	Lamp	Lap	Micro	Mug	Frid	Scis	Stora	Table	Trash	Vase
<b>P1</b>	57.9	42.5	32.0	33.8	58.0	64.6	33.2	76.0	86.8	<b>64.4</b>	53.2	58.6	55.9	65.6	62.2	29.7	96.5	49.4	80.0	49.6	86.4	51.9	50.5	55.2	54.7
<b>P2</b>	37.3	—	20.1	—	—	38.2	—	55.6	—	38.3	—	—	—	—	—	27.0	—	41.7	—	35.5	—	44.6	34.3	—	—
<b>P3</b>	35.6	—	13.4	29.5	—	27.8	28.4	48.9	76.5	30.4	33.4	47.6	—	—	32.9	18.9	—	37.2	—	33.5	—	38.0	29.0	34.8	44.4
<b>Avg</b>	51.2	42.5	21.8	31.7	58.0	43.5	30.8	60.2	81.7	44.4	43.3	53.1	55.9	65.6	47.6	25.2	96.5	42.8	80.0	39.5	86.4	44.8	37.9	45.0	49.6
<b>P<sup>+</sup>1</b>	<b>65.5</b>	59.7	51.8	53.2	67.3	68.0	<b>48.0</b>	80.6	89.7	59.3	<b>68.5</b>	64.7	62.4	62.2	<b>64.9</b>	<b>39.0</b>	96.6	55.7	83.9	51.8	87.4	58.0	<b>69.5</b>	64.3	<b>64.4</b>
<b>P<sup>+</sup>2</b>	44.5	—	38.8	—	—	43.6	—	55.3	—	49.3	—	—	—	—	—	<b>32.6</b>	—	48.2	—	41.9	—	49.6	<b>41.1</b>	—	—
<b>P<sup>+</sup>3</b>	42.5	—	30.3	41.4	—	39.2	<b>41.6</b>	50.1	80.7	32.6	38.4	52.4	—	—	<b>34.1</b>	<b>25.3</b>	—	48.5	—	36.4	—	40.5	<b>33.9</b>	46.7	49.8
<b>Avg</b>	58.1	59.7	40.3	<b>47.3</b>	67.3	50.3	<b>44.8</b>	62.0	85.2	47.1	53.5	58.6	62.4	62.2	<b>49.5</b>	<b>32.3</b>	96.6	50.8	83.9	43.4	87.4	49.4	<b>48.2</b>	55.5	57.1
<b>S1</b>	60.4	57.2	55.5	<b>54.5</b>	<b>70.6</b>	67.4	33.3	70.4	90.6	52.6	46.2	59.8	63.9	64.9	37.6	30.2	<b>97.0</b>	49.2	83.6	50.4	75.6	<b>61.9</b>	50.0	62.9	63.8
<b>S2</b>	41.7	—	40.8	—	—	39.6	—	59.0	—	48.1	—	—	—	—	—	24.9	—	47.6	—	34.8	—	46.0	34.5	—	—
<b>S3</b>	37.0	—	36.2	32.2	—	30.0	24.8	50.0	80.1	30.5	37.2	44.1	—	—	22.2	19.6	—	43.9	—	39.1	—	44.6	20.1	42.4	32.4
<b>Avg</b>	53.6	57.2	44.2	43.4	<b>70.6</b>	45.7	29.1	59.8	85.4	43.7	41.7	52.0	63.9	64.9	29.9	24.9	<b>97.0</b>	46.9	83.6	41.4	75.6	50.8	34.9	52.7	48.1
<b>C1</b>	64.3	<b>66.5</b>	<b>55.8</b>	49.7	61.7	<b>69.6</b>	42.7	<b>82.4</b>	<b>92.2</b>	63.3	64.1	<b>68.7</b>	<b>72.3</b>	<b>70.6</b>	62.6	21.3	<b>97.0</b>	<b>58.7</b>	<b>86.5</b>	<b>55.2</b>	<b>92.4</b>	61.4	17.3	<b>66.8</b>	63.4
<b>C2</b>	<b>46.5</b>	—	<b>42.6</b>	—	—	<b>47.4</b>	—	<b>65.1</b>	—	<b>49.4</b>	—	—	—	—	—	22.9	—	<b>62.2</b>	—	<b>42.6</b>	—	<b>57.2</b>	29.1	—	—
<b>C3</b>	<b>46.4</b>	—	<b>41.9</b>	<b>41.8</b>	—	<b>43.9</b>	36.3	<b>58.7</b>	<b>82.5</b>	<b>37.8</b>	<b>48.9</b>	<b>60.5</b>	—	—	<b>34.1</b>	20.1	—	<b>58.2</b>	—	<b>42.9</b>	—	<b>49.4</b>	21.3	<b>53.1</b>	<b>58.9</b>
<b>Avg</b>	<b>59.8</b>	<b>66.5</b>	<b>46.8</b>	45.8	61.7	<b>53.6</b>	39.5	<b>68.7</b>	<b>87.4</b>	<b>50.2</b>	<b>56.5</b>	<b>64.6</b>	<b>72.3</b>	<b>70.6</b>	48.4	21.4	<b>97.0</b>	<b>59.7</b>	<b>86.5</b>	<b>46.9</b>	<b>92.4</b>	<b>56.0</b>	22.6	<b>60.0</b>	<b>61.2</b>

Table 3. **Fine-grained semantic segmentation results (part-category mIoU %)**. Algorithm **P**, **P<sup>+</sup>**, **S** and **C** refer to PointNet [29], PointNet++ [30], SpiderCNN [40] and PointCNN [24], respectively. The number **1**, **2** and **3** refer to the three levels of segmentation: coarse-, middle- and fine-grained. We put short lines for the levels that are not defined.

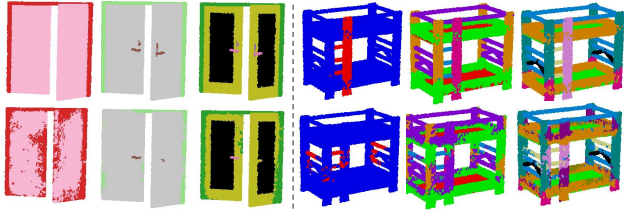


Figure 5. **Qualitative results for semantic segmentation.** The top row shows the ground-truth and the bottom row shows the PointCNN prediction. The black points indicate unlabeled points.

images to 3D point clouds by parameterizing a family of convolutional filters. To organize the unordered points into latent canonical order, PointCNN [24] proposes to learn  $\mathcal{X}$ -transformation, and applies  $\mathcal{X}$ -convolution operations on the canonical points.

We train the four methods on the dataset, using the default network architectures and hyperparameters described in their papers. Instead of training a single network for all object categories as done in most of these papers, we train a network for each category at each segmentation level. We input only the 3D coordinates for fair comparison<sup>4</sup> and train the networks until convergence. More training details are described in the supplementary material.

**Evaluation and Results.** We evaluate the algorithms at three segmentation levels for each object category: coarse-, middle- and fine-grained. The coarse level approximately corresponds to the granularity in Yi *et al.* [43]. The fine-grained level refers to the segmentation down to leaf levels in the segmentation hierarchies. For structurally deep hierarchies, we define the middle level in between. Among 24 object categories, all of them have the coarse levels, while 9

<sup>4</sup>PointNet++ [30] and SpiderCNN [40] use point normals as additional inputs. For fair comparison, we only input the 3D coordinates.

have the middle levels and 17 have the fine levels. Overall, we define 50 segmentation levels for 24 object categories.

In Table 3, we report the semantic segmentation performances at multiple levels of granularity on PartNet. We use the mean Intersection-over-Union (mIoU) scores as the evaluation metric. After removing unlabeled ground-truth points, for each object category, we first calculate the IoU between the predicted point set and the ground-truth point set for each semantic part category across all test shapes. Then, we average the per-part-category IoUs to compute the mIoU for the object category. We further calculate the average mIoU across different levels for each object category and finally report the average cross all object categories.

Unsurprisingly, performance for all four algorithms drop by a large margin from the coarse level to the fine-grained level. Figure 5 shows qualitative results from PointCNN. The method does not perform well on small parts, such as the door handle on the door example, and visually similar parts, such as stair steps and the horizontal bars on the bed frame. How to learn discriminative features that better capture both local geometry and global context for these issues would be an interest topic for future works.

## 5.2. Hierarchical Semantic Segmentation

Shape segmentation is hierarchical by its nature. From coarse semantics to fine-grained details, hierarchical understanding on 3D objects develops a holistic and comprehensive reasoning on the shape components. For this purpose, we study hierarchical semantic segmentation problem that predicts semantic part labels in the entire shape hierarchies that cover both coarse- and fine-grained part concepts.

A key problem towards hierarchical segmentation is how to leverage the rich part relationships on the given shape templates in the learning procedure. Recognizing a chair base as a finer-level swivel base significantly reduces the

	Avg	Bed	Bott	Chair	Clock	Dish	Disp	Door	Ear	Fauc	Knife	Lamp	Micro	Frid	Stora	Table	Trash	Vase
<b>Bottom-Up</b>	51.2	40.8	56.1	47.2	38.3	<b>61.5</b>	84.1	<b>52.6</b>	54.3	<b>63.4</b>	52.3	36.8	48.2	41.0	46.8	38.3	<b>53.6</b>	54.4
<b>Top-Down</b>	50.8	41.1	<b>56.2</b>	46.5	34.3	54.5	<b>84.7</b>	50.6	<b>59.5</b>	61.4	<b>55.6</b>	37.1	<b>48.8</b>	<b>41.6</b>	45.2	37.0	53.5	<b>55.6</b>
<b>Ensemble</b>	<b>51.7</b>	<b>42.0</b>	54.7	<b>48.1</b>	<b>44.5</b>	58.8	<b>84.7</b>	51.4	57.2	61.9	51.9	<b>37.6</b>	47.5	41.4	<b>47.3</b>	<b>44.0</b>	52.8	53.1

Table 4. **Hierarchical segmentation results (part-category mIoU %)**. We present the hierarchical segmentation performances for three baseline methods: bottom-up, top-down and ensemble. We conduct experiments on 17 out of 24 categories with tree depth bigger than 1.

solution space for detecting more fine-grained parts such as central supporting bars, star-base legs and wheels. On the other hand, the lack of a chair back increases the possibility that the object is a stool. Different from Sec. 5.1 where we consider the problem at each segmentation level separately, hierarchical segmentation requires a holistic understanding on the entire part hierarchy.

**Benchmark Algorithms.** We propose three baseline methods to tackle hierarchical segmentation: bottom-up, top-down and ensemble. The bottom-up method considers only the leaf-node parts during training and groups the prediction of the children nodes to parent nodes as defined in the hierarchies in the bottom-up inference. The top-down method learns a multi-labeling task over all part semantic labels on the tree and conducts a top-down inference by classifying coarser-level nodes first and then finer-level ones. For the ensemble method, we train flat segmentation at multiple levels as defined in Sec. 5.1 and conduct joint inference by calculating the average log-likelihood scores over all the root-to-leaf paths on the tree. We use PointNet++ [30] as the backbone network in this work, and other methods listed in Sec. 5.1 can also be used. More architecture and training details are described in the supplementary material.

**Evaluation and Results.** Table 8 demonstrates the performances of the three baseline methods. We calculate mIoU for each part category and compute the average over all the tree nodes as the evaluation metric. The experimental results show that the three methods perform similarly with small performance gaps. The ensemble method performs slightly better over the other two, especially for the categories with rich structural and sub-categorization variation, such as chair, table and clock.

The bottom-up method only considers leaf-node parts in the training. Although the template structure is not directly used, the parent-node semantics of leaf nodes are implicitly encoded in the leaf-node part definitions. For example, the vertical bars for chair backs and chair arms are two different leaf nodes. The top-down method explicitly leverages the tree structures in both the training and the testing phases. However, prediction errors are accumulated through top-down inference. The ensemble method decouples the hierarchical segmentation task into individual tasks at multiple levels and performs joint inference, taking the predictions at all levels into consideration. Though demonstrating better performances, it has more hyper-parameters and requires longer training time for the multiple networks.

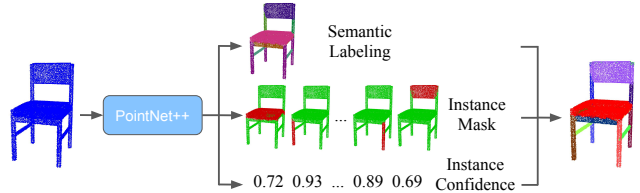


Figure 6. **The proposed detection-by-segmentation method for instance segmentation.** The network learns to predict three components: the semantic label for each point, a set of disjoint instance masks and their confidence scores for part instances.

### 5.3. Instance Segmentation

The goal of instance segmentation is to detect every individual part instance and segment it out from the context of the shape. Many applications in computer graphics, vision and robotics, including manufacturing, assembly, interaction and manipulation, require the instance-level part recognition. Compared to detecting objects from scenes, parts on objects usually have stronger and more intertwined structural relationships. The existence of many visually-similar but semantically-different parts makes the part detection problem challenging. To the best of our knowledge, this work is the first to provide a large-scale shape part instance-level segmentation benchmark.

Given a shape point cloud as input, the task of part instance segmentation is to provide several disjoint masks over the entire point cloud, each of which corresponds to an individual part instance on the object. We adopt the part semantics from the defined segmentation levels in Sec. 5.1. The detected masks should have no overlaps, but they together do not necessarily cover the entire point cloud, as some points may not belong to any part of interests.

**Benchmark Algorithms.** We propose a novel detection-by-segmentation network to address instance segmentation. We illustrate our network architecture in Figure 6. We use PointNet++ [30] as the backbone network for extracting features and predicting both semantic segmentation for each point and  $K$  instance segmentation masks  $\{\hat{y}_i \in [0, 1]^N | i = 1, 2, \dots, K\}$  over the input point cloud of size  $N$ . Moreover, we train a separate mask  $\hat{y}_{other}$  for the points without semantic labels in the ground-truth. A softmax activation layer is applied to encourage the mutual exclusiveness among different masks so that  $\hat{y}_1 + \hat{y}_2 + \dots + \hat{y}_K + \hat{y}_{other} = \mathbf{1}$ . To train the network, we apply Hungarian algorithm [20] to find a bipartite matching  $\mathcal{M} : \{i \rightarrow \mathcal{M}(i) | i = 1, 2, \dots, T\}$  between the prediction

	Avg	Bag	Bed	Bott	Bowl	Chair	Clock	Dish	Disp	Door	Ear	Fauc	Hat	Key	Knife	Lamp	Lap	Micro	Mug	Frid	Scis	Stora	Table	Trash	Vase
S1	55.7	38.8	29.8	61.9	56.9	72.4	20.3	72.2	89.3	49.0	57.8	63.2	68.7	20.0	63.2	32.7	100	50.6	82.2	50.6	71.7	32.9	49.2	56.8	46.6
S2	29.7	—	15.4	—	—	25.4	—	58.1	—	25.4	—	—	—	—	21.7	—	—	49.4	—	22.1	—	30.5	18.9	—	—
S3	29.5	—	11.8	45.1	—	19.4	18.2	38.3	78.8	15.4	35.9	37.8	—	—	38.3	14.4	—	32.7	—	18.2	—	21.5	14.6	24.9	36.5
Avg	46.8	38.8	19.0	53.5	56.9	39.1	19.3	56.2	84.0	29.9	46.9	50.5	68.7	20.0	50.7	22.9	100	44.2	82.2	30.3	71.7	28.3	27.5	40.9	41.6
O1	62.6	64.7	48.4	63.6	59.7	74.4	42.8	76.3	93.3	52.9	57.7	69.6	70.9	43.9	58.4	37.2	100	50.0	86.0	50.0	80.9	45.2	54.2	71.7	49.8
O2	37.4	—	23.0	—	—	35.5	—	62.8	—	39.7	—	—	—	—	26.9	—	—	47.8	—	35.2	—	35.0	31.0	—	—
O3	36.6	—	15.0	48.6	—	29.0	32.3	53.3	80.1	17.2	39.4	44.7	—	—	45.8	18.7	—	34.8	—	26.5	—	27.5	23.9	33.7	52.0
Avg	54.4	64.7	28.8	56.1	59.7	46.3	37.5	64.1	86.7	36.6	48.5	57.1	70.9	43.9	52.1	27.6	100	44.2	86.0	37.2	80.9	35.9	36.4	52.7	50.9

Table 5. **Instance segmentation results (part-category mAP %, IoU threshold 0.5)**. Algorithm S and O refer to SGPN [34] and our proposed method respectively. The number 1, 2 and 3 refer to the three levels of segmentation: coarse-, middle- and fine-grained. We put short lines for the levels that are not defined.

masks  $\{\hat{y}_i | i = 1, 2, \dots, K\}$  and the ground-true masks  $\{y_i | i = 1, 2, \dots, T\}$ , and regress each prediction  $\hat{y}_{M(t)}$  to the matched ground-truth mask  $y_t$ . We employ a relaxed version of IoU [19] defined as  $\text{IoU}(p, q) = \langle p, q \rangle / (\|p\|_1 + \|q\|_1 - \langle p, q \rangle)$ , as the metric for Hungarian algorithm. In the meanwhile, a separate branch is trained to predict confidence scores for the predicted masks  $\{C_i | i = 1, 2, \dots, K\}$ .

The loss function is defined as  $L = L_{sem} + \lambda_{ins} L_{ins} + \lambda_{other} L_{other} + \lambda_{conf} L_{conf} + \lambda_{l21} L_{l21}$ , combining five terms: a cross-entropy semantic segmentation loss  $L_{sem}$ , an IoU loss for mask regression  $L_{ins} = \sum_{i=1}^T \text{IoU}(\hat{y}_{M(i)}, y_i)$ , an IoU loss for the unlabeled points  $L_{other} = \text{IoU}(\hat{y}_{other}, y_{other})$ , a prediction-confidence loss  $L_{conf} = \sum_{i=1}^T (C_{M(i)} - \text{IoU}(\hat{y}_{M(i)}, y_i))^2$  and a  $l_{2,1}$ -norm regularization term  $L_{l21} = \sum_{i=1}^K \|\hat{y}_i\|_2 + \|\hat{y}_{other}\|_2$  to encourage unused prediction masks to vanish [32]. We use  $N = 10,000$ ,  $K = 200$ ,  $\lambda_{ins} = 1.0$ ,  $\lambda_{other} = 1.0$ ,  $\lambda_{conf} = 1.0$  and  $\lambda_{l21} = 0.1$  for all the experiments.

We compare the proposed method with SGPN [34], which learns similarity scores among all pairs of points and detect part instances by grouping points that share similar features. We follow most of the default settings and hyper-parameters described in their paper. We first pre-train PointNet++ semantic segmentation branch and then fine-tune it for improving the per-point feature similarity matrix and confidence maps. We use margin values of 1 and 2 for the double-hinge loss as suggested by the authors of [34], instead of 10 and 80 in the original paper. We feed 10,000 points to the network at a time, and use a batch-size of 32 in the pre-training and 1 in the fine-tuning.

**Evaluation and Results.** Table 9 reports the per-category mean Average Precision (mAP) scores for SPGN and our proposed method. For each object category, the mAP score calculates the AP for each semantic part category across all test shapes and averages the AP across all part categories. Finally, we take the average of the mAP scores across different levels of segmentation within each object category and then report the average over all object categories. We compute the IoU between each prediction mask and the closest ground-truth mask and regard a prediction mask as true positive when IoU is larger than 0.5.

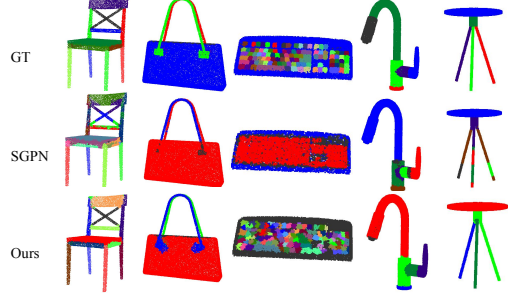


Figure 7. **Qualitative results for instance segmentation.** Our method produces more robust and cleaner results than SGPN.



Figure 8. **Learned instance correspondences.** The corresponding parts are marked with the same color.

Figure 7 shows qualitative comparisons for our proposed method and SGPN. Our proposed method produces more robust and cleaner instance predictions. After learning for point features, SGPN has a post-processing stage that merges points with similar features as one component. This process involves many thresholding hyper-parameters. Even though most parameters are automatically inferred from the validation data, SGPN still suffers from predicting partial or noisy instances in the case of bad thresholding. Our proposed method learns structural priors within each object category that is more instance-aware and robust in predicting complete instances. We observe that training for a set of disjoint masks across multiple shapes gives us consistent part instances. We show the learned part correspondence in Figure 8.

## 6. Conclusion

We introduce PartNet: a *large-scale* benchmark for *fine-grained, hierarchical, and instance-level* 3D shape segmentation. It contains 573,585 part annotations for 26,671



ShapeNet [3] models from 24 object categories. Based on the dataset, we propose three shape segmentation benchmarks: fine-grained semantic segmentation, hierarchical semantic segmentation and instance segmentation. We benchmark four state-of-the-art algorithms for semantic segmentation and propose a novel method for instance segmentation that outperforms the existing baseline method. Our dataset enables future research directions such as collecting more geometric and semantic annotation on parts, investigating shape grammars for synthesis and animating object articulation in virtual environments for robotic learning.

## Acknowledgements

This research was supported by NSF grants CRI-1729205 and IIS-1763268, a Vannevar Bush Faculty Fellowship, a Google fellowship, and gifts from Autodesk, Google and Intel AI Lab. We especially thank Zhe Hu from Hikvision for the help on data annotation and Linfeng Zhao for the help on preparing hierarchical templates. We appreciate the 66 annotators from Hikvision, Ytuu and Data++ on data annotation.

## References

- [1] M. Attene, S. Katz, M. Mortara, G. Patané, M. Spagnuolo, and A. Tal. Mesh segmentation—a comparative study. In *Shape Modeling and Applications, 2006. SMI 2006. IEEE International Conference on*, pages 7–7. IEEE, 2006. 2
- [2] H. Benhabiles, J.-P. Vandeborre, G. Lavoué, and M. Daoudi. A framework for the objective evaluation of segmentation algorithms using a ground-truth of human segmented 3D-models. In *IEEE International Conference on Shape Modeling and Applications (SMI)*, pages Session–5, 2009. 2
- [3] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, et al. ShapeNet: An information-rich 3D model repository. *arXiv preprint arXiv:1512.03012*, 2015. 1, 2, 5, 9, 12, 13
- [4] A. X. Chang, R. Mago, P. Krishna, M. Savva, and C. Fellbaum. Linking WordNet to 3D shapes. In *Global WordNet Conference*, 2018. 2, 3, 4, 11
- [5] X. Chen, A. Golovinskiy, and T. Funkhouser. A benchmark for 3D mesh segmentation. *ACM Transactions on Graphics (Proc. SIGGRAPH)*, 2009. 1, 2, 4, 12
- [6] B. Graham, M. Engelcke, and L. van der Maaten. 3D semantic segmentation with submanifold sparse convolutional networks. *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, pages 18–22, 2018. 1, 5
- [7] P. Hermosilla, T. Ritschel, P.-P. Vázquez, À. Vinacua, and T. Ropinski. Monte carlo convolution for learning on non-uniformly sampled point clouds. *arXiv preprint arXiv:1806.01759*, 2018. 1, 5
- [8] D. D. Hoffman and W. A. Richards. Parts of recognition. *Cognition*, 18(1-3):65–96, 1984. 1
- [9] R. Hu, L. Fan, and L. Liu. Co-segmentation of 3D shapes via subspace clustering. In *Computer graphics forum*, volume 31, pages 1703–1713. Wiley Online Library, 2012. 2
- [10] R. Hu, W. Li, O. Van Kaick, A. Shamir, H. Zhang, and H. Huang. Learning to predict part mobility from a single static snapshot. *ACM Transactions on Graphics (TOG)*, 36(6):227, 2017. 3
- [11] R. Hu, O. van Kaick, B. Wu, H. Huang, A. Shamir, and H. Zhang. Learning how objects function via co-analysis of interactions. *ACM Transactions on Graphics (TOG)*, 35(4):47, 2016. 1, 3
- [12] R. Hu, Z. Yan, J. Zhang, O. van Kaick, A. Shamir, H. Zhang, and H. Huang. Predictive and generative neural networks for object functionality. In *Computer Graphics Forum (Eurographics State-of-the-art report)*, volume 37, pages 603–624, 2018. 1
- [13] J. Huang, H. Su, and L. Guibas. Robust watertight manifold surface generation method for shapenet models. *arXiv preprint arXiv:1802.01698*, 2018. 4
- [14] A. Jain, T. Thormählen, T. Ritschel, and H.-P. Seidel. Exploring shape variations by 3d-model decomposition and part-based recombination. In *Computer Graphics Forum*, volume 31, pages 631–640. Wiley Online Library, 2012. 1
- [15] E. Kalogerakis, A. Hertzmann, and K. Singh. Learning 3D mesh segmentation and labeling. *ACM Transactions on Graphics (TOG)*, 29(4):102, 2010. 2
- [16] V. G. Kim, S. Chaudhuri, L. Guibas, and T. Funkhouser. Shape2pose: Human-centric shape analysis. *ACM Transactions on Graphics (TOG)*, 33(4):120, 2014. 1
- [17] R. Klokov and V. Lempitsky. Escape from cells: Deep kd-networks for the recognition of 3D point cloud models. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 863–872. IEEE, 2017. 1, 5
- [18] E. Kolve, R. Mottaghi, D. Gordon, Y. Zhu, A. Gupta, and A. Farhadi. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017. 1
- [19] P. Krähenbühl and V. Koltun. Parameter learning and convergent inference for dense random fields. In *International Conference on Machine Learning*, pages 513–521, 2013. 8
- [20] H. W. Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 7
- [21] T. Le and Y. Duan. PointGrid: A deep network for 3D shape understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9204–9214, 2018. 1, 5
- [22] J. Li, B. M. Chen, and G. H. Lee. SO-Net: Self-organizing network for point cloud analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9397–9406, 2018. 1, 5
- [23] J. Li, K. Xu, S. Chaudhuri, E. Yumer, H. Zhang, and L. Guibas. Grass: Generative recursive autoencoders for shape structures. *ACM Transactions on Graphics (TOG)*, 36(4):52, 2017. 1
- [24] Y. Li, R. Bu, M. Sun, and B. Chen. PointCNN: Convolution on  $\mathcal{X}$ -transformed points. *Advances in neural information processing systems (NIPS)*, 2018. 1, 5, 6, 13, 14
- [25] Z. Liu, W. T. Freeman, J. B. Tenenbaum, and J. Wu. Physical primitive decomposition. *arXiv preprint arXiv:1809.05070*, 2018. 1

- [26] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008. [1](#)
- [27] M. Ovsjanikov, W. Li, L. Guibas, and N. J. Mitra. Exploration of continuous variability in collections of 3d shapes. In *ACM Transactions on Graphics (TOG)*, volume 30, page 33. ACM, 2011. [1](#)
- [28] X. Puig, K. Ra, M. Boben, J. Li, T. Wang, S. Fidler, and A. Torralba. Virtualhome: Simulating household activities via programs. In *CVPR*, 2018. [1](#)
- [29] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, *IEEE*, volume 1, page 4, 2017. [1](#), [5](#), [6](#), [13](#), [14](#)
- [30] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems*, pages 5099–5108, 2017. [1](#), [5](#), [6](#), [7](#), [13](#), [14](#), [15](#)
- [31] H. Su, V. Jampani, D. Sun, S. Maji, E. Kalogerakis, M.-H. Yang, and J. Kautz. SplatNet: Sparse lattice networks for point cloud processing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2530–2539, 2018. [1](#), [5](#)
- [32] M. Sung, H. Su, R. Yu, and L. Guibas. Deep functional dictionaries: Learning consistent semantic structures on 3D models from functions. *Advances in neural information processing systems (NIPS)*, 2018. [8](#)
- [33] P.-S. Wang, Y. Liu, Y.-X. Guo, C.-Y. Sun, and X. Tong. OctNet: Octree-based convolutional neural networks for 3D shape analysis. *ACM Transactions on Graphics (TOG)*, 36(4):72, 2017. [1](#), [5](#)
- [34] W. Wang, R. Yu, Q. Huang, and U. Neumann. SGPN: Similarity group proposal network for 3D point cloud instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2569–2578, 2018. [8](#), [15](#)
- [35] X. Wang, B. Zhou, H. Fang, X. Chen, Q. Zhao, and K. Xu. Learning to group and label fine-grained shape components. *ACM Transactions on Graphics (SIGGRAPH Asia 2018)*, 37(6), 2018. [2](#)
- [36] Y. Wang, S. Asafi, O. Van Kaick, H. Zhang, D. Cohen-Or, and B. Chen. Active co-analysis of a set of shapes. *ACM Transactions on Graphics (TOG)*, 31(6):165, 2012. [2](#)
- [37] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon. Dynamic graph cnn for learning on point clouds. *arXiv preprint arXiv:1801.07829*, 2018. [1](#), [5](#)
- [38] Z. Wang and F. Lu. VoxSegNet: Volumetric CNNs for semantic part segmentation of 3D shapes. *arXiv preprint arXiv:1809.00226*, 2018. [1](#), [5](#)
- [39] Z. Wu, X. Wang, D. Lin, D. Lischinski, D. Cohen-Or, and H. Huang. Structure-aware generative network for 3d-shape modeling. *arXiv preprint arXiv:1808.03981*, 2018. [1](#)
- [40] Y. Xu, T. Fan, M. Xu, L. Zeng, and Y. Qiao. Spider-CNN: Deep learning on point sets with parameterized convolutional filters. *European Conference on Computer Vision (ECCV)*, 2018. [1](#), [5](#), [6](#), [13](#), [14](#)
- [41] C. Yan, D. Misra, A. Bennet, A. Walsman, Y. Bisk, and Y. Artzi. Chalet: Cornell house agent learning environment. *arXiv preprint arXiv:1801.07357*, 2018. [1](#)
- [42] L. Yi, L. Guibas, A. Hertzmann, V. G. Kim, H. Su, and E. Yumer. Learning hierarchical shape segmentation and labeling from online repositories. *ACM Transactions on Graphics (Proc. SIGGRAPH Asia)*, 2017. [1](#), [2](#)
- [43] L. Yi, V. G. Kim, D. Ceylan, I. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, L. Guibas, et al. A scalable active framework for region annotation in 3D shape collections. *ACM Transactions on Graphics (TOG)*, 35(6):210, 2016. [1](#), [2](#), [4](#), [5](#), [6](#), [12](#), [13](#), [14](#)
- [44] L. Yi, H. Su, X. Guo, and L. J. Guibas. SyncSpecCNN: Synchronized spectral CNN for 3D shape segmentation. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6584–6592, 2017. [1](#), [5](#)
- [45] Y. Zhu, D. Gordon, E. Kolve, D. Fox, L. Fei-Fei, A. Gupta, R. Mottaghi, and A. Farhadi. Visual semantic planning using deep successor representations. *arXiv preprint ArXiv:1705.08080*, pages 1–13, 2017. [1](#)

## A. Overview

This document provides additional dataset visualization and statistics (Sec [B](#)), hierarchical template design details and visualization (Sec [C](#)), and the architectures and training details for the three shape segmentation tasks (Sec [D](#)), to the main paper.

## B. More Dataset Visualization and Statistics

We present more visualization and statistics over the proposed PartNet dataset.

### B.1. More Fine-grained Segmentation Visualization

Figure [13](#) and [14](#) show more visualization for fine-grained instance-level segmentation annotations in PartNet. We observe the complexity of the annotated segmentation and the heterogeneous variation of shapes within each object category.

### B.2. More Hierarchical Segmentation Visualization

Figure [15](#), [16](#) and [17](#) show more visualization for example hierarchical instance-level segmentation annotations in PartNet. We visualize the tree-structure of the hierarchical segmentation annotation with the 2D part renderings associated to the tree nodes.

### B.3. Shape Statistics

We report the statistics for the number of annotations, unique shapes and shapes that we collect multiple human annotations in Figure [9](#).

### B.4. Part Statistics

We report the statistics for the number of part semantics for each object category in Figure [10](#). We also present the

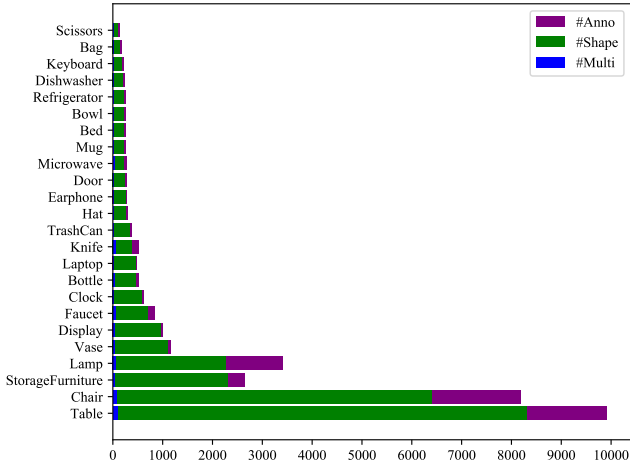


Figure 9. **PartNet shape statistics.** We report the statistics for the number of annotations, unique shapes and shapes that we collect multiple human annotations.

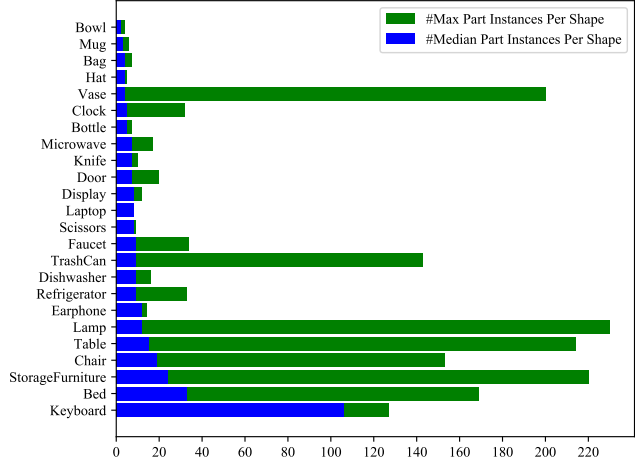


Figure 11. **PartNet part instance statistics.** We report the statistics for the maximum and median number of part instances per shape for each object category.

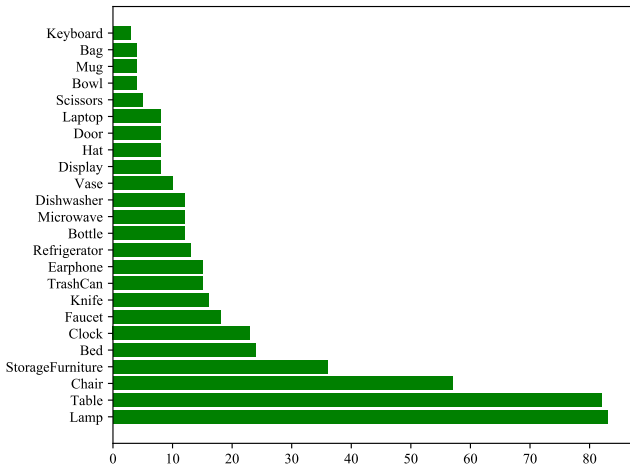


Figure 10. **PartNet part semantics statistics.** We report the statistics for the number of part semantics for each object category.

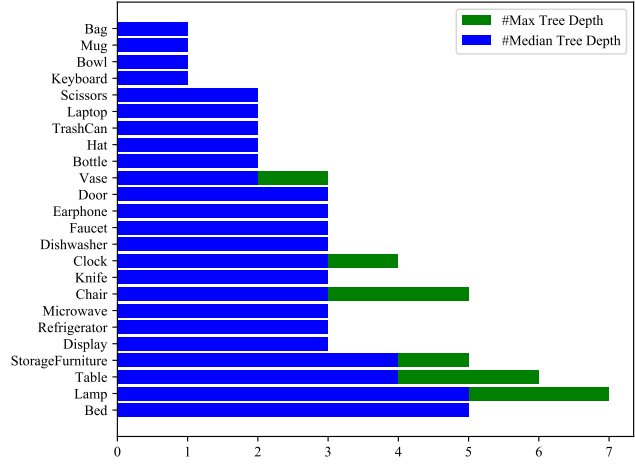


Figure 12. **PartNet tree depth statistics.** We report the statistics for the maximum and median tree depth for each object category.

statistics for the maximum and median number of part instances per shape for each object category in Figure 11. We report the statistics for the maximum and median tree depth for each object category in Figure 12.

## C. More Template Design Details and Visualization

We provide more details and visualization for the expert-defined hierarchical templates to guide the hierarchical segmentation annotation and the template refinement procedure to resolve annotation inconsistencies.

### C.1. Template Design Details

We design templates according to the rules of thumb that we describe in the main paper. We also consulted many on-

line references<sup>5</sup> that describe object parts (often for manufacturing and assembly) and previous works that relate language to the shapes [4] as guides for the design of our template. To ensure that our templates cover most of the shape variations and part semantics of each object category, we generated a t-SNE [26] visualization for the entire shape space to study the shape variation. We trained an auto-encoder based on the shape geometry within each object category to obtain shape embeddings for the t-SNE visualization.

Although we try to cover the most common part semantics in our templates, it is still not easy to cover all possible object parts. Thus, we allow annotators to deviate from the templates and define their own parts and segmentation structures. Among all the annotated part instances, 1.3% of

<sup>5</sup>E.g. <http://www.props.eric-hart.com/resources/parts-of-a-chair/>.

	Avg	Bag	Bed	Bott	Bowl	Chair	Clock	Dish	Disp	Door	Ear	Fauc	Hat	Key	Knife	Lamp	Lap	Micro	Mug	Frid	Scis	Stora	Table	Trash	Vase
<b>O</b> <sub>avg</sub>	69.8	54.8	70.0	87.5	87.7	59.0	62.1	67.3	85.2	64.4	74.1	69.9	86.8	77.0	75.0	44.6	61.6	71.0	91.0	65.3	88.7	68.0	40.1	51.4	72.6
<b>O</b> <sub>std</sub>	19.0	29.3	17.4	6.9	8.1	26.7	24.6	19.1	17.8	15.2	15.6	17.6	9.5	17.9	14.6	29.0	27.1	19.3	10.6	27.7	9.0	21.3	29.3	23.3	19.5
<b>R</b> <sub>avg</sub>	83.3	82.1	76.2	89.3	91.7	77.8	91.1	81.5	94.0	77.0	83.0	84.7	89.3	89.6	77.8	72.7	78.3	84.4	91.7	85.1	90.2	77.1	71.4	71.0	92.3
<b>R</b> <sub>std</sub>	10.4	11.1	9.2	6.0	7.4	15.2	7.0	7.2	2.8	11.2	13.5	8.6	10.3	3.5	14.9	17.3	14.6	9.1	9.7	6.2	8.6	12.8	22.6	13.2	7.5

Table 6. **The average confusion scores and the standard deviations for multiple annotations (%)**. We report the average confusion scores and the standard deviations by calculating over the entries on the diagonal of the confusion matrix for each object category using the small subset of shapes that we collect multiple human annotations. Rows **O** and **R** respectively refer to the scores before and after the template refinement process.

them are defined by the annotators. In the raw annotation, 13.1% of shapes contained user-defined part labels.

Our analysis shows that our template designs are able to cover most of the ShapeNet [3] shapes. Of the 27,260 shapes we collected in total, our annotators successfully labeled 26,671 of them, giving our templates a coverage rate of at least 97.8% for ShapeNet shapes. While template coverage is a potential issue, the remaining 2.2% were not annotated mainly due to other issues such as poor mesh quality, classification error, error during mesh splitting, *etc.*

We design hierarchical templates that cover both the coarse-level part semantics and fine-grained part details down to the primitive level, *e.g. chair back vertical bar and bed base surface panel*. Most primitive-level parts are atomic such that they are very unlikely to be further divided for end applications. If an application requires different segmentation hierarchy or level of segmentation than the ones we already provide in our template, developers and researchers can try to build up their own segmentation based upon the atomic primitives we obtain in PartNet.

Moreover, we try our best to make the shared part concepts among different shapes and even different object categories share the same part labels. For example, we use the part label *leg* for *table, chair, lamp base, etc.* and the part label *wheel* for both *chair swivel base wheel* and *refrigerator base wheel*. Such part concept sharing provides rich part correspondences within a specific object category and across multiple object categories.

## C.2. Template Refinement Details

Fine-grained shape segmentation is challenging to annotate due to the subtle concept gaps among similar part semantics. Even though we provide detailed textual and visual explanation for our pre-defined parts, we still observe some annotation inconsistencies across multiple annotators. To quantitatively diagnose such issues, we reserve a small subset of shapes for which we collect multiple human annotations. Then, we compute the confusion scores among the predefined parts across the multiple annotations and conduct careful template refinement to reduce the part ambiguity.

There are primarily three sources of such inconsistencies: boundary ambiguity, granularity ambiguity and part labeling ambiguity. Boundary ambiguity refers to the un-

clear boundary between two parts, which is also commonly seen in previous works [5, 43]. For example, the boundary between the bottle neck and the bottle body is not that clear for wine bottles. Granularity ambiguity means that different annotators have different understanding about the segmentation granularity of the defined parts. One example is that, for a curvy and continuous chair arm, one can regard it as a whole piece or imagine the separation of armrest and arm support. The most common type of ambiguity in our dataset is the part labeling ambiguity. The fine-grained part concepts, though intended to be different category-wise, may apply to the same part on a given object. For example, a connecting structure between the seat and the base of a chair can be considered as chair seat support or chair base connector.

We study the mutual human agreement on the multiple annotation subset. We consider the parts defined at the leaf node level of segmentation on the hierarchy and compute the confusion matrix across multiple human annotations<sup>6</sup>. The ideal confusion matrix should be close to the diagonal matrix without any part-level ambiguity. In our analysis, we observe human disagreement among some of our initial part definitions. To address the ambiguity, we either merge two similar concepts with high confusion scores or remove the hard-to-distinguish parts from evaluation. For example, we find our annotators often mix up the annotation for regular tables and desks due to the similarity in the two concepts. Thus, we merge the desk subtype into the regular table subtype to address this issue. In other cases, some small parts such as the buttons on the displays are very tricky to segment out from the main display frame. Since they may not be reliably segmented out, we decided to remove such unclear segmentation from evaluation.

Table 6 compares the annotation consistency before and after the template refinement process. We compute the confusion matrices at the most fine-grained segmentation level. After the template refinement, the data consistency score is 83.3% on average, having 13.5% improvement over the raw annotation. The template refinement process improves the annotation consistency by a clear margin. This also reflects the complexity of the task in terms of annotating fine-grained part concepts. Future works may investigate how to

<sup>6</sup>We consider the entire path labels as histories to the leaf nodes when computing the confusion matrix.

	Avg	Bag	Bed	Bott	Bowl	Chair	Clock	Dish	Disp	Door	Ear	Fauc	Hat	Key	Knife	Lamp	Lap	Micro	Mug	Frid	Scis	Stora	Table	Trash	Vase
<b>P1</b>	71.8	59.3	39.6	81.0	78.5	81.8	67.1	78.9	88.2	71.1	68.0	67.5	58.5	65.6	66.5	46.5	96.5	75.0	84.2	<b>79.6</b>	86.5	55.9	85.6	66.7	76.3
<b>P2</b>	50.1	—	21.3	—	—	52.4	—	60.0	—	47.1	—	—	—	—	—	43.5	—	<b>64.3</b>	—	<b>63.9</b>	—	48.8	50.0	—	—
<b>P3</b>	48.2	—	13.0	55.3	—	44.8	37.8	55.2	79.0	38.8	47.5	55.5	—	—	40.0	34.7	—	54.5	—	<b>53.2</b>	—	47.4	42.5	46.4	74.0
<b>Avg</b>	63.4	59.3	24.6	68.2	78.5	59.7	52.4	64.7	83.6	52.3	57.8	61.5	58.5	65.6	53.2	41.6	96.5	64.6	84.2	<b>65.6</b>	86.5	50.7	59.4	56.5	75.2
<b>PP1</b>	<b>76.8</b>	72.7	54.7	85.8	78.5	<b>84.5</b>	<b>74.1</b>	<b>81.9</b>	90.7	<b>73.5</b>	<b>77.8</b>	<b>73.6</b>	64.2	62.5	<b>75.0</b>	<b>65.5</b>	96.6	80.3	90.9	72.1	87.5	61.2	<b>86.7</b>	71.5	<b>81.4</b>
<b>PP2</b>	<b>54.7</b>	—	34.8	—	—	<b>54.9</b>	—	60.6	—	<b>57.0</b>	—	—	—	—	—	<b>56.8</b>	—	63.0	—	58.4	—	52.9	<b>53.6</b>	—	—
<b>PP3</b>	<b>53.4</b>	—	25.1	<b>61.0</b>	—	<b>49.6</b>	<b>46.1</b>	52.5	81.0	<b>48.0</b>	<b>56.1</b>	60.4	—	—	<b>49.1</b>	<b>46.0</b>	—	54.3	—	50.7	—	50.6	<b>47.0</b>	<b>54.7</b>	<b>75.1</b>
<b>Avg</b>	<b>68.1</b>	72.7	38.2	<b>73.4</b>	78.5	<b>63.0</b>	<b>60.1</b>	65.0	85.8	<b>59.5</b>	<b>67.0</b>	<b>67.0</b>	64.2	62.5	<b>62.0</b>	<b>56.1</b>	96.6	65.9	90.9	60.4	87.5	54.9	<b>62.4</b>	<b>63.1</b>	<b>78.2</b>
<b>S1</b>	73.9	<b>72.9</b>	<b>55.9</b>	<b>86.1</b>	<b>83.4</b>	83.8	72.1	73.3	90.4	60.4	70.6	71.5	<b>71.6</b>	64.6	42.1	59.1	<b>97.1</b>	78.6	91.6	68.7	77.0	<b>64.2</b>	83.8	<b>74.4</b>	79.5
<b>S2</b>	53.3	—	<b>37.8</b>	—	—	53.6	—	<b>65.3</b>	—	55.0	—	—	—	—	—	41.4	—	62.1	—	62.6	—	49.8	51.7	—	—
<b>S3</b>	48.0	—	27.2	52.8	—	44.7	44.2	51.1	77.2	40.7	47.5	53.7	—	—	27.3	35.7	—	54.4	—	52.4	—	<b>53.1</b>	43.3	48.0	62.3
<b>Avg</b>	65.1	<b>72.9</b>	40.3	69.4	<b>83.4</b>	60.7	58.1	63.2	83.8	52.0	59.0	62.6	<b>71.6</b>	64.6	34.7	45.4	<b>97.1</b>	65.0	91.6	61.2	77.0	<b>55.7</b>	59.6	61.2	70.9
<b>C1</b>	75.5	72.0	55.3	83.6	75.0	83.9	65.6	81.8	<b>91.9</b>	68.1	74.5	71.1	66.8	<b>70.4</b>	68.1	55.6	<b>97.1</b>	<b>83.1</b>	<b>92.7</b>	78.9	<b>92.6</b>	58.8	85.5	67.7	71.8
<b>C2</b>	52.1	—	36.6	—	—	52.9	—	63.4	—	54.9	—	—	—	—	—	42.4	—	64.1	—	57.7	—	<b>54.4</b>	42.7	—	—
<b>C3</b>	49.6	—	<b>29.1</b>	58.7	—	47.7	36.2	<b>55.3</b>	<b>81.5</b>	40.4	55.8	<b>60.7</b>	—	—	26.4	34.4	—	<b>58.7</b>	—	50.8	—	52.3	37.4	50.8	67.0
<b>Avg</b>	66.3	72.0	<b>40.3</b>	71.2	75.0	61.5	50.9	<b>66.8</b>	<b>86.7</b>	54.5	65.2	65.9	66.8	<b>70.4</b>	47.2	44.1	<b>97.1</b>	<b>68.6</b>	<b>92.7</b>	62.5	<b>92.6</b>	55.2	55.2	59.2	69.4

Table 7. **Fine-grained semantic segmentation results (shape mIoU %)**. Algorithm **P**, **P<sup>+</sup>**, **S** and **C** refer to PointNet [29], PointNet++ [30], SpiderCNN [40] and PointCNN [24], respectively. The number **1**, **2** and **3** refer to the three levels of segmentation: coarse-, middle- and fine-grained. We put short lines for the levels that are not defined.

further design better templates with less part ambiguities.

### C.3. More Visualization of Hierarchical Templates

Figure 18, 19 and 20 show more visualization for the expert-designed hierarchical templates after resolving the data inconsistency and conducting template refinements. We show the lamp template in the main paper.

## D. Tasks and Benchmarks

In this section, we provide more details about the architectures and training details for the benchmark algorithms. We also present additional evaluation metrics, shape mean Intersection-over-Union (shape mIoU) and shape mean Average-Precision (Shape mAP), and report the quantitative results using these metrics.

### D.1. Fine-grained Semantic Segmentation

**More Architecture and Training Details** We follow the default architectures and training hyper-parameters used in the original papers: PointNet [29], PointNet++ [30], SpiderCNN [40] and PointCNN [24], except the following few modifications:

- Instead of training one network for all object categories as done in the four prior works, we train separate networks for each object category at each segmentation level. This is mainly to handle the increase in the number of parts for fine-grained part segmentation. Originally, there are only 50 parts for all 16 object categories using the coarse ShapeNet Part dataset [43]. Now, using PartNet, there could be 480 different part semantics in total. Also, due to the data imbalance among different object categories, training a single network may overfit to the big categories.

- We change the input point cloud size to 10,000. The original papers usually sample 1,000, 2,000 or 4,000 points and input to the networks. PartNet suggests to use at least 10,000 to guarantee enough point sampling over small fine-grained parts, *e.g.* a door handle, or a small button.
- We reduce the batch sizes for training the networks if necessary. Since we use point cloud size 10,000, to fit the training in NVIDIA TITAN XP GPU 12G memory, we need to adjust the training batch size accordingly. For PointNet [29], PointNet++ [30], SpiderCNN [40] and PointCNN [24], we use batch size of 24, 24, 2 and 4 respectively.
- We only input 3D coordinates as inputs to all the networks for fair comparison. Although the 3D CAD models in ShapeNet [3] usually provide additional features, *e.g.* opacity, point normals, textures and material information, there is no guarantee for the quality of such information. Thus, we choose not to use them as the inputs. Also, only using pure geometry potentially increase the network generalizability to unseen objects or real scans [29]. PointNet++ [30] and SpiderCNN [40] by defaults take advantage of the point normals as additional inputs. In this paper, we remove such inputs to the networks. However, point normals can be estimated from the point clouds. We leave this as a future work.

**Shape mIoU Metric and Results** We introduce the shape mean Intersect-over-Union (Shape mIoU) evaluation metric as a secondary metric to the Part-category mIoU metric in the main paper. Shape mIoU metric considers shapes as evaluation units and measures how an algorithm segment

	Avg	Bed	Bott	Chair	Clock	Dish	Disp	Door	Ear	Fauc	Knife	Lamp	Micro	Frid	Stora	Table	Trash	Vase
<b>Bottom-Up</b>	65.9	42.0	74.3	63.8	64.1	<b>66.3</b>	84.2	61.4	70.0	<b>74.2</b>	67.1	62.7	<b>63.0</b>	60.8	57.8	65.7	62.8	80.9
<b>Top-Down</b>	65.9	42.0	73.7	62.3	<b>65.5</b>	64.0	85.5	63.1	71.1	73.5	<b>68.8</b>	63.3	62.7	58.8	57.6	66.2	63.0	79.3
<b>Ensemble</b>	<b>66.6</b>	<b>42.9</b>	<b>74.4</b>	<b>64.3</b>	<b>65.5</b>	62.7	<b>85.8</b>	<b>63.7</b>	<b>71.7</b>	74.0	66.7	<b>63.4</b>	61.9	<b>61.5</b>	<b>60.6</b>	<b>67.5</b>	<b>64.0</b>	<b>82.2</b>

Table 8. **Hierarchical segmentation results (shape mIoU %)**. We present the hierarchical segmentation performances for three baseline methods: bottom-up, top-down and ensemble. We conduct experiments on 17 out of 24 categories with tree depth bigger than 1.

an average shape in the object category. In contrast, Part-category mIoU reports the average performance over all part semantics and indicates how an algorithm performs for any given part category.

Shape mIoU is widely used on ShapeNet Part dataset [43] for 3D shape coarse semantic segmentation [29, 30, 40, 24]. We propose a slightly different version for fine-grained semantic segmentation. For each test shape, we first compute the IoU for each part semantics that either presents in the ground-truth or is predicted by the algorithm, and then we calculate the mean IoU for this shape. We remove the ground-truth unlabeled points from the evaluation. Finally, we calculate the Shape mIoU by averaging mIoU over all test shape instances.

We benchmark the four algorithms using Shape mIoU in Table 7. Besides the Shape mIoU scores for each object category at each segmentation level, we also report the average across levels for each object categories and further calculate the average over all object categories.

We observe that PointNet++ [30] achieves the best performance using the Shape mIoU metric, while PointCNN [24] performs the best using the Part-category mIoU metric. The Part-category mIoU metric considers all part semantics equally while the Shape mIoU metric considers all shapes equally. We observe an unbalanced counts for different part semantics in most object categories, *e.g.* there are much more chair legs than chair wheels. To achieve good numbers on Part-category mIoU, a segmentation algorithm needs to perform equally well on both frequent parts and rare parts, while the Shape mIoU metric bias over the frequently observed parts.

## D.2. Hierarchical Semantic Segmentation

We describe the architecture and training details for the three baseline methods we propose for hierarchical semantic segmentation in the main paper. All three methods use PointNet++ [30] segmentation network as the network backbone. The difference of the three methods is mainly on the training and inference strategies to enforce the tree knowledge to the final prediction.

**The Bottom-up Method** The bottom-up method learns a network to perform segmentation at the most fine-grained leaf part semantics. We use the PointNet++ [30] segmentation network with a softmax activation layer as the network architecture. At inference time, we use the ground-truth tree

hierarchy to gather the prediction for the parent nodes. The parent node prediction is the sum of all its children node predictions. Even though we only train for the leaf node parts, the parent history is implicitly encoded. For example, we define vertical bars for both chair back and chair arm, but they are two different leaf node parts: *chair back vertical bar* and *chair arm vertical bar*.

In the ground-truth annotation, all the points in the point cloud belong to the root node. Each point is assigned a path of labels from the root node to some intermediate node in the tree. The paths for most points are all the way down to the leaf levels while some points may not. For example, a point on a bed blanket (removed from evaluation since it cannot be distinguished without color information) may be assigned with labels  $\{bed, bed\ unit, sleeping\ area\}$  in the ground-truth annotation. The part *sleeping area* is not a leaf part. For such cases, we introduce an additional leaf node *other* for each parent node in the tree and consider them in the training.

**The Top-down Method** The top-down method learns a multi-labeling task for all the part semantics in the tree, considering both the leaf nodes and the parent nodes. Compared to the bottom-up method, the top-down method takes advantage of the tree structures at training time.

Assuming there are  $T$  tree nodes in the hierarchy, we train a PointNet++ [30] segmentation network for a  $T$ -way classification for each point. We apply a softmax activation layer to enforce label mutual exclusiveness. For a point with the ground-truth labels  $\{y_1, y_2, y_3\}$  and prediction softmax scores  $\{s_i | i = 1, 2, \dots, T\}$ , we train the labels using a multi-labeling cross-entropy loss

$$L = -\log(s_{y_1}) - \log(s_{y_2}) - \log(s_{y_3}) \quad (1)$$

to increase the values of all the three label predictions over the rest labels.

**The Ensemble Method** The ensemble method trains multiple neural networks at different levels of segmentation as defined in the fine-grained semantic segmentation task. The key idea is that conducting segmentation at the coarse-, middle- and fine-grained levels separately may learn different features that work the best at each level. Compared to the bottom-up method that we only train at the most fine-grained level, additional signal at the coarse level helps distinguish the coarse-level part semantics more easily. For

	Avg	Bag	Bed	Bott	Bowl	Chair	Clock	Dish	Disp	Door	Ear	Fauc	Hat	Key	Knife	Lamp	Lap	Micro	Mug	Frid	Scis	Stora	Table	Trash	Vase
<b>S1</b>	72.5	62.8	38.7	76.7	83.2	91.5	41.5	<b>81.4</b>	91.3	71.2	81.4	82.2	71.9	23.2	<b>78.0</b>	60.3	<b>100</b>	76.2	94.3	60.6	74.9	55.0	80.1	76.1	87.1
<b>S2</b>	50.2	—	22.7	—	—	51.1	—	<b>78.7</b>	—	43.3	—	—	—	—	—	49.1	—	68.6	—	42.9	—	51.9	43.7	—	—
<b>S3</b>	50.2	—	17.5	66.5	—	42.3	40.7	59.3	83.9	29.0	60.2	61.6	—	—	55.0	37.6	—	53.7	—	30.6	—	45.1	37.8	50.0	82.0
<b>Avg</b>	64.2	62.8	26.3	71.6	83.2	61.6	41.1	73.1	87.6	47.8	70.8	71.9	71.9	23.2	66.5	49.0	<b>100</b>	66.2	94.3	44.7	74.9	50.7	53.8	63.0	84.6
<b>O1</b>	<b>80.3</b>	<b>78.4</b>	<b>62.2</b>	<b>80.8</b>	<b>83.8</b>	<b>94.9</b>	<b>74.6</b>	<b>81.4</b>	<b>94.3</b>	<b>76.1</b>	<b>87.1</b>	<b>86.5</b>	<b>77.8</b>	<b>44.5</b>	76.6	<b>65.0</b>	<b>100</b>	<b>79.5</b>	<b>95.3</b>	<b>79.0</b>	<b>87.6</b>	<b>62.7</b>	<b>88.1</b>	<b>82.3</b>	<b>89.0</b>
<b>O2</b>	<b>60.5</b>	—	<b>29.4</b>	—	—	<b>64.7</b>	—	75.4	—	<b>61.1</b>	—	—	—	—	—	<b>56.8</b>	—	<b>78.2</b>	—	<b>61.7</b>	—	<b>57.4</b>	<b>59.4</b>	—	—
<b>O3</b>	<b>57.7</b>	—	<b>22.1</b>	<b>68.3</b>	—	<b>58.4</b>	<b>53.7</b>	<b>67.5</b>	<b>84.8</b>	<b>38.0</b>	<b>62.4</b>	<b>66.8</b>	—	—	<b>63.5</b>	<b>45.8</b>	—	<b>54.0</b>	—	<b>45.0</b>	—	<b>52.6</b>	<b>52.5</b>	<b>58.7</b>	<b>86.4</b>
<b>Avg</b>	<b>72.2</b>	<b>78.4</b>	<b>37.9</b>	<b>74.6</b>	<b>83.8</b>	<b>72.7</b>	<b>64.2</b>	<b>74.8</b>	<b>89.5</b>	<b>58.4</b>	<b>74.8</b>	<b>76.6</b>	<b>77.8</b>	<b>44.5</b>	<b>70.1</b>	<b>55.8</b>	<b>100</b>	<b>70.6</b>	<b>95.3</b>	<b>61.9</b>	<b>87.6</b>	<b>57.6</b>	<b>66.7</b>	<b>70.5</b>	<b>87.7</b>

Table 9. Instance segmentation results (shape mAP %, IoU threshold 0.5). Algorithm **S** and **O** refer to SGPN [34] and our proposed method respectively. The number **1**, **2** and **3** refer to the three levels of segmentation: coarse-, middle- and fine-grained. We put short lines for the levels that are not defined.

example, the local geometric features for both chair back vertical bars and chair arm vertical bars may be very similar, but the coarse-level semantics may distinguish chair backs and chair arms better.

During the training, we train 2~3 networks at multiple levels of segmentation. At the inference time, we perform a joint inference considering the prediction scores from all the networks. We use a path-voting strategy: for each path from the root node to the leaf node, we calculate the average log-likelihood over the network prediction scores after applying the softmax activations, and select the path with the highest score as the joint label predictions.

**Shape mIoU Metric and Results** Similar to Sec D.1, we define Shape mIoU for hierarchical segmentation. The mIoU for each shape is calculated over the part semantics in the entire hierarchical template that are either predicted by the network or included in the ground-truth. The unrelated parts are not taken into consideration. Table 8 shows the quantitative evaluation for the three baseline methods. We observe similar performance for the three methods, with the ensemble method works slightly better.

### D.3. Instance Segmentation

**More Architecture and Training Details** To train our proposed method, we use batch size 32, learning rate 0.001, and the default batch normalization settings used in the PointNet++ [30].

For SGPN [34], we use two-stage training as suggested by the authors of [34]. We first pretrain the PointNet++ semantic segmentation branch using batch size 32 and learning rate 0.001, with the default batch normalization as in PointNet++. And then, we jointly train for the semantic segmentation, similarity score matrix and confidence scores with batch size 1 and learning rate 0.0001. As suggested in the original SGPN paper, for the first five epochs of the joint training, we only turn on the loss for training the similarity scores matrix. The rest training epochs are done with the full losses switched on. We have to use batch size 1 because the input point cloud has the size of 10,000 and thus the similarity score matrix forms a  $10,000 \times 10,000$  ma-

trix, which occupies too much GPU memory. Our proposed method is more memory-efficient, compared to SGPN. We also observe that our training is much faster than SPGN. We train all the networks until convergence.

**Shape mAP Metric and Results** We define Shape mean Average-Precision (Shape mAP) metric as a secondary metric to the Part-category mAP metric in the main paper. Similar to the Shape mIoU scores we use in Sec D.1 and D.2, Shape mAP reports the part instance segmentation performance on an average shape in a object category. It averages across the test shapes, instead of averaging across all different part semantics, as benchmarked by Part-category mAP in the main paper.

To calculate Shape mAP for a test shape, we consider the AP for the part semantics that occur either in the ground-truth or the prediction for the given shape and compute their average as the mean AP score. Then, we average the mAP across all test shapes within a object category. Table 9 reports the part instance segmentation performance under the Shape mAP scores. We see a clear performance improvement of the proposed method over SGPN.

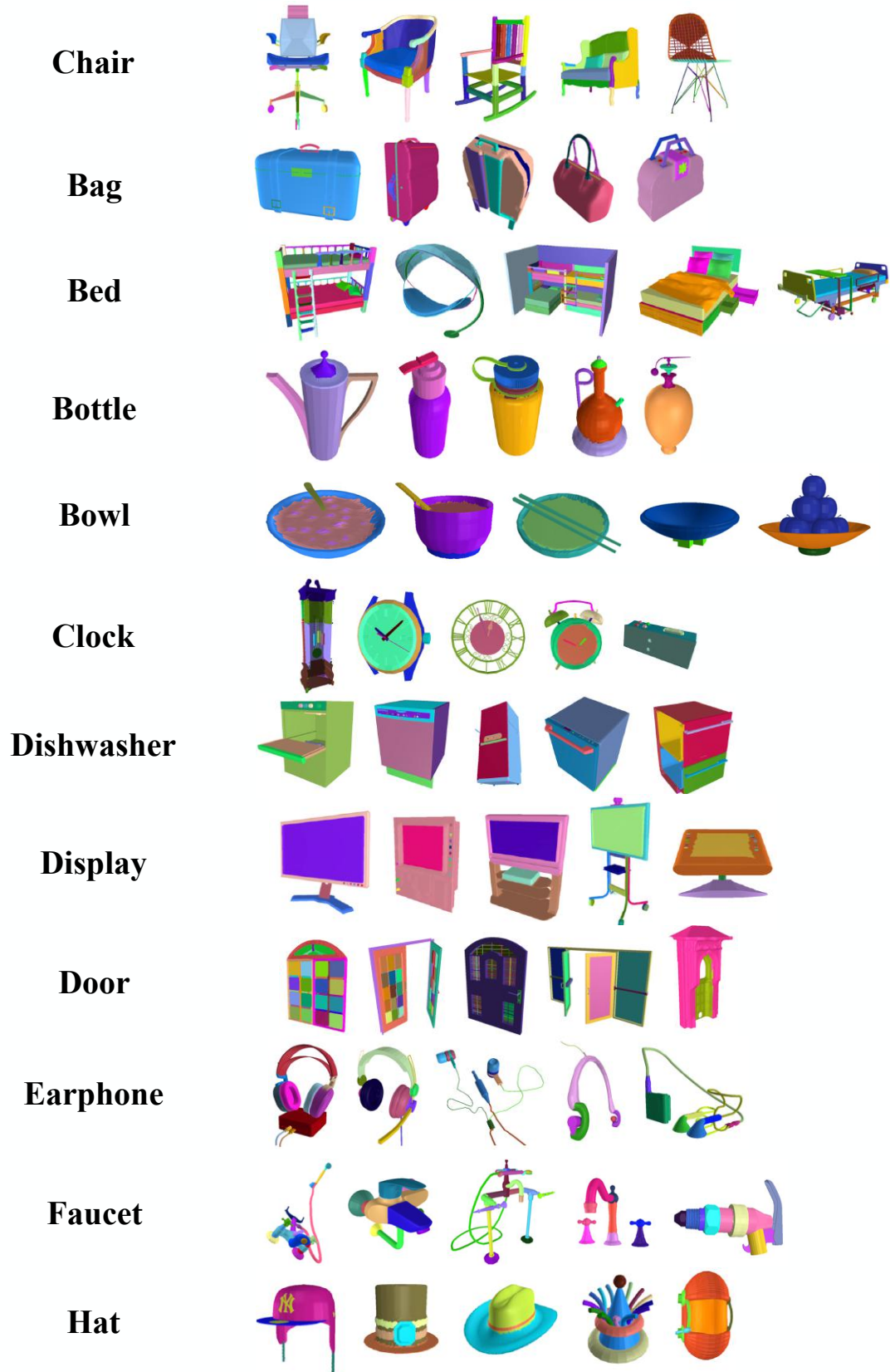


Figure 13. **Fine-grained instance-level segmentation visualization (1/2)**. We present visualization for example fine-grained instance-level segmentation annotations for chair, bag, bed, bottle, bowl, clock, dishwasher, display, door, earphone, faucet, and hat.



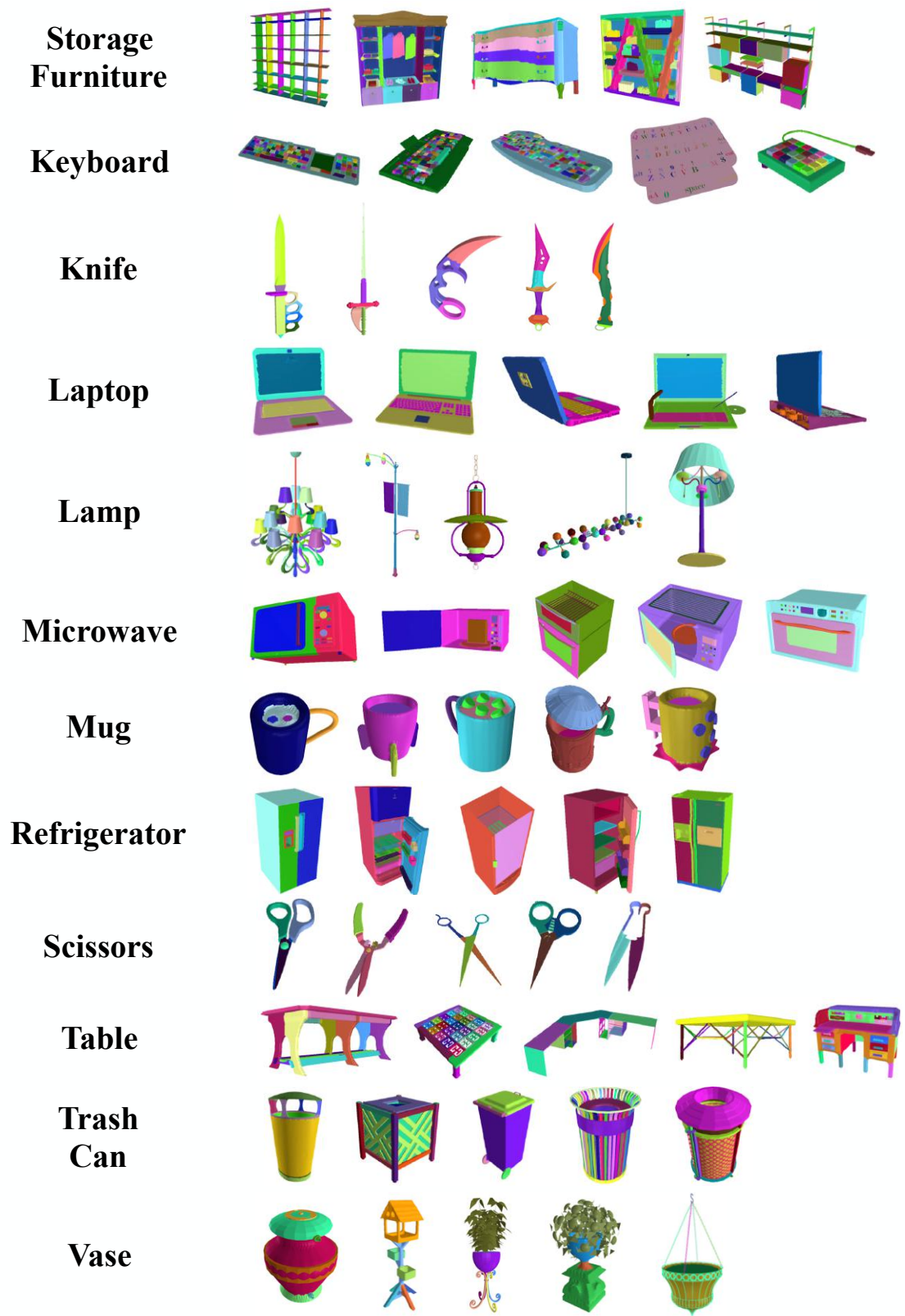


Figure 14. **Fine-grained instance-level segmentation visualization (2/2).** We present visualization for example fine-grained instance-level segmentation annotations for storage furniture, keyboard, knife, laptop, lamp, microwave, mug, refrigerator, scissors, table, trash can, and vase.

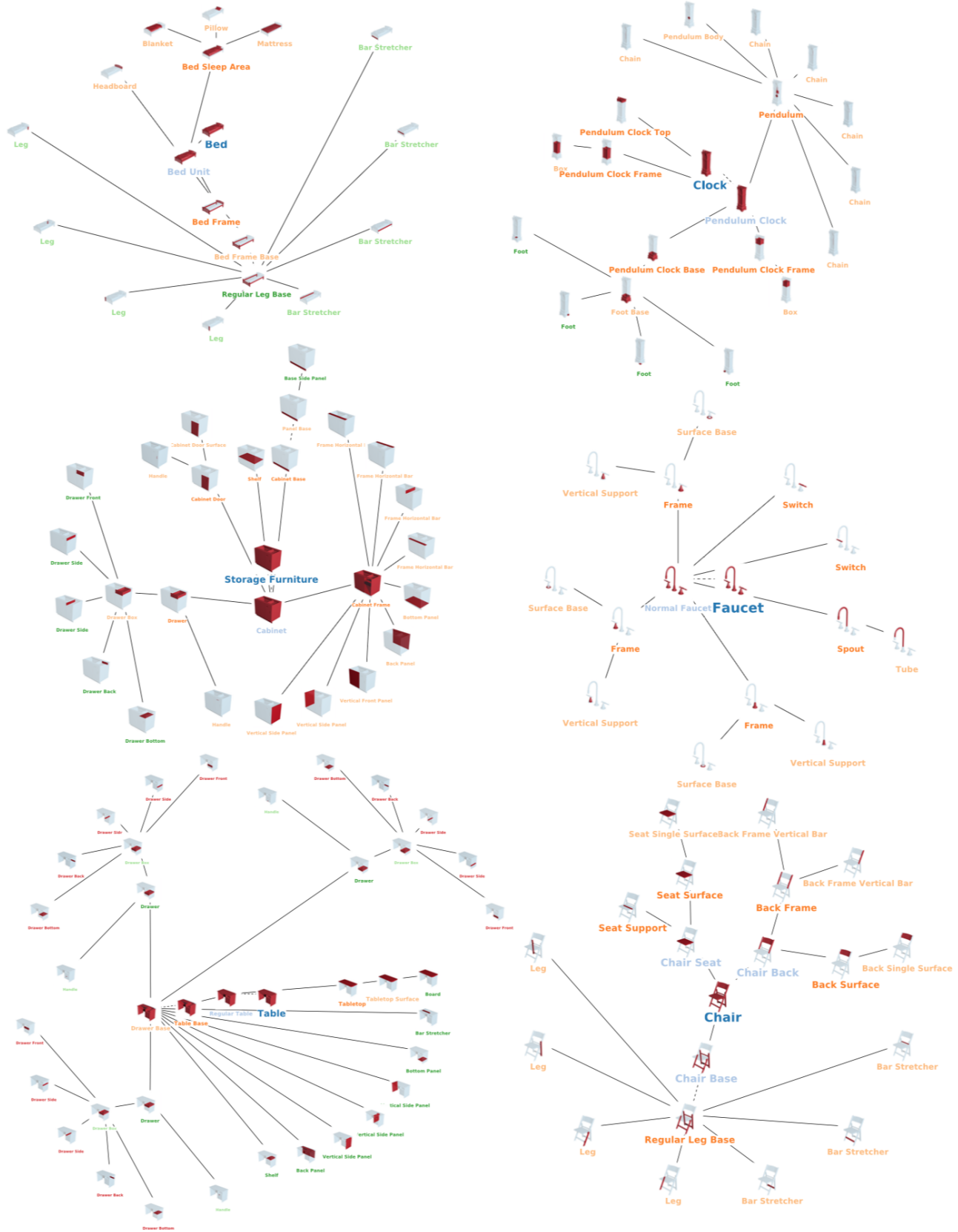


Figure 15. **Hierarchical instance-level segmentation visualization (1/3)**. We present visualization for example hierarchical instance-level segmentation annotations for bed, clock, storage furniture, faucet, table, and chair. The lamp examples are shown in the main paper. The And-nodes are drawn in solid lines and Or-nodes in dash lines.



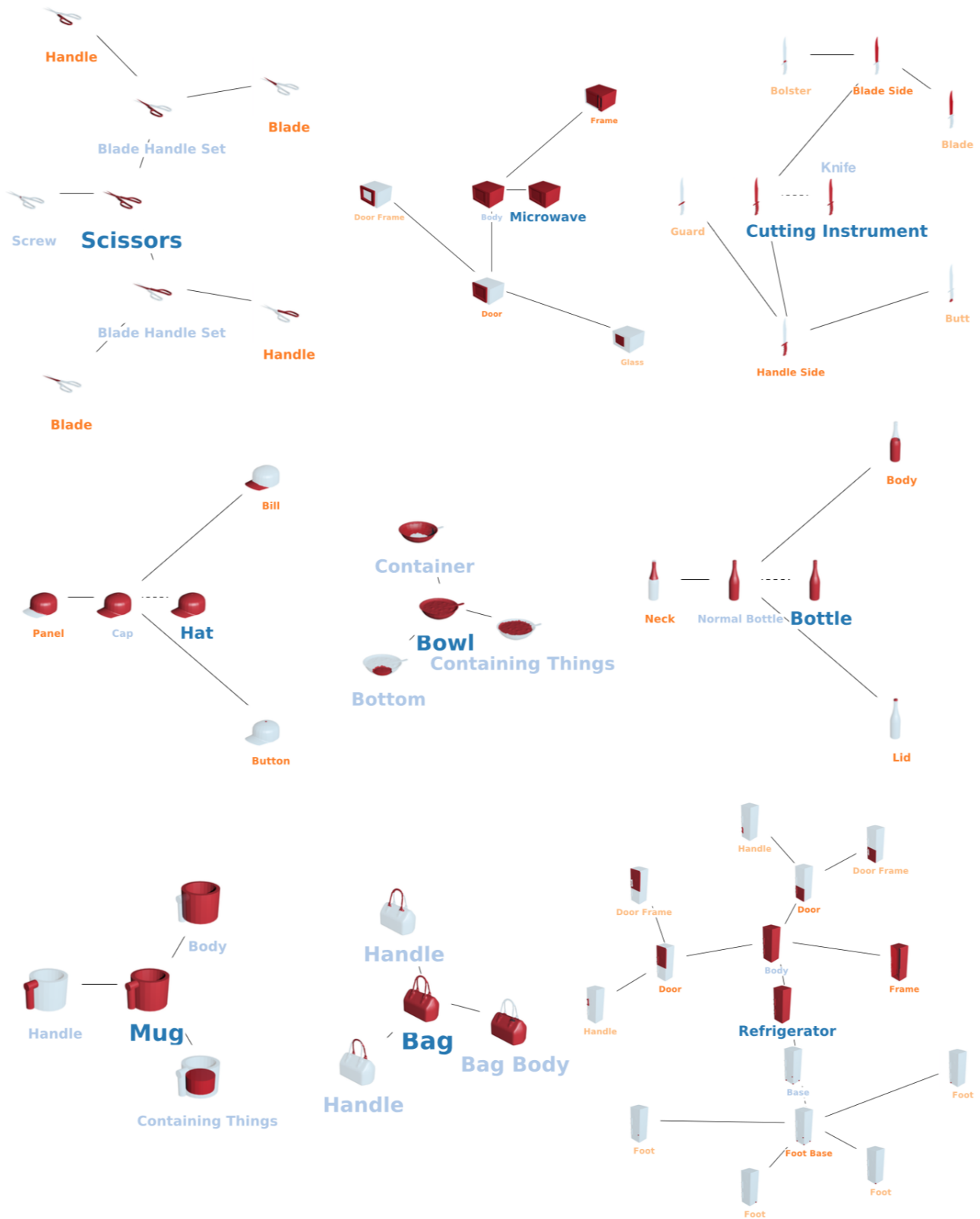


Figure 17. **Hierarchical instance-level segmentation visualization (3/3)**. We present visualization for example hierarchical instance-level segmentation annotations for scissors, microwave, knife (cutting instrument), hat, bowl, bottle, mug, bag, and refrigerator. The lamp examples are shown in the main paper. The And-nodes are drawn in solid lines and Or-nodes in dash lines.



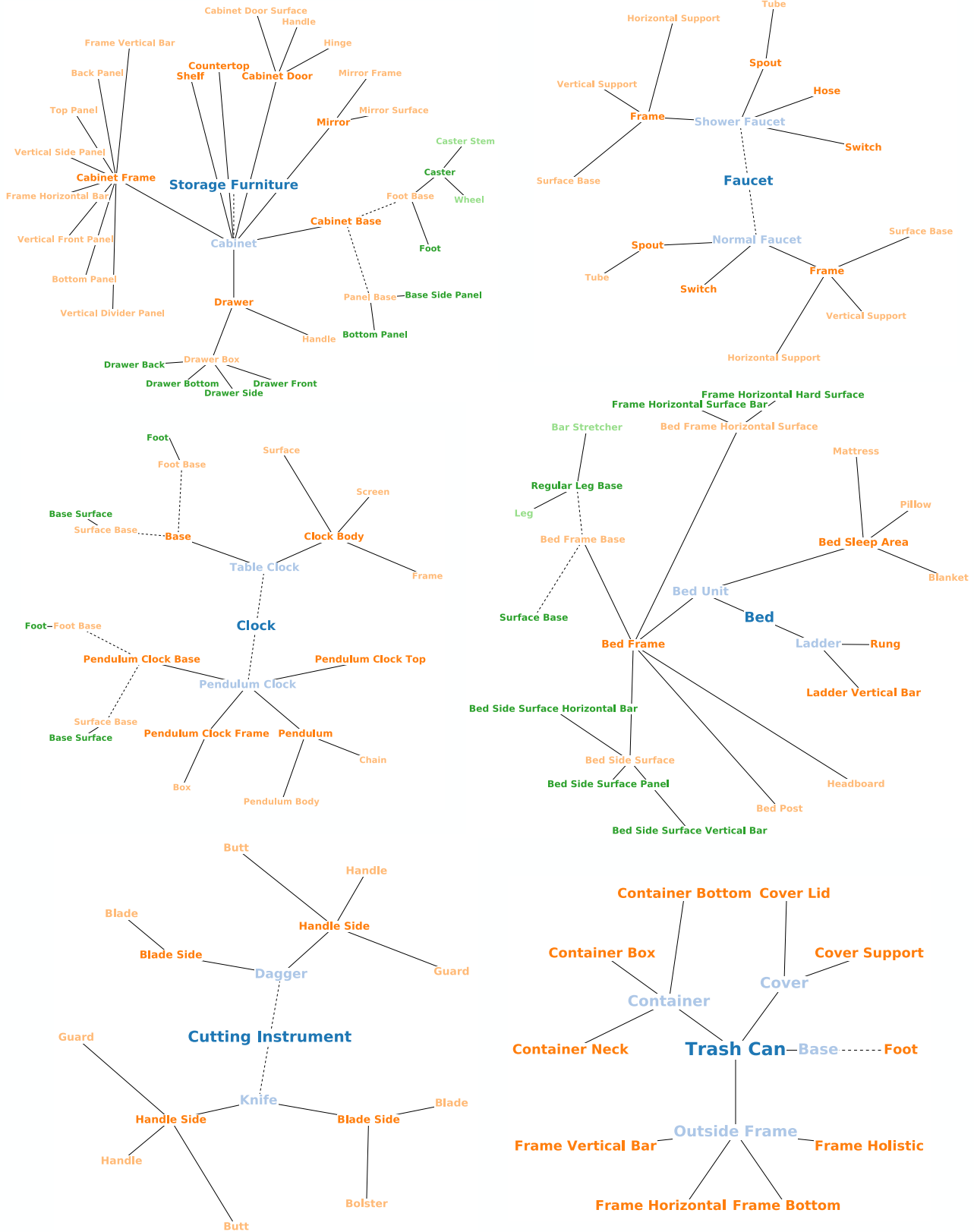


Figure 19. **Template visualization (2/3).** We present the templates for storage furniture, faucet, clock, bed, knife (cutting instrument), and trash can. The lamp template is shown in the main paper. The And-nodes are drawn in solid lines and Or-nodes in dashed lines.

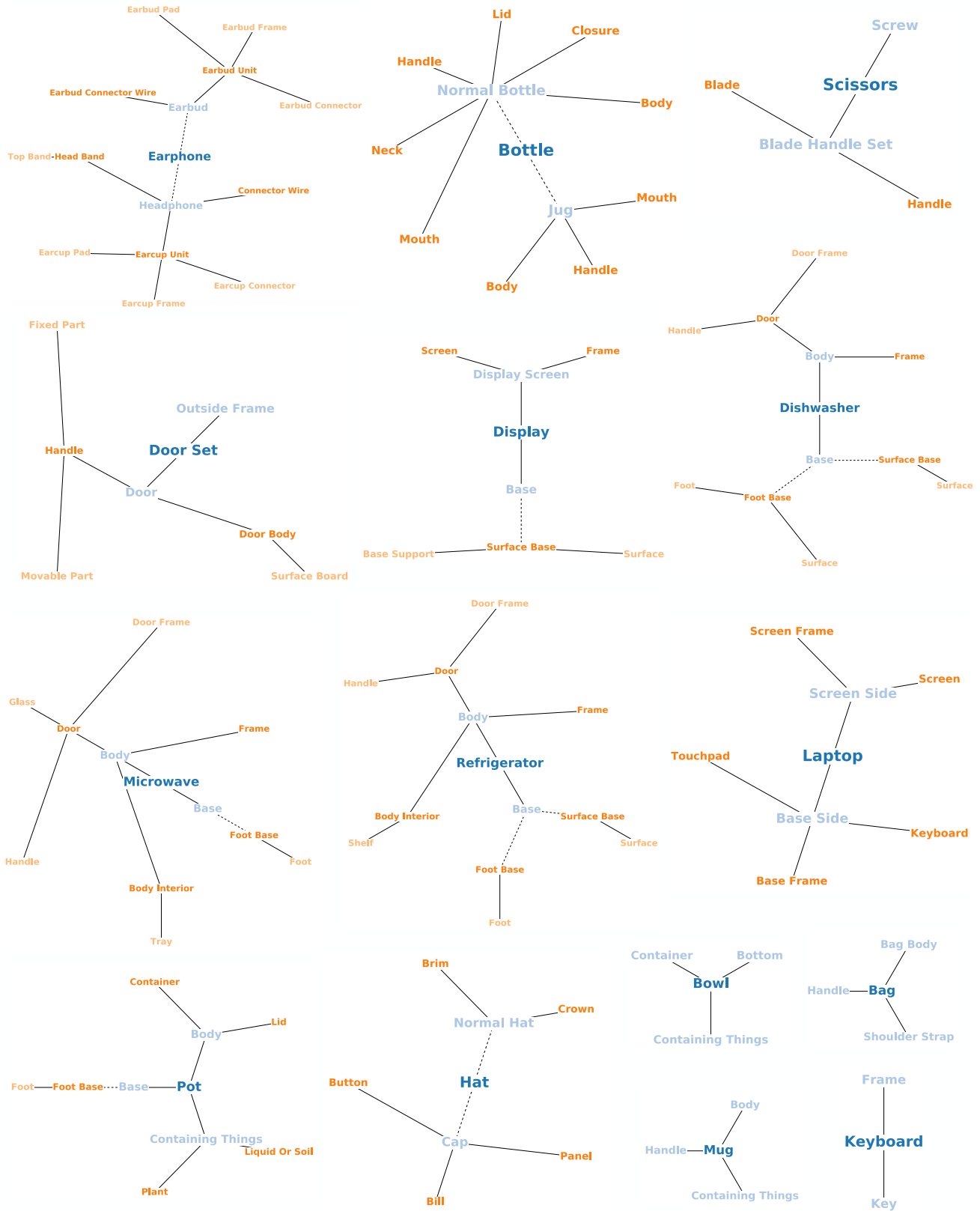


Figure 20. **Template visualization (3/3)**. We present the templates for earphone, bottle, scissors, door (door set), display, dishwasher, microwave, refrigerator, laptop, vase (pot), hat, bowl, bag, mug, and keyboard. The lamp template is shown in the main paper. The And-nodes are drawn in solid lines and Or-nodes in dash lines.