

Labeling 3D scenes for Personal Assistant Robots

Hema Swetha Koppula*, Abhishek Anand*, Thorsten Joachims, Ashutosh Saxena

Department of Computer Science, Cornell University.

{hema, aa755, tj, asaxena}@cs.cornell.edu

Abstract—Inexpensive RGB-D cameras that give an RGB image together with depth data have become widely available. We use this data to build 3D point clouds of a full scene. In this paper, we address the task of labeling objects in this 3D point cloud of a complete indoor scene such as an office. We propose a graphical model that captures various features and contextual relations, including the local visual appearance and shape cues, object co-occurrence relationships and geometric relationships. With a large number of object classes and relations, the model’s parsimony becomes important and we address that by using multiple types of edge potentials. The model admits efficient approximate inference, and we train it using a maximum-margin learning approach. In our experiments over a total of 52 3D scenes of homes and offices (composed from about 550 views, having 2495 segments labeled with 27 object classes), we get a performance of 84.06% in labeling 17 object classes for offices, and 73.38% in labeling 17 object classes for home scenes. Finally, we applied these algorithms successfully on a mobile robot for the task of finding an object in a large cluttered room.

I. INTRODUCTION

Inexpensive RGB-D sensors that augment an RGB image with depth data have recently become widely available. At the same time, years of research on SLAM (Simultaneous Localization and Mapping) now make it possible to merge multiple RGB-D images into a single point cloud, easily providing an approximate 3D model of a complete indoor scene (e.g., a room). In this paper, we explore how this move from part-of-scene 2D images to full-scene 3D point clouds can improve the richness of models for object labeling.

In the past, a significant amount of work has been done in semantic labeling of 2D images. However, a lot of valuable information about the shape and geometric layout of objects is lost when a 2D image is formed from the corresponding 3D world. A classifier that has access to a full 3D model, can access important geometric properties in addition to the local shape and appearance of an object. For example, many objects occur in characteristic relative geometric configurations (e.g., a monitor is almost always on a table), and many objects consist of visually distinct parts that occur in a certain relative configuration. More generally, a 3D model makes it easy to reason about a variety of properties, which are based on 3D distances, volume and local convexity.

In our work, we first use SLAM in order to compose multiple views from a Microsoft Kinect RGB-D sensor together into one 3D point cloud, providing each RGB pixel with an absolute 3D location in the scene. We then (over-)segment the scene and predict semantic labels for each segment (see Fig. 1). We predict not only coarse classes like in [1, 2] (i.e.,

wall, ground, ceiling, building), but also label individual objects (e.g., printer, keyboard, mouse). Furthermore, we model rich relational information beyond an associative coupling of labels [2].

In this paper, we propose and evaluate the first model and learning algorithm for scene understanding that exploits rich relational information derived from the full-scene 3D point cloud for object labeling. In particular, we propose a graphical model that naturally captures the geometric relationships of a 3D scene. Each 3D segment is associated with a node, and pairwise potentials model the relationships between segments (e.g., co-planarity, convexity, visual similarity, object occurrences and proximity). The model admits efficient approximate inference [3], and we show that it can be trained using a maximum-margin approach [4, 5, 6] that globally minimizes an upper bound on the training loss. We model both associative and non-associative coupling of labels. With a large number of object classes, the model’s parsimony becomes important. Some features are better indicators of label similarity, while other features are better indicators of non-associative relations such as geometric arrangement (e.g., “on top of,” “in front of”). We therefore model them using appropriate clique potentials rather than using general clique potentials. Our model is highly flexible and we have made our software available for download to other researchers in this emerging area of 3D scene understanding.

To empirically evaluate our model and algorithms, we perform several experiments over a total of 52 scenes of two types: offices and homes. These scenes were built from about 550 views from the Kinect sensor, and they will also be made available for public use. We consider labeling each segment (from a total of about 50 segments per scene) into 27 classes (17 for offices and 17 for homes, with 7 classes in common). Our experiments show that our method, which captures several local cues and contextual properties, achieves an overall performance of 84.06% on office scenes and 73.38% on home scenes. We also consider the problem of labeling 3D segments with multiple attributes meaningful to robotics context (such as small objects that can be manipulated, furniture, etc.). Finally, we successfully applied these algorithms on a mobile robot for the task of finding an object in a large cluttered room.

II. RELATED WORK

There is a huge body of work in the area of scene understanding and object recognition from 2D images. Previous works focus on several different aspects: designing good local features such as HOG (histogram-of-gradients) [7] and bag of words [8], designing good global (context) features such as

* indicates equal contribution.

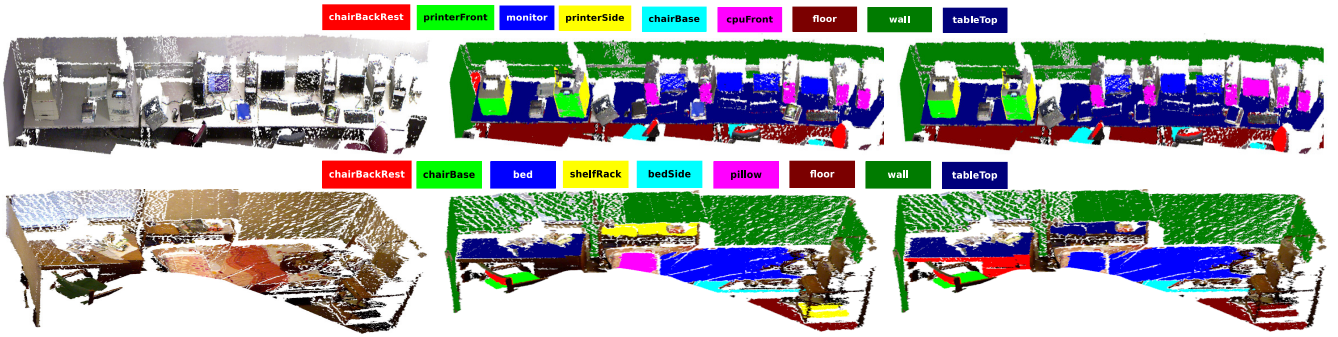


Fig. 1. Office scene (top) and Home scene (bottom with the corresponding label coloring above the images). The left-most is the original point cloud, the middle is the ground truth labeling and the right most is the point cloud with predicted labels.

GIST features [9], and combining multiple tasks [10]. However, these approaches do not consider the relative arrangement of the parts of the object or of multiple objects with respect to each other. A number of works propose models that explicitly capture the relations between different parts of the object [11], and between different objects in 2D images [12, 13]. However, a lot of valuable information about the shape and geometric layout of objects is lost when a 2D image is formed from the corresponding 3D world. In some recent works, 3D layout or depths have been used for improving object detection (e.g., [14, 15, 16, 17, 18]). Here a rough 3D scene geometry (e.g., the main surfaces in a scene) is inferred from a single 2D image or a stereo video, respectively. However, the estimated geometry is not accurate enough to give significant improvements. With 3D data, we can more precisely determine the shape, size and geometric orientation of the objects, and several other properties and therefore capture much stronger context.

The recent availability of synchronized videos of both color and depth obtained from RGB-D (Kinect-style) depth cameras, shifted the focus to making use of both visual as well as shape features for object detection [19, 20, 21, 22, 23] and 3D segmentation (e.g., [24]). These methods demonstrate that augmenting visual features with 3D information can enhance object detection in cluttered, real-world environments. However, these works do not make use of the contextual relationships between various objects which have been shown to be useful for tasks such as object detection and scene understanding in 2D images. Our goal is to perform semantic labeling of indoor 3D scenes by modeling and learning several contextual relationships.

There is also some recent work in labeling outdoor scenes obtained from LIDAR data into a few geometric classes (e.g., ground, building, trees, vegetation, etc.). [25, 26] capture context by designing node features and [27] do so by stacking layers of classifiers; however these methods do not model the correlation between the labels. Some of these works model some contextual relationships in the learning model itself. For example, [2, 28] use associative Markov networks in order to favor similar labels for nodes in the cliques. However, many relative features between objects are not associative in nature. For example, the relationship “on top of” does not hold in between two ground segments, i.e., a ground segment cannot be “on top of” another ground segment. Therefore, using an

associative Markov network is very restrictive for our problem. All of these works [2, 26, 27, 28] were designed for outdoor scenes with LIDAR data (without RGB values) and therefore would not apply directly to RGB-D data in indoor environments. Furthermore, these methods only consider very few geometric classes (between three to five classes) in outdoor environments, whereas we consider a large number of object classes for labeling the indoor RGB-D data.

The most related work to ours is [1], where they label the planar patches in a point-cloud of an indoor scene with four geometric labels (walls, floors, ceilings, clutter). They use a CRF to model geometrical relationships such as orthogonal, parallel, adjacent, and coplanar. The learning method for estimating the parameters was based on maximizing the pseudo-likelihood resulting in a sub-optimal learning algorithm. In comparison, our basic representation is 3D segments (as compared to planar patches) and we consider a much larger number of classes (beyond just the geometric classes). We capture a much richer set of relationships between pairs of objects, and use a principled max-margin learning method to learn the parameters of our model.

III. APPROACH

We now outline our approach, including the model, its inference methods, and the learning algorithm. Our input is multiple Kinect RGB-D images of an indoor scene stitched into a single 3D point cloud using RGBDSLAM [29]. Each such point cloud is then over-segmented based on smoothness (i.e., difference in the local surface normals) and continuity of surfaces (i.e., distance between the points). These segments are the atomic units in our model. Our goal is to label each of them.

Before getting into the technical details of the model, the following outlines the properties we aim to capture:

Visual appearance. The reasonable success of object detection in 2D images shows that visual appearance is a good indicator for labeling scenes. We therefore model the local color, texture, gradients of intensities, etc. for predicting the labels. In addition, we also model the property that if nearby segments are similar in visual appearance, they are more likely to belong to the same object.

Local shape and geometry. Objects have characteristic shapes—for example, a table is horizontal, a monitor is vertical, a keyboard is uneven, and a sofa is usually smoothly

curved. Furthermore, parts of an object often form a convex shape. We compute 3D shape features to capture this.

Geometrical context. Many sets of objects occur in characteristic relative geometric configurations. For example, a monitor is always *on-top-of* a table, chairs are usually found *near* tables, a keyboard is *in-front-of* a monitor. This means that our model needs to capture *non-associative* relationships (i.e., that neighboring segments differ in their labels in specific patterns).

Note that the examples given above are just illustrative. For any particular practical application, there will likely be other properties that could also be included. As demonstrated in the following section, our model is flexible enough to include a wide range of features.

A. Model Formulation

We model the three-dimensional structure of a scene using a model isomorphic to a Markov Random Field with log-linear node and pairwise edge potentials. Given a segmented point cloud $\mathbf{x} = (x_1, \dots, x_N)$ consisting of segments x_i , we aim to predict a labeling $\mathbf{y} = (y_1, \dots, y_N)$ for the segments. Each segment label y_i is itself a vector of K binary class labels $y_i = (y_i^1, \dots, y_i^K)$, with each $y_i^k \in \{0, 1\}$ indicating whether a segment i is a member of class k . Note that multiple y_i^k can be 1 for each segment (e.g., a segment can be both a “chair” and a “movable object”). We use such multi-labelings in our attribute experiments where each segment can have multiple attributes, but not in segment labeling experiments where each segment can have only one label.

For a segmented point cloud \mathbf{x} , the prediction $\hat{\mathbf{y}}$ is computed as the argmax of a discriminant function $f_{\mathbf{w}}(\mathbf{x}, \mathbf{y})$ that is parameterized by a vector of weights \mathbf{w} .

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} f_{\mathbf{w}}(\mathbf{x}, \mathbf{y}) \quad (1)$$

The discriminant function captures the dependencies between segment labels as defined by an undirected graph $(\mathcal{V}, \mathcal{E})$ of vertices $\mathcal{V} = \{1, \dots, N\}$ and edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$. We describe in Section III-B how this graph is derived from the spatial proximity of the segments. Given $(\mathcal{V}, \mathcal{E})$, we define the following discriminant function based on individual segment features $\phi_n(i)$ and edge features $\phi_t(i, j)$ as further described below.

$$f_{\mathbf{w}}(\mathbf{y}, \mathbf{x}) = \sum_{i \in \mathcal{V}} \sum_{k=1}^K y_i^k \left[w_n^k \cdot \phi_n(i) \right] + \sum_{(i,j) \in \mathcal{E}} \sum_{T_t \in \mathcal{T}} \sum_{(l,k) \in T_t} y_i^l y_j^k \left[w_t^{lk} \cdot \phi_t(i, j) \right] \quad (2)$$

The node feature map $\phi_n(i)$ describes segment i through a vector of features, and there is one weight vector for each of the K classes. Examples of such features are the ones capturing local visual appearance, shape and geometry. The edge feature maps $\phi_t(i, j)$ describe the relationship between segments i and j . Examples of edge features are the ones capturing similarity in visual appearance and geometric context.¹ There may be multiple types t of edge feature maps

¹Even though it is not represented in the notation, note that both the node feature map $\phi_n(i)$ and the edge feature maps $\phi_t(i, j)$ can compute their features based on the full \mathbf{x} , not just x_i and x_j .

$\phi_t(i, j)$, and each type has a graph over the K classes with edges T_t . If T_t contains an edge between classes l and k , then this feature map and a weight vector w_t^{lk} is used to model the dependencies between classes l and k . If the edge is not present in T_t , then $\phi_t(i, j)$ is not used.

We say that a type t of edge features is modeled by an associative edge potential if $T_t = \{(k, k) | \forall k = 1..K\}$. And it is modeled by a non-associative edge potential if $T_t = \{(l, k) | \forall l, k = 1..K\}$. Finally, it is modeled by an object-associative edge potential if $T_t = \{(l, k) | \exists \text{object}, l, k \in \text{parts}(\text{object})\}$

Parsimonious model. In our experiments we distinguished between two types of edge feature maps—“object-associative” features $\phi_{oa}(i, j)$ used between classes that are parts of the same object (e.g., “chair base”, “chair back” and “chair back rest”), and “non-associative” features $\phi_{na}(i, j)$ that are used between any pair of classes. Examples of features in the object-associative feature map $\phi_{oa}(i, j)$ include similarity in appearance, co-planarity, and convexity—i.e., features that indicate whether two adjacent segments belong to the same class or object. A key reason for distinguishing between object-associative and non-associative features is parsimony of the model. In this parsimonious model (referred to as `svm_mrf_parsimon`), we model object associative features using object-associative edge potentials and non-associative features as non-associative edge potentials. As not all edge features are “non-associative”, we avoid learning weight vectors for relationships which do not exist. Note that $|T_{na}| \gg |T_{oa}|$ since, in practice, the number of parts of an object is much less than K . Due to this, the model we learn with both type of edge features will have much lesser number of parameters compared to a model learnt with all edge features as “non-associative features”.

B. Features

Table I summarizes the features used in our experiments. $\lambda_{i0}, \lambda_{i1}$ and λ_{i2} are the 3 eigen-values of the scatter matrix computed from the points of segment i in increasing order. c_i is the centroid of segment i . r_i is the ray vector to the c_i from the camera in which it was captured. rh_i is the projection of r_i on horizontal plane. \hat{n}_i is the unit normal of segment i which points towards the camera ($r_i \cdot \hat{n}_i < 0$).

The node features $\phi_n(i)$ consist of visual appearance features based on histogram of HSV values and the histogram of gradients (HOG), as well as local shape and geometry features that capture properties such as how planar a segment is, its absolute location above ground, and its shape. Some features capture spatial location of an object in the scene (e.g., N9).

We connect two segments (nodes) i and j by an edge if there exists a point in segment i and a point in segment j which are less than *context_range* distance apart. This captures the closest distance between two segments (as compared to centroid distance between the segments)—we study the effect of context range more in Section IV. The edge features $\phi_t(i, j)$ (Table I-right) consist of associative features (E1) based on visual appearance and local shape, as well as non-associative

TABLE I
Node features for segment i .

Description	Count
Visual Appearance	48
N1. Histogram of HSV color values	14
N2. Average HSV color values	3
N3. Average of HOG features of the blocks in image spanned by the points of a segment	31
Local Shape and Geometry	8
N4. linearness ($\lambda_{i0} - \lambda_{i1}$), planariness ($\lambda_{i1} - \lambda_{i2}$), Scatter: λ_{i0}	3
N6. Vertical component of the normal: \hat{n}_{iz}	1
N7. Vertical position of centroid: c_{iz}	1
N8. Vert. and Hor. extent of bounding box	2
N9. Dist. from the scene boundary	1
Features for edge (segment i , segment j).	
Visual Appearance (associative)	3
E1. Difference of avg HSV color values	3
Local Shape and Geometry (associative)	2
E2. Coplanarity and convexity	2
Geometric context (non-associative)	6
E3. Horizontal distance b/w centroids.	1
E4. Vertical Displacement b/w centroids: $(c_{iz} - c_{jz})$	1
E5. Angle between normals (Dot product): $\hat{n}_i \cdot \hat{n}_j$	1
E6. Diff. in angle with vert.: $\cos^{-1}(n_{iz}) - \cos^{-1}(n_{jz})$	1
E7. Dist. between closest points	1
E8. rel. position from camera (in front of/behind).	1

features (E3-E8) that capture the tendencies of two objects to occur in certain configurations. Note that our features are insensitive to horizontal translation and rotation of the camera. However, our features place a lot of emphasis on the vertical direction because gravity influences the shape and relative positions of objects to a large extent.

C. Learning and Inference

Solving the argmax in Eq. 1 for the discriminant function in Eq. ?? is NP hard. It can be formulated as the following mixed-integer program, which can be solved by a general-purpose MIP solver² in about 20 minutes on a typical scene.

$$\begin{aligned}
\hat{y} = \underset{y}{\operatorname{argmax}} \max_z \sum_{i \in \mathcal{V}} \sum_{k=1}^K y_i^k \left[w_n^k \cdot \phi_n(i) \right] \\
+ \sum_{(i,j) \in \mathcal{E}} \sum_{T_t \in \mathcal{T}} \sum_{(l,k) \in T_t} z_{ij}^{lk} \left[w_t^{lk} \cdot \phi_t(i,j) \right] \quad (3) \\
\forall i, j, l, k : z_{ij}^{lk} \leq y_i^l, z_{ij}^{lk} \leq y_j^k, \quad y_i^l + y_j^k \leq z_{ij}^{lk} + 1 \\
z_{ij}^{lk}, y_i^l \in \{0, 1\}, \quad \forall i : \sum_{j=1}^K y_i^j = 1 \quad (4)
\end{aligned}$$

However, if we remove the last constraint (??), and relax the variables z_{ij}^{lk} and y_i^l to the interval $[0, 1]$, we get a linear relaxation that can be shown to always have half-integral solutions (i.e. y_i^l only take values $\{0, 0.5, 1\}$ at the solution) [30]. Furthermore, this relaxation can also be solved as a quadratic pseudo-Boolean optimization problem using a graph-cut method [3], which is orders of magnitude faster than using a general purpose LP solver (i.e., 2 sec for labeling a typical full scene in our experiments, and 0.2 sec for a single view). For training, we use the *SVM^{struct}*³ software which uses the cutting plane method to jointly learn values of w_n and w_t 's so as to minimize a regularized upper bound on the training error.

IV. EXPERIMENTS

A. Data

We consider labeling object segments in full 3D scene (as compared to 2.5D data from a single view). For this

purpose, we collected data of 24 office and 28 home scenes. Each scene was reconstructed from about 8-9 RGB-D views from a Microsoft Kinect sensor and we have a total of about 550 views. Each scene contains about a million colored points. We first over-segment the 3D scene (as described earlier) to obtain the atomic units of our representation. For training, we manually labeled the segments, and we selected the labels which were present in a minimum of 5 scenes in the dataset. Specifically, the office labels are: *{wall, floor, tableTop, tableDrawer, tableLeg, chairBackRest, chairBase, chairBack, monitor, printerFront, printerSide, keyboard, cpuTop, cpuFront, cpuSide, book, paper}*, and the home labels are: *{wall, floor, tableTop, tableDrawer, tableLeg, chairBackRest, chairBase, sofaBase, sofaArm, sofaBackRest, bed, bedSide, quilt, pillow, shelfRack, laptop, book}*. This gave us a total of 1108 labeled segments in the office scenes and 1387 segments in the home scenes. Often one object may be divided into multiple segments because of over-segmentation. We have made this data available at: <http://pr.cs.cornell.edu/sceneunderstanding>.

B. Results

Table II shows the results, performed using 4-fold cross-validation and averaging performance across the folds for the models trained separately on home and office datasets. We use both the macro and micro averaging to aggregate precision and recall over various classes. Since our algorithm can only predict one label per segment, micro precision and recall are same as the percentage of correctly classified segments. Macro precision and recall are respectively the averages of precision and recall for all classes. The optimal C value is determined separately for each of the algorithms by cross-validation.

Figure 1 shows the original point cloud, ground-truth and predicted labels for one office (top) and one home scene (bottom). We see that on majority of the classes our model predicts the correct label. It makes mistakes on some tricky cases, such as a pillow getting confused with the bed, and table-top getting confused with the shelf-rack.

One of our goals is to study the effect of various factors, and therefore we compared various versions of the algorithms with various settings. We discuss them in the following.

Do Image and Point-Cloud Features Capture Complementary Information? The RGB-D data contains both image and depth information, and enables us to compute a wide variety of features. In this experiment, we compare the two kinds of features: Image (RGB) and Shape (Point Cloud) features. To show the effect of the features independent of the effect of context, we only use the node potentials from our model, referred to as *svm_node_only* in Table II. The *svm_node_only* model is equivalent to the multi-class SVM formulation [31]. Table II shows that Shape features are more effective compared to the Image, and the combination works better on both precision and recall. This indicates that the two types of features offer complementary information and their combination is better for our classification task.

How Important is Context? Using our *svm_mrf_parsimon* model as described in Section III-A, we show significant

²<http://www.tfinley.net/software/pyglpk/readme.html>

³http://svmlight.joachims.org/svm_struct.html

TABLE II
AVERAGE MICRO PRECISION/RECALL, AVERAGE MACRO PRECISION AND RECALL FOR HOME AND OFFICE SCENES.

features	algorithm	Office Scenes			Home Scenes		
		micro	macro	Recall	micro	macro	Recall
None	chance	26.23	5.88	5.88	29.38	5.88	5.88
Image Only	svm_node_only	46.67	35.73	31.67	38.00	15.03	14.50
Shape Only	svm_node_only	75.36	64.56	60.88	56.25	35.90	36.52
Image+Shape	svm_node_only	77.97	69.44	66.23	56.50	37.18	34.73
Image+Shape & context	single_frames	84.32	77.84	68.12	69.13	47.84	43.62
Image+Shape & context	svm_mrf_assoc	75.94	63.89	61.79	62.50	44.65	38.34
Image+Shape & context	svm_mrf_nonassoc	81.45	76.79	70.07	72.38	57.82	53.62
Image+Shape & context	svm_mrf_parsimon	84.06	80.52	72.64	73.38	56.81	54.80

improvements in the performance over using svm_node_only model on both datasets. In office scenes, the micro precision increased by 6.09% over the best svm_node_only model that does not use any context. In home scenes the increase is much higher, 16.88%.

The type of contextual relations we capture depend on the type of edge potentials we model. To study this, we compared our method with models using only associative (svm_mrf_assoc) or only non-associative (svm_mrf_nonassoc) edge potentials. We observed that modeling all edge features using associative potentials is poor compared to our full model. In fact, using only associative potentials showed a drop in performance compared to svm_nodeonly model on the office dataset. This indicates it is important to capture the relations between regions having *different* labels. Our svm_mrf_non_assoc model does so, by modeling all edge features using non-associative potentials, which can favour or disfavour labels of different classes for nearby segments. It gives higher precision and recall compared to svm_nodeonly and svm_mrf_assoc.

However, not all the edge features are non-associative in nature, modeling them using only non-associative potentials could be an overkill (each non-associative feature adds K^2 more parameters to be learnt). Therefore using our svm_mrf_parsimon model to model these relations achieves higher performance in both datasets.

How Large should the Context Range be?

Context relationships of different objects can be meaningful for different spatial distances. This range may vary depending on the environment as well. For example, in an office, keyboard and monitor go together, but they may have little

relation with a sofa that is slightly farther away. In a house, sofa and table may go together even if they are farther away.

In order to study this, we compared our svm_mrf_parsimon with varying context range for determining the neighborhood (see Figure 2 for average micro precision vs range plot). Note that the context range is determined from the boundary of one segment to the boundary of the other, and hence it is somewhat

independent of the size of the object. We note that increasing the context range increases the performance to some level, and then it drops slightly. We attribute this to the fact that with increasing the context range, irrelevant objects may get an edge in the graph, and with limited training data, spurious relationships may be learned. We observe that the optimal context range for office scenes is around 0.3 meters and 0.6 meters for home scenes.

How does a Full 3D Model Compare to a 2.5D Model? In Table II, we compare the performance of our full model with a model that was trained and tested on single views of the same scene. During the comparison, the training folds were consistent with other experiments, however the segmentation of this point-cloud was different (because the input point-cloud itself is from single view). This makes the micro precision values not meaningful because the distribution of labels is not same for the two cases. In particular, many large object in scenes (e.g., wall, ground) get split up into multiple segments in single views. We observed that the macro precision and recall are higher when multiple views are combined to form the scene. We attribute the improvement in macro precision and recall to the fact that larger scenes have more context, and models are more complete because of multiple views.

What is the Effect of the Inference Method? The results for svm_mrf algorithms Table II were generated using the MIP solver. The graph-cut algorithm however, gives a higher precision and lower recall on both datasets. For example, on office data, the graphcut inference for our svm_mrf_parsimon gave a micro precision of 90.25 and micro recall of 61.74. Here, the micro precision and recall are not same as some of the segments might not get any label. Since it is orders of magnitude faster, it is ideal for realtime robotic applications.

C. Robotic experiments

The ability to label segments is very useful for robotics applications, for example, in detecting objects (so that a robot can find/retrieve an object on request) or for other robotic tasks such as manipulation. We therefore performed two relevant robotic experiments.

Attribute Learning: In some robotic tasks, such as robotic grasping [32] or placing [33], it is not important to know the exact object category, but just knowing a few attributes of an object may be useful. For example, if a robot has to clean a floor, it would help if it knows which objects it can move and which it cannot. If it has to place an object, it should place them on horizontal surfaces, preferably where

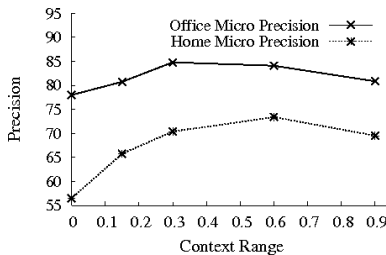


Fig. 2. Effect of context range on precision (=recall here).

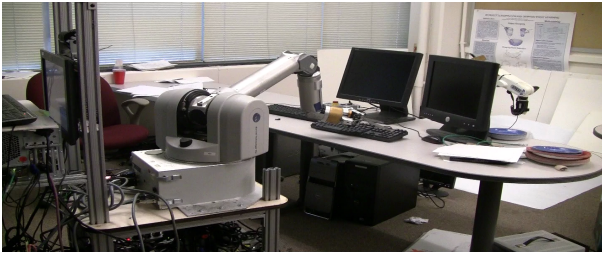


Fig. 3. Cornell's POLAR (PersOnAL Assistant Robot) using our classifier for detecting a keyboard in a cluttered room.

humans do not sit. With this motivation we have designed 8 attributes, each for the home and office scenes, giving a total of 10 unique attributes, comprised of: *wall, floor, flat-horizontal-surfaces, furniture, fabric, heavy, seating-areas, small-objects, table-top-objects, electronics*. Note that each segment in the point cloud can have multiple attributes and therefore we can learn these attributes using our model which naturally allows multiple labels per segment. We compute the precision and recall over the attributes by counting how many attributes were correctly inferred. In home scenes we obtained a precision of 83.12% and 70.03% recall, and in the office scenes we obtain 87.92% precision and 71.93% recall.

Robotic Object Detection: We finally use our algorithm on a mobile robot, mounted with a Kinect, for completing the goal of finding an object such as a keyboard in an extremely cluttered room (Fig. 3). The following video shows our robot successfully finding the keyboard in an office: <http://pr.cs.cornell.edu/sceneunderstanding>

In conclusion, we have proposed and evaluated the first model and learning algorithm for scene understanding that exploits rich relational information from full-scene 3D point clouds. We applied this technique to object labeling problem, and studied affects of various factors on a large dataset. Our robotic applications shows that such inexpensive RGB-D sensors can be quite useful for scene understanding by robots.

Acknowledgements. We thank Gaurab Basu, Matthew Cong and Yun Jiang for the robotic experiments.

REFERENCES

- [1] X. Xiong and D. Huber, "Using context to create semantic 3d models of indoor environments," in *BMVC*, 2010.
- [2] D. Anguelov, B. Taskar, V. Chatalbashev, D. Koller, D. Gupta, G. Heitz, and A. Ng, "Discriminative learning of markov random fields for segmentation of 3d scan data," in *CVPR*, 2005.
- [3] C. Rother, V. Kolmogorov, V. Lempitsky, and M. Szmur, "Optimizing binary mrfs via extended roof duality," in *CVPR*, 2007.
- [4] B. Taskar, V. Chatalbashev, and D. Koller, "Learning associative markov networks," in *ICML*, 2004.
- [5] I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun, "Support vector machine learning for interdependent and structured output spaces," in *ICML*, 2004.
- [6] T. Finley and T. Joachims, "Training structural svms when exact inference is intractable," in *ICML*, 2008.
- [7] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *CVPR*, 2005.
- [8] G. Csurka, C. Dance, L. Fan, J. Willamowski, and C. Bray, "Visual categorization with bags of keypoints," in *Workshop on statistical learning in computer vision, ECCV*, 2004.
- [9] A. Torralba, "Contextual priming for object detection," *IJCV*, vol. 53, no. 2, pp. 169–191, 2003.
- [10] C. Li, A. Kowdle, A. Saxena, and T. Chen, "Towards holistic scene understanding: Feedback enabled cascaded classification models," in *NIPS*, 2010.
- [11] P. Felzenszwalb, D. McAllester, and D. Ramanan, "A discriminatively trained, multiscale, deformable part model," in *CVPR*, 2008.
- [12] K. P. Murphy, A. Torralba, and W. T. Freeman, "Using the forest to see the trees: a graphical model relating features, objects and scenes," *NIPS*, 2003.
- [13] G. Heitz and D. Koller, "Learning spatial context: Using stuff to find things," in *ECCV*, 2008.
- [14] A. Saxena, S. H. Chung, and A. Y. Ng, "Learning depth from single monocular images," in *NIPS*, 2005.
- [15] D. Hoiem, A. A. Efros, and M. Hebert, "Putting objects in perspective," in *In CVPR*, 2006.
- [16] A. Saxena, M. Sun, and A. Y. Ng, "Make3d: Learning 3d scene structure from a single still image," *IEEE PAMI*, vol. 31, no. 5, pp. 824–840, 2009.
- [17] D. C. Lee, A. Gupta, M. Hebert, and T. Kanade, "Estimating spatial layout of rooms using volumetric reasoning about objects and surfaces," in *NIPS*, 2010.
- [18] B. Leibe, N. Cornelis, K. Cornelis, and L. V. Gool, "Dynamic 3d scene analysis from a moving vehicle," in *CVPR*, 2007.
- [19] S. Gould, P. Baumstarck, M. Quigley, A. Y. Ng, and D. Koller, "Integrating Visual and Range Data for Robotic Object Detection," in *ECCV workshop Multi-camera Multi-modal (M2SFA2)*, 2008.
- [20] M. Quigley, S. Batra, S. Gould, E. Klingbeil, Q. V. Le, A. Wellman, and A. Y. Ng, "High-accuracy 3d sensing for mobile manipulation: Improving object detection and door opening," in *ICRA*, 2009.
- [21] R. B. Rusu, Z. C. Marton, N. Blodow, M. Dolha, and M. Beetz, "Towards 3d point cloud based object maps for household environments," *Robot. Auton. Syst.*, vol. 56, pp. 927–941, 2008.
- [22] K. Lai, L. Bo, X. Ren, and D. Fox, "A Large-Scale Hierarchical Multi-View RGB-D Object Dataset," in *ICRA*, 2011.
- [23] —, "Sparse Distance Learning for Object Recognition Combining RGB and Depth Information," in *ICRA*, 2011.
- [24] A. Collet Romea, S. Srinivasa, and M. Hebert, "Structure discovery in multi-modal data : a region-based approach," in *ICRA*, 2011.
- [25] A. Golovinskiy, V. G. Kim, and T. Funkhouser, "Shape-based recognition of 3d point clouds in urban environments," *ICCV*, 2009.
- [26] R. Shapovalov, A. Velizhev, and O. Barinova, "Non-associative markov networks for 3d point cloud classification," in *ISPRIS Comm III Symp PCV*, 2010.
- [27] X. Xiong, D. Munoz, J. A. Bagnell, and M. Hebert, "3-d scene analysis via sequenced predictions over points and regions," in *ICRA*, 2011.
- [28] D. Munoz, N. Vandapel, and M. Hebert, "Onboard contextual classification of 3-d point clouds with learned high-order markov random fields," in *ICRA*, 2009.
- [29] "RGBDSLAM," <http://openslam.org/rgbdslam.html>.
- [30] P. Hammer, P. Hansen, and B. Simeone, "Roof duality, complementation and persistency in quadratic 0–1 optimization," *Mathematical Programming*, vol. 28, no. 2, pp. 121–155, 1984.
- [31] T. Joachims, T. Finley, and C. Yu, "Cutting-plane training of structural SVMs," *Machine Learning*, vol. 77, no. 1, p. 27, 2009.
- [32] Y. Jiang, S. Moseson, and A. Saxena, "Efficient grasping from rgbd images: Learning using a new rectangle representation," in *ICRA*, 2011.
- [33] Y. Jiang, C. Zheng, M. Lim, and A. Saxena, "Learning to place new objects," in *RSS workshop on mobile manipulation*, 2011.