

Data-Driven Contextual Modeling for 3D Scene Understanding

Yifei Shi^a, Pinxin Long^c, Kai Xu^{a,c,*}, Hui Huang^{b,c}, Yueshan Xiong^a

^aHPCL, National University of Defense Technology

^bShenzhen University

^cShenzhen VisuCA Key Lab / SIAT

Abstract

The recent development of fast depth map fusion technique enables the realtime, detailed scene reconstruction using commodity depth camera, making the indoor scene understanding more possible than ever. To address the specific challenges in object analysis at subscene level, this work proposes a data-driven approach to modeling contextual information covering both intra-object part relations and inter-object object layouts. Our method combines the detection of individual objects and object groups within the same framework, enabling contextual analysis without knowing the objects in the scene *a priori*. The key idea is that while contextual information could benefit the detection of either individual objects or object groups, both can contribute to object extraction when objects are unknown.

Our method starts with a robust segmentation and partitions a subscene into segments, each of which represents either an independent object or a part of some object. A set of classifiers are trained for both individual objects and object groups, using a database of 3D scene models. We employ the multiple kernel learning (MKL) to learn per-category optimized classifiers for objects and object groups. Finally, we perform a graph matching to extract objects using the classifiers, thus grouping the segments into either an object or an object group. The output is an object-level labeled segmentation of the input subscene. Experiments demonstrate that the unified contextual analysis framework achieves robust object detection and recognition over cluttered subscenes.

Keywords: Scene understanding, object recognition, contextual modeling, data-driven approach

*Corresponding author: kaixu@nudt.edu.cn

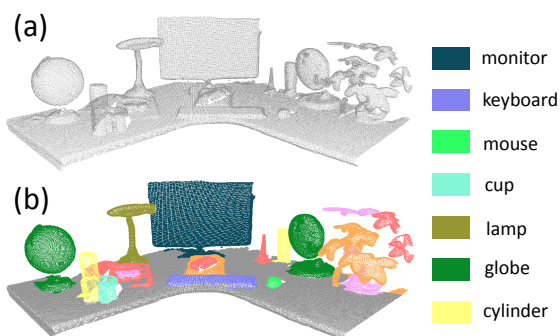


Figure 1: Scene understanding by our method. (a): The input point cloud of a table-top scene. (b): The labeling result (legends show semantic labels in color).

1. Introduction

2 With the rapid development of 3D sensing techniques,
3 the digitalization of large-scale indoor scenes has be-
4 come unprecedentedly accessible to a wide range of
5 applications. Among the most exciting and promising
6 applications, robot-operated exploration and interaction
7 over unknown indoor environment would benefit signif-
8 icantly from the availability of high-quality and realtime
9 acquired 3D geometry information [1]. Such 3D in-
10 formation can not only improve robot navigation and
11 exploration, but more importantly, facilitate efficient
12 robot-scene interaction with fine-grained understanding
13 of scene objects. The latter may support highly complex
14 robot tasks such as room cleaning.

15 Motivated by the high demand, extensive research has
16 been devoted to the understanding of scanned indoor
17 scenes. Most existing works on scene understanding
18 focus on large-scale objects, such as furniture, as well
19 as their spatial layout [2, 3, 4, 5, 6], since the analysis

20 is usually limited by the quality and resolution of input
21 scans. Recent advances in volumetric scan fusion tech-
22 nique (such as KinectFusion [7]) has made it possible
23 to reconstruct quality and detailed scenes from scans
24 captured by commodity depth camera (e.g. Microsoft
25 Kinect and Asus Xtion). The dense point clouds pro-
26 cessed by KinectFusion can well capture small scale ob-
27 jects such as household objects, which enables detailed
28 understanding at a subscene level, e.g. many objects
29 placed a tabletop; see Figure 1.

30 Object analysis at subscene level is arguably much more
31 challenging than that at whole scene level. *Firstly*, un-
32 like furnitures which are usually sparsely distributed in
33 an indoor scene, household objects are often highly clut-
34 tered due to the limited space of supporting surfaces [8].
35 For example, a tabletop scene is typically cluttered with
36 many on-table objects. *Secondly*, repetition of objects,
37 which is ubiquitous among furnitures and has been ex-
38 tensively exploited in previous works [3, 5], may not be
39 as commonly seen among household objects. For ex-
40 ample, the objects placed on a table are mostly unique
41 within the subscene. *Thirdly*, from the acquisition point
42 of view, smaller objects are often more sensitive to scan-
43 ning imperfection. These challenges make the existing
44 methods, dealing with large-scale furniture layout, un-
45 suitable for the object analysis of small-scale subscenes.

46 To address these challenges, it seems a natural option
47 is to fully utilize the inter-object relations, or contextual
48 information. However, a key prerequisite for contextual
49 scene analysis is that all objects are segmented and la-
50 beled with semantic tags [9], which is apparently infea-
51 sible for an unsegmented scene. Essentially, context is
52 defined with objects. Without knowing objects, how can
53 we utilize contextual information to help the identifica-
54 tion of objects? In this work, we try to tackle this prob-
55 lem through integrating the discovery of both individ-
56 ual objects and object groups into a unified framework.
57 While the former involves grouping parts into an object,
58 which detects individual objects, the latter amounts to
59 finding structure groups [10] composed of multiple ob-
60 jects, which can actually enhance or reinforce the de-
61 tection and recognition of objects within the structure
62 group. The key idea is that contextual information could
63 benefit the detection of either individual objects or ob-
64 ject groups, when objects are unknown. However, both
65 can contribute to object extraction.

66 To enable such unified framework, we take a data-driven
67 approach equipped with several key procedures. First,
68 we propose a robust segmentation method to partition a
69 indoor scene into segments which each represents either

70 an independent object or a part of some object. We then
71 train a set of classifiers for both individual objects and
72 object groups, based on a database of 3D scene models.
73 To improve the classification accuracy, we employ mul-
74 tiple kernel learning (MKL) [11] to learn per-category
75 optimized SVM classifiers for various objects and ob-
76 ject groups. Finally, we perform a graph matching to
77 extract objects using the classifiers, thus grouping the
78 segments into either an object or an object group. The
79 input of our algorithm is an indoor scene point cloud,
80 and the output is an object-level labeled segmentation
81 of the input scene. Experiments demonstrate the ro-
82 bust performance for both segment extraction and object
83 recognition on several subscenes.

84 Our approach possesses two key features compared with
85 previous methods. First, we perform a segmentation
86 process before recognition, which leads to robust han-
87 dling of cluttered scenes. Second, instead of solving the
88 recognition of individual objects and object groups as
89 two separate problems, we encode features of both indi-
90 vidual objects and object layout into a unified classifier
91 via contextual modeling.

92 2. Related Work

93 Scene understanding is a long-standing research topic
94 which has received extensive research from both com-
95 puter vision and computer graphics community. We
96 mainly review those works which take 3D point clouds
97 as input.

98 *Point cloud segmentation.* Mesh segmentation is a fun-
99 damental shape analysis problem in computer graphics,
100 for which both heuristic methods [12] and data-driven
101 approach [13] have been extensively studied over the
102 years. On the other hand, the segmentation of 3D point
103 clouds remains to be a challenging problem.

104 There are three kinds of methods for point cloud seg-
105 mentation [14]. The first type is based on primitive
106 fitting [3, 15, 5]. It is hard for these methods to deal
107 with objects with complex shape. The second kind
108 of techniques is the region growing method. Nan et
109 al. [2] propose a controlled region growing process
110 which searches for meaningful objects in the scene by
111 accumulating surface patches with high classification
112 likelihood. Berner et al. [16] detect symmetric regions
113 using region growing. Another line of methods formu-
114 lates the point cloud segmentation as a Markov Ran-
115 dom Field (MRF) or Conditional Random Field (CRF)

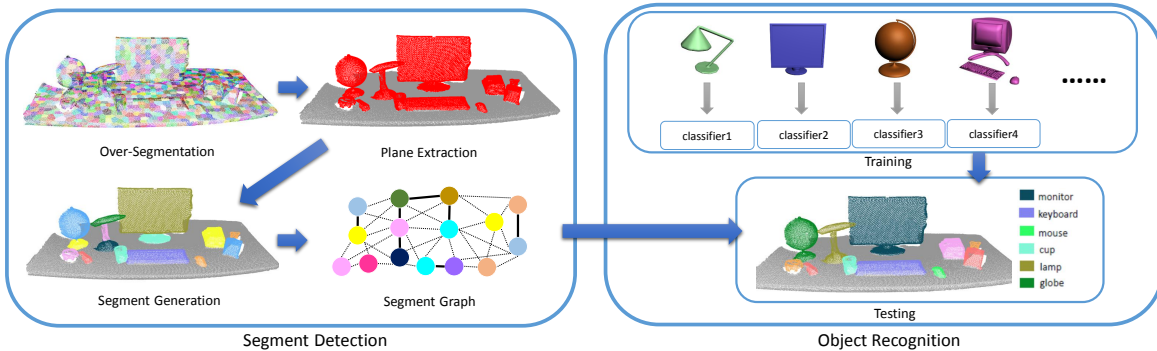


Figure 2: An overview of our algorithm. We first over-segment the scene and extract the supporting plane on the patch graph, then segment the scene into segments and represent the whole scene using a segment graph (a). To obtain the contextual information, we train a set of classifiers for both single objects and object groups using multiple kernel learning (b). The classifiers are used to group the segments into objects or object groups (c).

116 problem [4, 17, 14]. A representative random field seg-
 117 mentation method is the min-cut algorithm [17]. The
 118 method extracts foreground from background through
 119 building a KNN graph over which min-cut is performed.
 120 The shortcoming of min-cut algorithm is that the se-
 121 lection of seed points relies on human interaction. We
 122 extend the min-cut algorithm by first generating a set
 123 of object hypotheses via multiple binary min-cuts and
 124 then selecting the most probable ones based on a voting
 125 scheme, thus avoiding the seed selection.

126 *Object recognition.* Recently, the development of com-
 127 munity RGB-D cameras has opened many new oppor-
 128 tunities for 3D object recognition and scene recogni-
 129 tion [18, 19]. With the ever-growing amount of 3D mod-
 130 els becoming available, data-driven approach starts to
 131 play an important role in 3D object recognition and has
 132 gained great success [20].

133 Nan et al. [2] propose a search-classify approach to
 134 scene understanding by interleaving segmentation and
 135 classification in an iterative process. Li et al. [6] propose
 136 scene reconstruction by retrieving objects from a 3D
 137 model database. Song et al. [21] render database mod-
 138 els from hundreds of viewpoints and train an exemplar-
 139 SVM classifier for each of them to achieve object recog-
 140 nition. Their method overcomes several difficulties in
 141 object recognition, such as the variations of texture, il-
 142 lumination, etc. Chen et al. [22] utilize contextual infor-
 143 mation for indoor scene understanding. Small objects
 144 and incomplete scans can be recognized with the help of
 145 contextual relationships learned from database objects.
 146 Our method lends itself to cluttered indoor scene anal-
 147 ysis through integrating segmentation and recognition

148 into a single framework, which leads to a better per-
 149 formance when dealing with close-by objects than the
 150 contour-based method of [22].

151 Another line of analysis method is unsupervised learn-
 152 ing based on the presence of repetitions or symmetries
 153 in indoor scenes [3, 5, 23]. A limitation of such ap-
 154 proaches is that such repetitive patterns are less com-
 155 mon in subscenes dominated by household objects, e.g.,
 156 a tabletop scene.

157 *Plane extraction.* Plane extraction from point cloud is
 158 another important topic in scene understanding. For ex-
 159 ample, planes can be used to improve the reconstruc-
 160 tion of arbitrary objects containing both planar and non-
 161 planar regions [24].

162 Perhaps the most widely used approach for plane extrac-
 163 tion is RANSAC based plane fitting [15]. This method
 164 scales well with respect to the size of the input point
 165 cloud and the number of planes. Mattausch et al. [5]
 166 utilize planar patches as a compact representation of the
 167 point cloud of an indoor scene, which facilitates effi-
 168 cient repetition detection in a large-scale scene point
 169 cloud. Zhang et al. [24] perform plane extraction to
 170 delineate non-planar objects. Plane extraction has also
 171 been performed in the analysis of RGB-D data [25, 26].
 172 These works trim the plane boundary and convert the
 173 input data into a compact polygonal representation. Re-
 174 cently, Monzpart et al. [27] propose to reconstruct the
 175 raw scan of man-made scenes into an arrangement of
 176 planes with both local fitting and global regularization.

177 3. Overview

178 The input of our algorithm is a 3D point cloud of indoor
179 scene acquired and fused by KinectFusion. Our goal is
180 to detect objects in the scene and recognize their seman-
181 tic categories automatically. Our method proceeds in
182 two stages. First, we segment the point cloud into seg-
183 ments representing potential objects. Second, to achieve
184 object extraction and recognition, we propose a joint es-
185 timation of individual objects and object groups, as well
186 as their semantic categories.

187 *Segment detection.* In the first stage, we segment the
188 input scene (Figure 2 (a)). Specifically, we first over-
189 segment the entire scene and build a patch graph. We
190 then extract the supporting plane with a method inte-
191 grating RANSAC primitive fitting into graph-cut. Af-
192 ter plane extraction, the remaining points are grouped
193 into isolated groups. Within each group, we generate
194 segments via a robust segmentation algorithm, which
195 takes both geometry and appearance information into
196 account. Based on the segmentation, we represent the
197 entire scene as a segment graph with two types of edges
198 representing direct spatial adjacency (solid lines in Fig-
199 ure 2) and spatial proximity (dashed lines) between two
200 segments, respectively.

201 *Object extraction and recognition.* In the second phase,
202 we extract objects via recognizing both individual ob-
203 jects and object groups within a unified framework,
204 based on the above segment graph representation.

205 In an off-line stage, we train per-category optimized
206 SVM classifiers with multiple kernel learning for both
207 objects and object groups. The classifiers are trained
208 using 3D database models. Each 3D model is first con-
209 verted into 3D point cloud using virtual scanning and
210 segmented using the method mentioned above. We then
211 extract features from the corresponding segment graph
212 and train classifiers based on the graph.

213 In the online stage, we extract objects or object groups
214 from the segment graph of the input scene, through
215 searching for the subgraph matching corresponding to
216 the occurrence of database objects and object groups.
217 Once a matched subgraph is found, we use the cor-
218 responding SVM classifier to estimate the probability
219 of the match. Finally, we solve a labeling optimiza-
220 tion which minimizes the overall matching cost for all
221 matching probabilities.

222 4. Segment detection

223 Our goal is to partition the input scene into segments
224 which each represents either an independent object or
225 a part of an object. In order to segment objects from
226 cluttered scenes, we propose an unsupervised segment
227 detection approach to detect segments in 3D scene.

228 Specifically, we first over-segment the input point cloud
229 into a set of patches (Sec. 4.1) and detect the supporting
230 plane (Sec. 4.2). We then group the remaining patches
231 to extract potential objects or parts (Sec. 4.3) and rep-
232 resent them as a segment graph (Sec. 4.4). See Algo-
233 rithm 1 for an overview of our method.

234 4.1. Patch graph generation

235 We first over-segment the entire scene S into sev-
236 eral patches, using the method in [28]. We build a
237 patch graph based on the patches, denoted with $G_p =$
238 $(\mathcal{V}_p, \mathcal{E}_p)$, where \mathcal{V}_p and \mathcal{E}_p represent the patches and
239 the near-by relations within the patches, respectively.
240 Specifically, the near-by relations are determined by
241 comparing the nearest distance between two patches
242 with a threshold.

Essentially, our segment detection algorithm is a graph-
cut based approach. The most vital component for
graph-cut method is the definition of smooth term. In
this section, the smooth terms for all graph-cut opti-
mization are identical, which we first define here:

$$E_s(x_u, x_v) = w_c \cdot E_c + w_p \cdot E_p + w_n \cdot E_n, \quad (1)$$

243 where x_u, x_v are two adjacent patches. E_c, E_p, E_n are
244 the differences between two adjacent patches in terms of
245 color, planarity and normal. w_c, w_p, w_n are the weights.

E_c and E_p are computed based on the chi-square dis-
tance of the color and planarity histogram between u
and v , we normalize them to $(0, 1)$. It is worth mention-
ing that the planarity histogram are computed as fol-
low: first compute the least-square plane for a patch,
then built a histogram for distances of all points in the
patch to the plane. The formulation for E_n is different
for convex and concave situations. Specifically, the for-
mulation is:

$$E_n(x_u, x_v) = 1 - \eta(1 - \cos \theta_{u,v}), \quad (2)$$

246 where $\theta_{u,v}$ is the angle between the average normals of
247 patch P_u and P_v . For η , we take 0.01 (a small value) if
248 the two adjacent patches form a convex dihedral angle

Algorithm 1 :Segment Detection.

Input: scene S **Output:** segment graph G_s

- 1: $G_p \leftarrow \text{OverSegment}(S)$;
 - 2: $S \leftarrow \text{PlaneExtract}(S, G_p)$; //extract plane
 - 3: $\mathcal{H} \leftarrow \text{SegHypGen}(S, G_p)$; //generate seg. hypo.
 - 4: $T \leftarrow \text{SegHypSelect}(\mathcal{H})$; //select seg. hypo.
 - 5: $G_s \leftarrow \text{SegGraGen}(T)$; //generate seg. graph
 - 6: **return** G_s ;
-

249 and 1 otherwise, to encourage cuts around a concave
250 region [29].

251 Our smooth term takes both geometry (planarity and
252 normal) and appearance (color) factors into considera-
253 tion, thus makes the patches belong to different objects
254 can be detected easily.

255 4.2. Supporting plane extraction

256 Supporting plane is usually the largest object in most
257 subscenes of an indoor scene, such as tables, beds,
258 shelves, etc. The extraction of supporting plane is espe-
259 cially useful since it makes the detection of objects
260 on top of the supporting plane easier. Therefore, the
261 first step of our segment generation is supporting plane
262 extraction. For this task, perhaps the most straightfor-
263 ward approach is RANSAC based primitive fitting [15].
264 Since the objects placed on the supporting plane may be
265 very small or thin, setting a hard threshold for point-to-
266 plane distance may cause a lot of false positives. We
267 therefore improve this method by adding a graph-cut
268 optimization, to robustly segment on-top objects from
269 the supporting plane.

We try to assign each patch a binary label, denoted by
 $X = [x_1, \dots, x_n]$ with $x_i \in \{0, 1\}$. $x_i = 1$ if patch P_i lies
in the plane, and $x_i = 0$ otherwise. We formulate the
labeling problem as graph cuts over the patch graph:

$$E(X) = \sum_{u \in \mathcal{V}_p} E_d(x_u) + \sum_{(u,v) \in \mathcal{E}_p} E_s(x_u, x_v), \quad (3)$$

where the data term is defined as:

$$E_d(x_u) = \begin{cases} \delta, & \text{if } x_u = 1 \\ (1 - \frac{p}{p_{max}}) \cdot (1 - \frac{d}{d_{max}}) \cdot \cos \theta_{u,l}, & \text{if } x_u = 0 \end{cases}$$

270 where δ is a constant value, d the distance between the
271 center of u to the plane, and p the planarity of the patch.
272 d_{max} and p_{max} is the maximum distance and planarity,
273 respectively. We compute p as the average distance

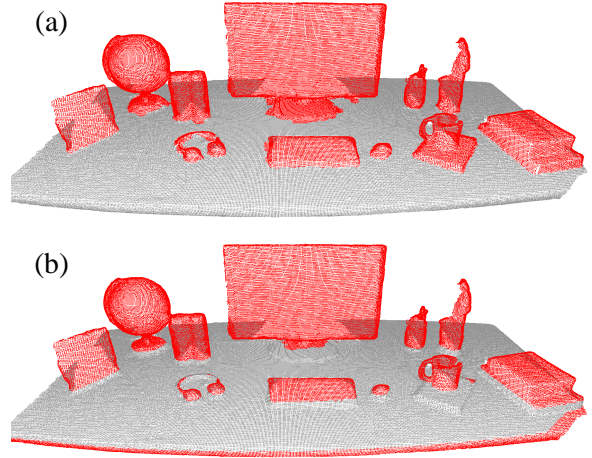


Figure 3: Plane extraction from the point cloud of a
tabletop scene by using our method (a) and RANSAC
based primitive fitting (b), respectively. While our
method can segment out the supporting plane accu-
rately, RANSAC missed some points due to the thin ob-
jects.

274 of all the points in patch P_u to its corresponding least-
275 square fitting plane. $\theta_{u,l}$ is the angle between the average
276 normal of P_u and the normal of the plane.

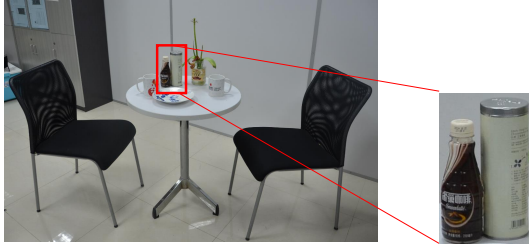
277 Figure 3 (a) demonstrates the segmentation results of
278 our method. As a comparison, the RANSAC based
279 primitive fitting can also get the majority of points cor-
280 rectly, but it fails when dealing with small and thin ob-
281 jects, as is shown in Figure 3 (b).

282 4.3. Segment generation

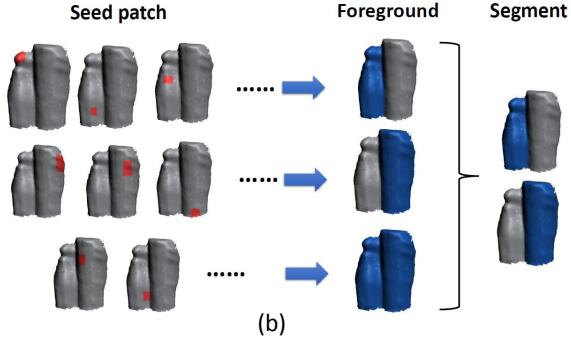
283 *Segment hypothesis generation.* After plane removal,
284 object extraction only amounts to segmenting the iso-
285 lated groups of patches on top of the supporting plane
286 into individual objects. To solve the problem, we pro-
287 pose to first generate a set of segment hypotheses and
288 then select the most prominent ones based on a voting
289 algorithm.

290 We first update the patch graph G_p by removing the
291 nodes belonging to the extracted plane. Based on the
292 updated patch graph, we generate segment hypothe-
293 ses by performing several times of binary graph cut,
294 where the foreground corresponds to potential objects
295 or prominent parts.

Different from other graph cut method, we do not se-
lect foreground seed heuristically. Instead, we use every
patch as seed and perform binary graph cuts for multiple



(a)



(b)

Figure 4: Illustration of our segment detection method. The scene is composed of two bottles stuck together on a round table (a). We use every patch as seed to generate many foreground hypotheses and then select the most prominent ones (b).

times, generating many candidate foregrounds. In each binary cut, we select one patch as foreground seed but do not prescribe any seed for background. This is performed by introducing a background penalty for each non-seed patch [30]. Specifically, we select one patch, denoted by P_s , labeling it as foreground $x_s = 1$, and minimize over binary patch labels $X = [x_1, \dots, x_n]$, $x_i \in \{0, 1\}$ (n is the number of patches) the following parametric energy function:

$$E^\lambda(X) = \sum_{u \in \mathcal{V}_p} E_d^\lambda(x_u) + \sum_{(u,v) \in \mathcal{E}_p} E_s(x_u, x_v), \quad (4)$$

where the data term is defined as:

$$E_d^\lambda(x_u) = \begin{cases} \infty, & \text{if } x_u = 0 \text{ and } u = s \\ 0, & \text{if } x_u = 1 \text{ and } u = s \\ 0, & \text{if } x_u = 0 \text{ and } u \neq s \\ f_u, & \text{if } x_u = 1 \text{ and } u \neq s \end{cases}$$

$$f_u = \begin{cases} k(d(P_s, P_u) - \lambda), & \text{if } d(P_s, P_u) > \lambda \\ 0, & \text{otherwise.} \end{cases}$$

²⁹⁶ f_u is the background penalty which penalizes a non-
²⁹⁷ seed patch which is distant from the foreground seed.

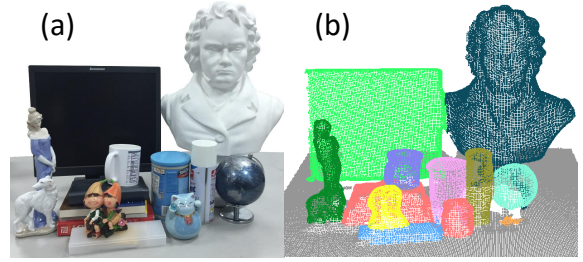


Figure 5: Segment detection from the point cloud of a highly cluttered scene (a) by using our method (b). The input data has a lot of close-by objects and the back view is not scanned, which makes the segmentation quite challenging. Our method can segment out most objects accurately.

²⁹⁸ $d(P_s, P_u)$ is the distance between the centers of patch P_s
²⁹⁹ and P_u . We use $k = 2.0$ for a steep penalty to quickly
³⁰⁰ reject those patches whose distance to P_s is larger than
³⁰¹ λ to be labeled as foreground. The parameter λ controls
³⁰² the range, centered around the foreground seed, within
³⁰³ which one seeks for foreground patches. Instead of us-
³⁰⁴ ing a hard threshold on this range, we slide λ from 0 to
³⁰⁵ ℓ_d (the diagonal length of the bounding box of the entire
³⁰⁶ scene) and find the first point where the total cut cost
³⁰⁷ drops significantly (up to 50%) and take the resulting
³⁰⁸ cuts as the segmentation result. The smooth term is the
³⁰⁹ same as the one used for plane extraction in Eq. (1).

³¹⁰ Once we select every patch as seed and perform graph
³¹¹ cut for each of them, we can obtain a set of foreground
³¹² segments. To filter out the redundancy, we cluster
³¹³ the foreground segments using non-parametric mean-
³¹⁴ shift [31]. The similarity between two segments, de-
³¹⁵ noted as S and T , is measured by the Jaccard index, i.e.,
³¹⁶ $s(S, T) = |S \cap T| / |S \cup T|$. For example, as is shown
³¹⁷ in Figure 4 (b), selecting the seed patches in the same
³¹⁸ row will led to identical foregrounds, thus these fore-
³¹⁹ grounds will cluster together after the mean-shift pro-
³²⁰ cessing, as the Jaccard index is high. For each cluster,
³²¹ we choose the cluster center as the segment hypothesis
³²² for that cluster. As a result, we obtain a pool of k hypo-
³²³ thetic segments, $\mathcal{H} = \{H_i\}_{i=1}^k$.

Segment hypothesis selection. The set of hypotheses may overlap with each other, making the labeling of patches ambiguous. To select good hypotheses without relying on heuristics or supervision, we propose a multi-class Markov random field (MRF) segmentation with object label selection, which minimizes the follow-

ing energy function:

$$E(L) = \sum_{u \in \mathcal{V}_p} E_d(l_u; P_u) + \sum_{(u,v) \in \mathcal{E}_p} E_s(l_u, l_v), \quad (5)$$

over the labeling for all patches: $L = [l_1, \dots, l_n]$, $l_i \in \{1, \dots, k\}$.

The data term $E_d(l_u; P_u)$ is defined as the likelihood that the patch P_u belongs to a particular segment hypothesis. For instance, for patch P_u and hypothesis H_i , we define the data term as the frequency of P_u being covered by the hypotheses in H_i :

$$E_d(H_i; P_u) = -\ln(t(P_u, C_i) / \sum_j t(P_u, C_j)), \quad (6)$$

where $t(P_u, C_i) = |\{P_u \subset H_j | H_j \in C_i\}|$ is the presence times of patch P_u in cluster C_i . The smooth term is also the same as the one in Eq. (1).

The data term selects a label for each patch based on a consensus voting by all foreground clusters: The larger a foreground cluster is, the more probable that its corresponding segment hypothesis represents an independent object, since the object is proposed by many binary segmentations. Figure 4 depicts our segment detection algorithm and Figure 5 demonstrates the segmentation results over a highly cluttered scene.

4.4. Segment graph generation

To deal with the recognition for both object and object group, we represent the entire scene as a segment graph $G_s = (\mathcal{V}_s, \mathcal{E}_s)$, where \mathcal{V}_s represents the segments we detected in the input scene and \mathcal{E}_s encodes the relationship between two segments. We use two kinds of edges to describe relations in G_s . If the shortest distance between two segments is less than a small threshold t_s , we use a *connection edge* to link them, that means the two segments contact with each other and probably belongs to the same object. If the shortest distance between two segments is large than the small threshold but less than a larger threshold t_l , we use a *proximity edge* to connect them, which means they are in the same supporting plane and has the potential to constitute a object group. The two kinds of edges represent the contextual information for intra-object part relations and inter-object object layouts, respectively. t_l is selected as slightly larger than the largest bounding box diagonal length of all object groups in the database. Figure 2 shows an illustration the segment graph of the given input scene.

5. Object Recognition

5.1. Training

When recognizing a scene containing multiple objects, human perception is predominantly affected by three levels of prior knowledge [32]: the shape information of individual parts, the part composition of individual objects, and the contextual relationship among object groups. In our object recognition procedure, we encode all these knowledge in an unified model and recognize objects and object groups simultaneously. Specifically, we train per-category optimized SVM classifiers for all kinds of objects and object groups, and then utilize these classifiers to test the category of the input segments. Here, an object group is refer to a group of objects whose co-occurrence is frequently seen in an indoor scene category [33]. For example, the monitor-keyboard-mouse combo is frequently seen in office.

Data Preparation. To learn the model from the database of 3D scene models, the first step is to convert the database models (training data) into point cloud representation, which is compatible against the input (test data), and extract features from the point clouds.

First, we download a set of 3D CAD models of household objects, denoted by $\{\Gamma_i\}$, from the internet. Each Γ_i contains the models belonging to the same shape category. Second, we collect indoor scene models from the dataset of [9] and [10]. In order to obtain object groups which are not only frequently occurring but also semantically significant, we extract local substructures $\{\Phi_i\}$ from the dataset as the focal points defined in [33]. Each Φ_i contains the substructures belonging to the same semantic group.

We then perform virtual scanning for all models/groups in $\{\Gamma_i\}$ and $\{\Phi_i\}$, similar to [2]. Such virtual scan could mimic the real situation of object clutter or incomplete scan, making the training data more suitable for learning a generalizable recognition model. After the virtual scanning, we compute segment graphs using the method described in Sec. 4 for object groups in $\{\Phi_i\}$. For individual objects in $\{\Gamma_i\}$, we perform the same process except for table extraction. The label of each virtually scanned point is determined by aligning the point cloud with the original 3D CAD models and transferring the labels based on closest point search.

Classifier Learning. We compute two kinds of features for our SVM classifier: node features and edge features.

Algorithm 2 :Training.

Input: object set $\{\Gamma_i\}$ and object group set $\{\Phi_i\}$ **Output:** classifiers C

```
1: for all  $\Gamma_i$  do
2:   for all  $\gamma_j$  in  $\Gamma_i$  do
3:      $\gamma_j \leftarrow \text{VirtualScan}(\gamma_j)$ ;
4:   end for
5:    $g_i \leftarrow \text{ConstructSegGraph}(\Gamma_i)$ ;
6:    $c_i^\gamma \leftarrow \text{MKL}(g_i)$ ;
   //train SVM for each single object category
7: end for
8: for all  $\Phi_i$  do
9:   for all  $\phi_j$  in  $\Phi_i$  do
10:     $\phi_j \leftarrow \text{VirtualScan}(\phi_j)$ ;
11:   end for
12:    $g_i \leftarrow \text{ConstructSegGraph}(\Phi_i)$ ;
13:    $c_i^\phi \leftarrow \text{MKL}(g_i)$ ;
   //train SVM for each object group category
14: end for
15: return  $C = \{c_i^\gamma\}_{i=1}^m + \{c_i^\phi\}_{i=1}^n$ ;
```

For each node, we voxelize its bounding box and extract features of shape, normal and volume as described in [21]. In addition, we estimate the oriented bounding box (OBB) for each object and measure its anisotropy:

$$c_l = \frac{s_1 - s_2}{(s_1 + s_2 + s_3)}, c_p = \frac{2(s_2 - s_3)}{(s_1 + s_2 + s_3)}, c_s = \frac{3s_3}{(s_1 + s_2 + s_3)} \quad (7)$$

where s_1, s_2, s_3 are the three scales of the OBB with $s_1 > s_2 > s_3 \geq 0$. For each edge, we compute the layout similarity [33] as its feature:

$$\gamma(p, q) = \frac{d_H(\text{obb}(p), \text{obb}(q))}{dl(p) + dl(q)}, \quad (8)$$

$$\rho(p, q) = \text{angle}(\mathbf{v}_{dir}(p, q), \mathbf{v}_{upright}), \quad (9)$$

402 The two features measure the distance and direction between two objects, respectively.

404 We compute features and learn pre-category optimized SVM classifiers for each category of individual objects 405 in $\{\Gamma_i\}$ and object groups in $\{\Phi_i\}$. Positive examples are 406 the models from the two datasets, while negative ones are 407 generated by using the method in [21] for individual 408 objects and the method in [22] for object groups. 409 In addition, we associate a triplet (n_n, n_c, n_p) with each 410 classifier, where n_n, n_c and n_p represent the number of 411 segments, edges and proximity edges, respectively. This 412 triplet is used to perform a coarse matching based on the 413 triplet, before testing with the classifier. 414

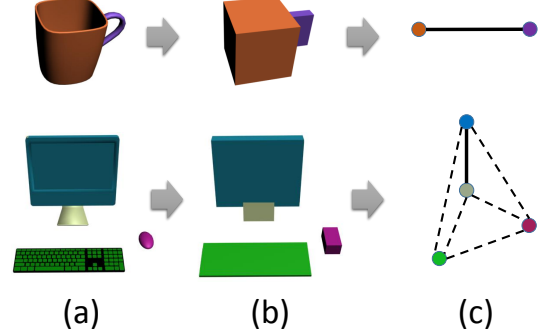


Figure 6: The generation of object and object group. The input is a segmented object or object group (a). We compute the OBB for each part (b) and connect them into a graph (c). The solid and the dashed lines in (c) are connection and proximity edge, respectively.

415 **Multiple Kernel Learning.** Kernel method has been 416 successfully applied into many learning areas, while the 417 results of these methods are heavily dependent on the 418 selection of kernels. Instead of choosing a single kernel, 419 it is better to have a set of kernels and use the combination 420 of them [11]. Since our features are computed 421 for both individual objects and their relations, it is especially 422 desirable to combine several kernels and to allow 423 the classifiers to choose their optimized kernels, in order 424 to reduce their bias [34]. The idea is to use a combination 425 of basic kernels $k(\mathbf{x}, \mathbf{y}) = \sum w_i \cdot k_i(\mathbf{x}, \mathbf{y})$ rather than 426 a specific kernel in SVM. The basic kernels could be 427 linear kernel, Gaussian kernel, polynomial kernel, etc.

428 Figure 7 illustrates the architecture of our MKL-based 429 classification. Given the segment graph of an individual 430 object or an object group, we first represent it in the feature 431 space spanned with six kinds of features. We then 432 transform the data from feature space to kernel space using 433 several predefined kernels. By computing the optimized 434 weights for each kernel space, we obtain the final 435 MKL classifier. The procedure for training the classifiers 436 is detailed in Algorithm 2.

437 5.2. Testing

438 **Data Preprocessing.** The segments in scenes acquired 439 by Kinect or any other commodity depth camera are 440 usually noisy and low-quality, making the recognition 441 quite difficult. Therefore, we first surface reconstruction 442 [35] to form a watertight surface for each segment, 443 and then compute features as described in Sec. 5.1.

Algorithm 3 :Testing.

Input: classifiers C and segments T
Output: segments label X

```

1: for all  $c_i$  in  $C$  do
2:   if Matching( $c_i, T$ ) then
3:      $cost_i \leftarrow$  ComputeProbability( $c_i, T$ );
4:   end if
5: end for;
6:  $X \leftarrow$  ComputeLabel( $\{cost_i\}_{i=1}^k$ );
   //compute label for all segments
7: return  $X$ ;

```

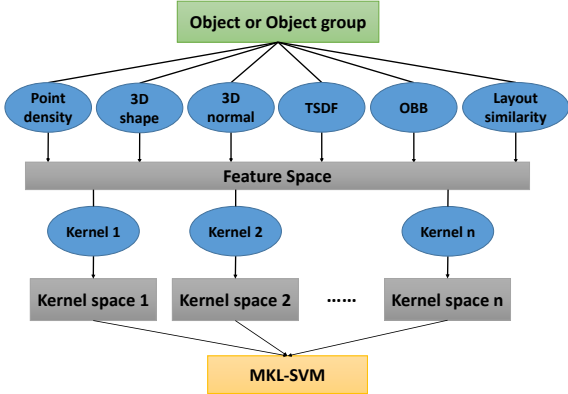


Figure 7: The architecture of our MKL-based classifier. Given an object or a object group, we compute its features and map it into several kernel spaces with several basic kernels. The MKL-SVM classifier is learned by computing the optimized weight for each kernel.

444 *Labeling Optimization.* To extract objects and object
445 groups from the segment graph, we search from the seg-
446 ment graph of the input scene for the subgraphs cor-
447 responding to the occurrences of database objects and
448 object groups. Graph matching can be formulated as
449 quadratic assignment problem, which is known to be
450 NP-hard, so an exhaustive search over the whole graph
451 leads to high computational cost.

452 In our method, the graph matching is performed as fol-
453 lows. For each MKL classifier, we first use the associ-
454 ated triplet (n_n, n_c, n_p) to filter subgraph matchings. A
455 subgraph is filtered if any one of the three terms is dif-
456 ferent from that of the classifier. For the remaining sub-
457 graphs, we use the learned MKL classifiers to test if it
458 belongs to the corresponding category and record the
459 probability if yes. The probability will be used as the
460 labeling cost which penalizes the mislabeling in the fol-
461 lowing optimization.

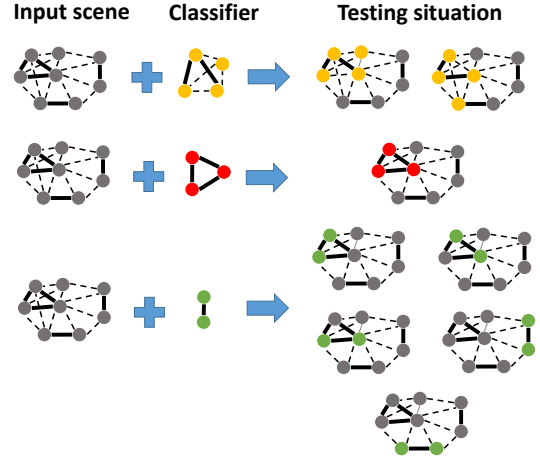


Figure 8: The matching strategy of our algorithm. Given a segment graph of the input scene on the left, we use all the three classifiers to test the occurrence of the corresponding subgraph. The testing samples are shown on the right. Note that some connection edge in the first row can be turned into a proximity one to allow more matches.

462 After applying all classifiers, we detect all the potential
463 objects or object groups in the input scene. The graph
464 matching strategy is illustrated in Figure 8. Note that we
465 allow a connection edge to be converted into a proximity
466 one to produce more matchings. The rationale of this is
467 that some segments not belonging to the same object
468 could be linked by connection edges mistakenly due to
469 small mutual distance.

Next, we solve a labeling optimization which minimizes the overall matching cost computed from all the matching probability. The final labeling, X , for all segments of the input scene is computed by:

$$X = \operatorname{argmin}_X \sum_{c_i \in C} D(X, c_i) \quad (10)$$

where:

$$D(X, c_i) = \begin{cases} 0, & \text{if recognized subgraph by} \\ & c_i \text{ is labeled correctly in } X \\ \text{cost}(X, c_i), & \text{otherwise.} \end{cases}$$

470 where $\text{cost}(X, c_i)$ is the labeling cost penalizing the
471 wrong labeling of the subgraph detected by the classifier
472 c_i . We found it suffices to solve this labeling optimiza-
473 tion using a combinatorial search over all labeling posi-
474 bilities since the possible labeling for each segment
475 is limited after the classifier filtering and testing. The

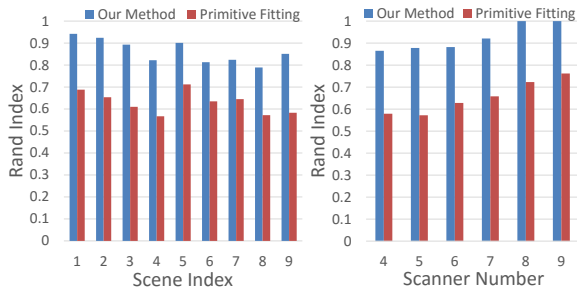


Figure 9: Segmentation comparison against the RANSAC based primitive fitting method [15]. Left: Comparison over nine test scenes. Right: Results of our method and the RANSAC-based one over scene #2 with increasing number of scans.

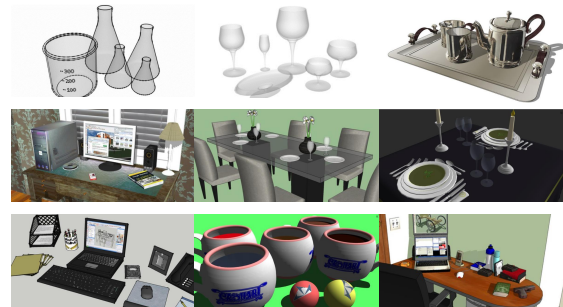


Figure 10: The test scenes used in segmentation evaluation.

476 whole testing process for object and object group detection
477 tion is described in Algorithm 3.

478 6. Results and Evaluation

479 We test our method on both real-world and virtually
480 scanned scenes. A gallery of results is shown in Fig-
481 ure 18. We first describe the experimental setting of our
482 method and then evaluate our method in two aspects,
483 i.e., the segment detection and the object recognition.

484 *Experimental Setting.* Our method is implemented us-
485 ing C++ and run on a desktop PC with an Intel I5-3750
486 CPU (quad core, 3.4GHz) and Nvidia GeForce GTX
487 460 graphics card. We scan a few indoor scenes using a
488 Microsoft Kinect. We also use the Washington scene
489 dataset [36] acquired by an ASUS Xtion PRO LIVE
490 RGB-D sensor. The parameter settings are provided be-
491 low. Patch size (diameter): 8cm for NYU-Depth V2
492 dataset and 4cm for others; w_c , w_p , and w_n in 1: 0.2,
493 0.3, and 0.5, respectively; δ for table extraction: 0.95
494 for all datasets; t_s and t_l for segment graph construc-
495 tion: 3cm and 50cm, respectively; Poisson iso-point
496 sampling density: 2cm; basic kernels for MKL (we use
497 SimpleMKL [37]): five Gaussian kernels and two poly-
498 nomial kernels.

499 *Segment Detection.* We test our segment detection al-
500 gorithm on nine tabletop scenes downloaded from the
501 Internet (Figure 10) and virtually scanned. We compare
502 our method with the RANSAC-based primitive fitting
503 method in [15]. The Rand Index [38] is used as the
504 evaluation criterion. We perform six tests on each scene

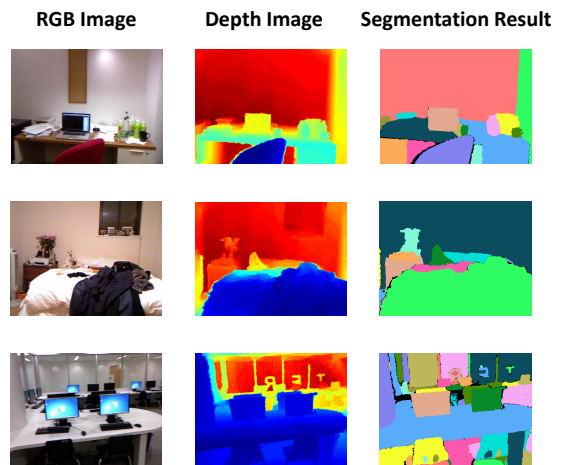


Figure 11: The segmentation results our algorithm over the scenes from the NYU-Depth V2 dataset. Our method can segment most objects correctly in the highly cluttered scenes.

505 with different number of scan and quality and take the
506 average Rand Index. In the virtual scanning, the virtual
507 scanners are positioned around the scene being scanned
508 and oriented to the center of the scene. The plot in Fig-
509 ure 9 (left) show that the Rand Index of our method is
510 higher than that of the RANSAC-based method over the
511 nine test scenes. We also evaluate how scan quality
512 would affect the segmentation results with the varying
513 number of scans for scene #2; see Figure 9 (right).

514 We also test our segmentation approach on NYU-Depth
515 V2 dataset. A significant feature of the depth images is
516 that the point cloud is of low resolution, making our seg-
517 mentation infeasible. In order to tackle this kind of in-
518 put, we made some changes over our algorithm. Given
519 an RGB-D image and its camera parameters, we first

	Primitive fitting[16]	Support relation[9]	Our method
Rand Index	61.8%	78.7%	76.4%

Figure 12: A comparison of the segmentation accuracy (Rand Index) of the methods in [15] and [8] and ours on the NYU-Depth V2 dataset.

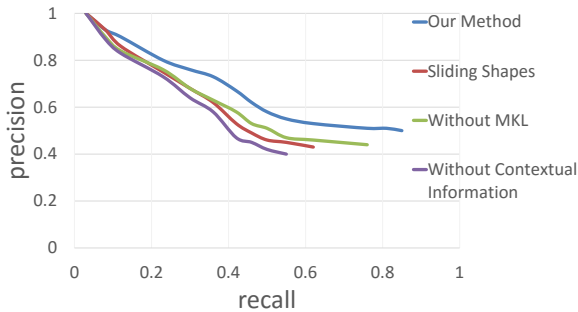


Figure 13: Precision-recall curves for object recognition. Comparison is made between our method and the other three methods by testing on the database in [36].

520 project the 2D points into 3D space to reconstruct a 3D
521 scene. We skip the table extraction process and detect
522 the segments for the near-camera points (distance less
523 than 2m) using our method, and cluster the rest distant
524 points using Euclidean cluster extraction [39].

525 We test our method on a selected subset of the NYU-
526 Depth V2 dataset as in [22], which contains 45 living
527 rooms and offices. Some results of our algorithm are
528 shown in Figure 11. We compare our method with the
529 support relation based method in [8] and the RANSAC-
530 based one in [15]. The segmentation Rand Index mea-
531 sures for the three methods are shown in Figure 12. The
532 support relation based method slightly outperforms our
533 method, due to the incorporation of the high-level prior.

534 *Object Recognition.* Our recognition database contains
535 900 objects in 18 categories and 10 kinds of object
536 groups. We test our object recognition method on two
537 scanned scene datasets. The first one is several real-
538 world scenes such as office, meeting room, and labo-
539 ratory, scanned by ourselves and the second dataset
540 from [36]. The scenes contain a variety of object cat-
541 egories with noisy and low quality scans.

542 Figure 18 demonstrates the results on six indoor scenes.
543 The semantic labels are shown using distinct colors,
544 while the contextual information is illustrated with red

545 dots and dashed lines. The majority of objects can be
546 recognized correctly, benefiting from the contextual in-
547 formation. The geometric ambiguity between different
548 categories of objects, such as a keyboard and a book, are
549 resolved with the help of contextual information. Some
550 segments are correctly segmented but not successfully
551 recognized due to capability of our recognition model
552 learned from the limited model database. This can be
553 improved by collecting more data and training a more
554 powerful model.

555 We evaluate our method on the database of [36] con-
556 taining 58 indoor scenes collected using KinectFusion.
557 We compare to three alternative methods: the sliding
558 shapes [21], a reduced version of our method by us-
559 ing linear SVM classifiers, and a reduced method with-
560 out using contextual information. The precision-recall
561 curves for recognition are plotted in Figure 13. It is
562 obvious that our method outperforms sliding shapes,
563 thanks to the object-group-level analysis and the MKL
564 classifiers in our method. The reduced method with-
565 out contextual information is slightly inferior to sliding
566 shapes. This is because sliding shapes use a plethora
567 of classifiers, which is three orders of magnitude more
568 than what our method uses.

569 As demonstrated in Figure 14, our method benefits from
570 the contextual information in two ways. First, context
571 helps to eliminate recognition ambiguity. For example,
572 the object in Figure 14 (a) can either be a book or a
573 keyboard, which is correctly recognized with the help
574 of the monitor-keyboard-mouse combo. Second, con-
575 text can enhance the recognition ability under low data
576 quality. For example, the cup in Figure 14 (b) is hard
577 to be recognized due to the low data quality, where the
578 cup-cup group helps recognize it.

579 We make two observations from the results. (1) The
580 precision is consistently high with the increasing of the
581 recall. (2) The recall converges to a high value but never
582 reaches 1 with the precision decreasing. These obser-
583 vations can be explained by the inter-restriction of the
584 multiple MKL classifiers. Our method finds a labeling
585 that tries to satisfy all the MKL detectors as much as
586 possible, leading to more reliable labeling result.

587 To evaluate the performance our method on cluttered
588 scenes, we scan six desktop scenes with an increasing
589 degree of object clutter. The objects we recognized are
590 highlighted with boxes in Figure 15. It is clear that our
591 method achieves robust recognition on these cluttered
592 scenes, especially the one in Figure 15 (c). As a com-
593 parison, the method in [22] cannot recognize the mouse
594 in (c), because the contour-based approach fails when

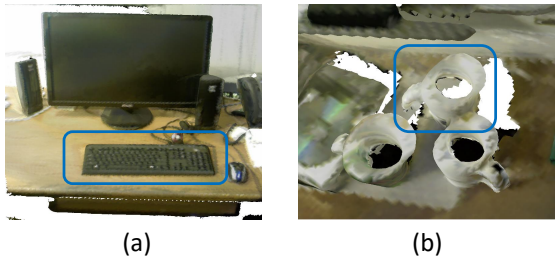


Figure 14: The contextual knowledge could benefit object recognition in two ways. (a): Resolving recognition ambiguity: The keyboard in blue box is recognized correctly due to the contextual information of the monitor-keyboard-mouse combo. (b): Enhancing recognition ability: The cup in blue box is in low scan quality but can be recognized based on the cup-cup combo.

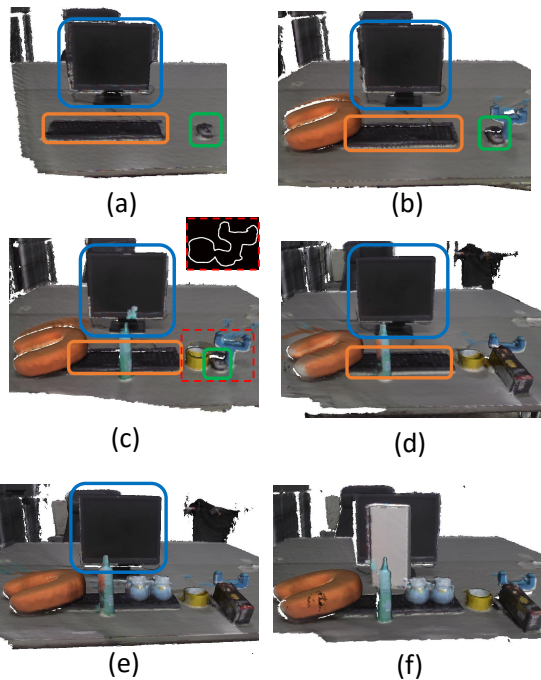


Figure 15: Our recognition results on several scenes with increasing degree of object from (a) to (f). The monitors, keyboards and mouses are correctly recognized by our method and labeled with blue, orange and green boxes.

595 dealing with cluttered scenes due to the incorrect con-
 596 tour extraction. The contour of the red box area in (c) is
 597 shown on the top-right corner.

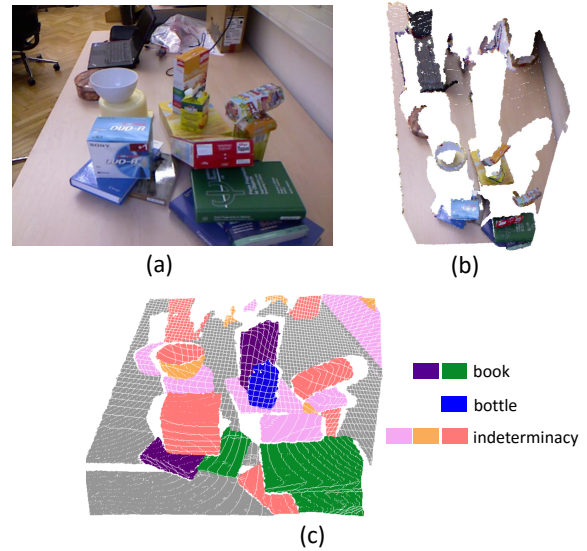


Figure 16: A failure case of our method. Our method cannot recognize most of the objects in a cluttered scene (c). This is due to the fact that the scene point cloud is only a single-view scan (b).

598 *Time Performance.* For a scene with 100K points, the
 599 segment detection takes 20 seconds. The training proced-
 600 ure of our object recognition is determined by the num-
 601 ber of individual object and object group categories. In
 602 our case, it takes about 1 hour to train a classifier us-
 603 ing SimpleMKL averagely. The training process takes
 604 about 32 hours in total for the 18 objects and the 10 ob-
 605 ject groups. The testing time is determined by the num-
 606 ber of segments and the degree of object clutter. The
 607 testing time for the scenes in Figure 18 (a) to (f) are 7.8,
 608 19.1, 39.5, 20.3, 1.7 and 12.9 minutes, respectively.

609 *Limitations.* Our method has the following limitations.
 610 First, our method does not provide a mechanism to deal
 611 with input data with severe missing parts. For example,
 612 if the input contains only a single-view scan, our method
 613 would not be able to produce meaningful segments for
 614 further analysis. A failure case of this is shown in Fig-
 615 ure 16. Second, our method can tolerate only moder-
 616 ate shape variation. It might fail when recognizing ob-
 617 jects with too special structure of segment graph, such
 618 as the case shown in Figure 17. Last, our method works
 619 the best for a scene containing a planar support. Al-
 620 though quite commonly seen in everyday indoor envi-
 621 ronments, the assumption does not generalize well for
 622 outdoor scenes.

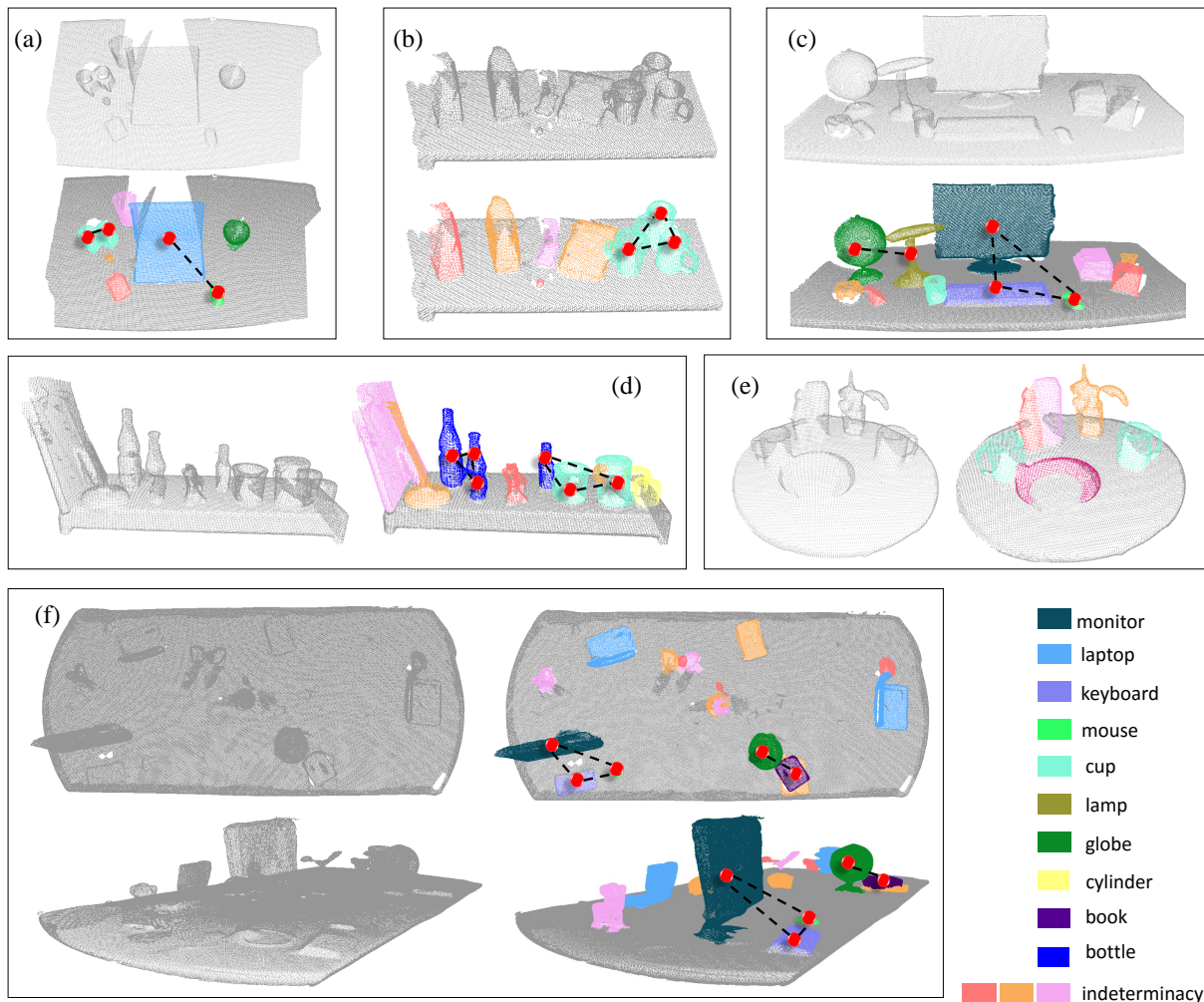


Figure 18: A gallery of scene understanding results by our method.

623 7. Discussion and future work

624 To achieve object analysis from clustered subscenes, we
 625 have developed a unified framework for the discovery of
 626 both individual objects and object groups, both of which
 627 are based on the contextual information learned from a
 628 database of 3D scene models. Our method makes the
 629 contextual information applicable even without know-
 630 ing the object segmentation of the input scene. The lat-
 631 ter has so far been predominantly assumed by existing
 632 methods, e.g., [22].

633 We see three venues for future work. First, our current
 634 work focuses on subscene analysis. It would be inter-
 635 esting to extend our method to deal with whole scene,

636 leading to multi-scale scene analysis in a unified frame-
 637 work. Currently, the contextual information is based on
 638 spatial proximity. As another future work, we would
 639 like to expand our contextual features with multi-modal
 640 object interaction, such as dynamic motion, to address
 641 more complex mutual relations among objects. Finally,
 642 it is natural to utilize our framework in robot-operated
 643 autonomous scene scanning and understanding.

644 Acknowledgements

645 We thank all the reviewers for their comments and feed-
 646 back. We would also like to acknowledge our research
 647 grants: NSFC (61572507, 61202333, 61379103),

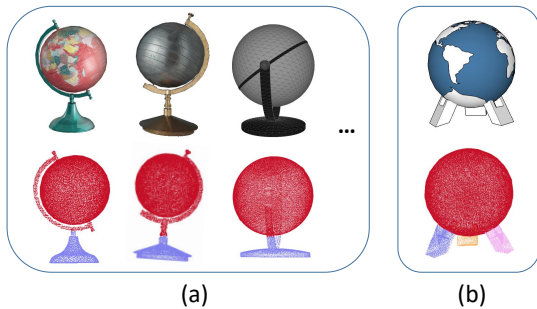


Figure 17: *The object classifier for a globe is trained using the examples containing two components (a). The recognition may fail when testing an exceptional instance of globe with three legs (b).*

648 973 Program (2014CB360503), Guangdong Science and Technology Program (2015A030312015, 649 2014B050502009, 2014TX01X033), Shenzhen VisuCA Key Lab (CXB201104220029A).

652 References

- 653 [1] K. Xu, H. Huang, Y. Shi, H. Li, P. Long, J. Caichen, W. Sun, B. Chen, Autoscanning for coupled scene reconstruction and proactive object analysis, *ACM Trans. on Graphics (Proc. of SIGGRAPH Asia)* 34 (6) (2015) 177:1–177:14.
- 654 [2] L. Nan, K. Xie, A. Sharf, A search-classify approach for cluttered indoor scene understanding, *ACM Trans. on Graphics (Proc. of SIGGRAPH Asia)* 31 (6) (2012) 137:1–137:10.
- 655 [3] Y. M. Kim, N. J. Mitra, D.-M. Yan, L. Guibas, Acquiring 3d indoor environments with variability and repetition, *ACM Trans. on Graphics (Proc. of SIGGRAPH Asia)* 31 (6) (2012) 138:1–138:11.
- 656 [4] T. Shao, W. Xu, K. Zhou, J. Wang, D. Li, B. Guo, An interactive approach to semantic modeling of indoor scenes with an RGBD camera, *ACM Trans. on Graphics (Proc. of SIGGRAPH Asia)* 31 (6) (2012) 136:1–136:11.
- 657 [5] O. Matusch, D. Panozzo, C. Mura, O. Sorkine-Hornung, R. Pajarola, Object detection and classification from large-scale cluttered indoor scans, *Computer Graphics Forum (Proc. of Eurographics)* 33 (2) 11–21.
- 658 [6] Y. Li, A. Dai, L. Guibas, M. Nießner, Database-assisted object retrieval for real-time 3d reconstruction, *Computer Graphics Forum (Proc. of Eurographics)* 34 (2) (2015) 435–446.
- 659 [7] R. A. Newcombe, A. J. Davison, S. Izadi, P. Kohli, O. Hilliges, J. Shotton, D. Molyneaux, S. Hodges, D. Kim, A. Fitzgibbon, KinectFusion: Real-time dense surface mapping and tracking, in: *Proc. IEEE Int. Symp. on Mixed and Augmented Reality*, 2011, pp. 127–136.
- 660 [8] N. Silberman, P. Kohli, D. Hoiem, R. Fergus, Indoor segmentation and support inference from rgb-d images, in: *Proc. Euro. Conf. on Computer Vision*, 2012, pp. 746–760.
- 661 [9] M. Fisher, M. Savva, P. Hanrahan, Characterizing structural relationships in scenes using graph kernels, *ACM Trans. on Graphics (Proc. of SIGGRAPH)* 30 (4) (2011) 34:1–34:11.
- 662 [10] K. Xu, K. Chen, H. Fu, W.-L. Sun, S.-M. Hu, Sketch2Scene: Sketch-based co-retrieval and co-placement of 3D models, *ACM Trans. on Graphics (Proc. of SIGGRAPH)* 32 (4) (2013) 123:1–123:10.
- 663 [11] M. Gönen, E. Alpaydın, Multiple kernel learning algorithms, *Journal of Machine Learning Research* 12 (2011) 2211–2268.
- 664 [12] A. Shamir, A survey on mesh segmentation techniques, *Computer Graphics Forum* 27 (6) (2008) 1539–1556.
- 665 [13] E. Kalogerakis, A. Hertzmann, K. Singh, Learning 3d mesh segmentation and labeling, *ACM Trans. Graph.* 29 (2010) 102:1–102:12.
- 666 [14] M. Johnson-Roberson, J. Bohg, M. Björkman, D. Kragic, Attention-based active 3d point cloud segmentation., in: *Proc. IEEE Int. Conf. on Intelligent Robots & Systems*, 2010, pp. 1165–1170.
- 667 [15] R. Schnabel, R. Wahl, R. Klein, Efficient RANSAC for point-cloud shape detection, *Computer Graphics Forum* 26 (2) (2007) 214–226.
- 668 [16] A. Berner, M. Bokeloh, M. Wand, A. Schilling, H.-P. Seidel, A graph-based approach to symmetry detection., in: *Volume Graphics*, Vol. 40, 2008, pp. 1–8.
- 669 [17] A. Golovinskiy, T. Funkhouser, Min-cut based segmentation of point clouds, in: *Proc. Int. Conf. on Computer Vision*, IEEE, 2009, pp. 39–46.
- 670 [18] L. A. Alexandre, 3d descriptors for object and category recognition: a comparative evaluation, in: *Proc. IEEE Int. Conf. on Intelligent Robots & Systems*, Vol. 1, pp. 1–6.
- 671 [19] K. Lai, L. Bo, X. Ren, D. Fox, Rgb-d object recognition: Features, algorithms, and a large scale benchmark, in: *Consumer Depth Cameras for Computer Vision*, Springer, 2013, pp. 167–192.
- 672 [20] K. Xu, V. G. Kim, Q. Huang, E. Kalogerakis, Data-driven shape analysis and processing, *Computer Graphics Forum* (2015) to appear.
- 673 [21] S. Song, J. Xiao, Sliding shapes for 3d object detection in depth images, in: *Proc. Euro. Conf. on Computer Vision*, Springer, 2014, pp. 634–651.
- 674 [22] K. Chen, Y.-K. Lai, Y.-X. Wu, R. Martin, S.-M. Hu, Automatic semantic modeling of indoor scenes from low-quality rgb-d data using contextual information, *ACM Trans. on Graphics (Proc. of SIGGRAPH Asia)* 33 (6) (2014) 208:1–208:15.
- 675 [23] J. S. Rudolph Triebel, Roland Siegwart, Unsupervised discovery of repetitive objects, in: *Proc. IEEE Int. Conf. on Robotics & Automation*, 2010, pp. 5041 – 5046.
- 676 [24] Y. Zhang, W. Xu, Y. Tong, K. Zhou, Online structure analysis for real-time indoor scene reconstruction, *ACM Trans. on Graphics* 32 (3) 159:1–159:12.
- 677 [25] J. Biswas, M. Veloso, Planar polygon extraction and merging from depth images, in: *Proc. IEEE Int. Conf. on Intelligent Robots & Systems*, IEEE, 2012, pp. 3859–3864.
- 678 [26] M. Dou, L. Guan, J.-M. Frahm, H. Fuchs, Exploring high-level plane primitives for indoor 3d reconstruction with a hand-held rgb-d camera, in: *Computer Vision-ACCV 2012 Workshops*, Springer, 2013, pp. 94–108.
- 679 [27] A. Monszpart, N. Mellado, G. Brostow, N. Mitra, RAPter: Rebuilding man-made scenes with regular arrangements of planes 34 (2015) 103:1–103:12.
- 680 [28] J. Papon, A. Abramov, M. Schoeler, F. Worgotter, Voxel cloud connectivity segmentation-supervoxels for point clouds, in: *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*, IEEE, 2013, pp. 2027–2034.
- 681 [29] S. Katz, A. Tal, Hierarchical mesh decomposition using fuzzy clustering and cuts, *ACM Trans. on Graphics (Proc. of SIGGRAPH)* 22 (3) (2003) 954–961.
- 682 [30] A. Golovinskiy, V. G. Kim, T. A. Funkhouser, Shape-based

- 751 recognition of 3d point clouds in urban environments, in: Proc.
752 Int. Conf. on Computer Vision, 2009, pp. 2154–2161.
- 753 [31] Y. Cheng, Mean shift, mode seeking, and clustering, IEEE
754 Trans. Pattern Analysis & Machine Intelligence 17 (8) (1995)
755 790–799.
- 756 [32] N. J. Mitra, M. Wand, H. Zhang, D. Cohen-Or, V. Kim, Q.-
757 X. Huang, Structure-aware shape processing, in: ACM SIG-
758 GRAPH 2014 Courses, 2014.
- 759 [33] K. Xu, R. Ma, H. Zhang, C. Zhu, A. Shamir, D. Cohen-Or,
760 H. Huang, Organizing heterogeneous scene collection through
761 contextual focal points, ACM Trans. on Graphics (Proc. of SIG-
762 GRAPH) 33 (4) (2014) 35:1–35:12.
- 763 [34] C. Zhu, X. Liu, Q. Liu, Y. Ming, J. Yin, Distance based multiple
764 kernel elm: A fast multiple kernel learning approach, Mathe-
765 matical Problems in Engineering 2015.
- 766 [35] M. Kazhdan, H. Hoppe, Screened poisson surface reconstruc-
767 tion, ACM Trans. on Graphics 32 (3) (2013) 29:1–29:13.
- 768 [36] A. Karpathy, S. Miller, L. Fei-Fei, Object discovery in 3d scenes
769 via shape analysis, in: Proc. IEEE Int. Conf. on Robotics &
770 Automation, IEEE, 2013, pp. 2088–2095.
- 771 [37] A. Rakotomamonjy, F. Bach, S. Canu, Y. Grandvalet, Sim-
772 pleMKL, Journal of Machine Learning Research 9 (2008) 2491–
773 2521.
- 774 [38] J. Chen, D. Bautembach, S. Izadi, Scalable real-time volumet-
775 ric surface reconstruction, ACM Trans. on Graphics (Proc. of
776 SIGGRAPH) 32 (4) (2013) 113:1–113:16.
- 777 [39] D. Sparks, Euclidean cluster analysis, Applied Statistics (1973)
778 126–130.