# Behavior Recovery/Analysis in Cluttered and Dynamic Indoor Scene
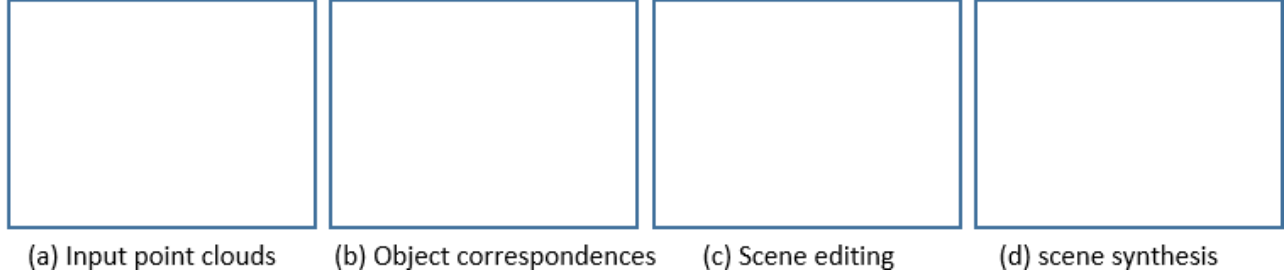


(a) Input point clouds    (b) Object correspondences    (c) Scene editing    (d) scene synthesis

**Figure 1:** *Behavior recovery and behavior-based editing of dynamic cluttered indoor scenes. From a set of dense scans at different times (a), our system first extract the object correspondences (b) and behavior model (*(xuejin: optionally show a graph)*). The recovered behavior model can be applied to many applications, such as scene editing (c) and indoor scene synthesis (d).*

## 1 Abstract

While modeling static indoor scenes using RGBD cameras has been extensively studied in recent years, we introduce a *behavior recovery* system to investigate the behavior of objects in an cluttered indoor scene. In a daily indoor environment, because of object functions and human behaviors, the spatial placements of objects presents non-unique but statistically regular displacements. We take the *inverse problem* to recover object behaviors from a collection of point clouds captured at different times in a dynamic indoor scene. Our system consists of two key parts, *extraction of object correspondence* and *behavior model*. Given a collection of dense point clouds of an indoor scene in daily use at different times, the correspondence between objects are extracted using an iterative segmentation-and-registration process. Our algorithm is robust to noise and incomplete parts in imperfect scans. In the second part is to recover object behaviors, which represent the spatial arrangement of objects and interrelations between objects. Using our method, the correlation between concrete geometry and semantic behaviors of an indoor scene can be established. Therefore, the recovered behavior model can be adopted to many applications that requires labeled 3D database, such as scene synthesis, scene arrangement and so on. We evaluate our algorithm on a number of indoor scenes including office, bedroom and so on. The results demonstrate that our algorithm build accurate object correspondence from imperfect scans of cluttered indoor scenes, based on which, the recovered behavior provides natural principles for many other applications oriented to indoor scenes.

**Keywords:** Indoor scene, behavior analysis, dynamic, object correspondence

## 1 Introduction

Modeling indoor scenes has attracted a large amount of attentions for decades in computer graphics. Recently, many techniques have been presented to generate static 3D models for indoor scenes, including dense modeling from RGBD data [Henry et al. 2012; Izadi et al. 2011; Xiao and Furukawa 2012; Yan et al. 2014], combing object classification and modeling [Shao et al. 2012; Nan et al. 2012; Kim et al. 2012], and synthesizing of 3D indoor scenes from large collection of examples [Fisher et al. 2012; Xu et al. 2013]. While visually appealing models are obtained for rendering, it remains challenging to extract the semantics that the geometric representation essentially encodes. On the other hand, the semantics are required in many applications such as indoor scene understanding, scene editing, and etc.

Comparing with static scenes, dynamic scene analysis has significant value in interior design, animation making, etc. The manners of how furniture objects interact with each other and how furniture objects interact with users play a very important rule in interior design. Typically, the geometric representation including object model and spatial placements of objects at different times implicitly encode the object functions and human behavior in that environment. However, the dynamic indoor scene analysis has not been investigated much in computer graphics.

Though the geometry data collection carries the behavior information, it is non-trivial to extract behavior from imperfect scans by consumer-level RGBD cameras. The challenges are in two-fold because the scanned point clouds are noisy, incomplete and with errors. First, segmentation of objects with accurate boundary is tedious because the objects in a cluttered indoor environment are in a large variety of scales. Moreover, the massive occlusions and self-occlusions in a cluttered scene makes the segmentation more challenging. Second, it is an arduous task to figure out exact object correspondences given imperfect segmentations. Furthermore, there are a great of deal of similar structures in man-made objects in indoor scenes, which leads to large ambiguities of shape correspondences.

In this paper, we present a novel algorithm to explore object behaviors in cluttered indoor scenes from a set of point clouds scanned using consumer-level RGBD cameras without any training data. The consistency and difference between frames simultaneously provide valuable hints for recovering object correspondences. First, each frame is roughly segmented into patches, which are then clustered into objects as initial correspondences hypothesis. While we project the object models back to each frame with the corresponding transformation, the consistency between captured data and the project data strengthens the initial hypothesis while the difference indicates wrong correspondences. Based on this validation, each frame is resegmented into patches. By iteratively perform the segmentation and registration steps, our work converges to coherent segmentations and correct correspondences between a bunch of objects at different scales in a large collection of points clouds.

In summary, the contributions of our system are three-fold:

1. To the best of our knowledge, our system, for the first time, performs behavior analysis in a dynamic indoor scene from point clouds scanned using consumer-level RGBD cameras without any database.

2. We present a global optimization framework to combine object segmentation, correspondence extraction and behavior analysis in an iterative scheme.

3. We present a novel behavior model in dynamic indoor scenes, which can be applied directly to many appealing applications.

## 2  Related Work

Many techniques have been proposed to generate static 3D indoor scenes in computer graphics. Though none of them focus on dynamic scene analysis like our system, they provide valuable reference on the underlying techniques.

**Reconstruction from RGBD Images.**   For static scenes, Kinect-Fusion [Izadi et al. 2011] enables the real-time reconstruction by holding and moving a depth camera. For large-scale indoor scenes with multiple rooms, reconstructing a dense 3D model from the noisy and incomplete scanned range data typically involves registration of point clouds in different views and a global optimization to reduce gaps in a large scene [Xiao and Furukawa 2012; Henry et al. 2012]. Their goal is mainly to generate high-quality point clouds but without semantic analysis of the objects appear in the scene. Recently, object classification is employed to assist modeling for massive indoor scenes that containing many instances of chairs, desks, etc. Koppula et al. [2011] first introduce the learning algorithm to understand the RGBD data of an indoor scene. To further reconstruct the 3D model for a cluttered indoor scene, 3D model databases can be used as template by searching for similar 3D model and then fitting the template to the scanned data [Shao et al. 2012; Nan et al. 2012]. [Kim et al. 2012] do not manually collect 3D models to build the database. The template model is reconstructed by scanning the same object in different configuration. Each model has an additional presentation by geometric primitives. [Shao et al. 2012] trains the class model based on geometry and appearance features to segment and label the RGBD data captured under sparse views. By learned an initial model for each class of object in indoor environments from a pre-labelled database, the model are refined progressively with user-refined segmentation results. The 3D model can be generated by placing the most similar model in the database according to the RGBD data. If objects move in a scene, they can be detected and reposed by segmented and classified based on the learned model from previously reconstructed model [Liu et al. 2014]. Different with these techniques, we pay more attentions on analyzing the object behaviors from the dynamic range data.

**Reconstruction from Sequential Point clouds.**   Many techniques have been proposed to reconstruct the object surfaces from the range data sequences. [Wand et al. 2007] uses a *statistical framework* to reconstruct the geometry from real-time range scanning. Each frame is divided into 3D pieces. A statistical model is used to iteratively merge adjacent frames by aligning pieces and optimizing their shapes. However, some geometric artifacts remain due to structured outliers and in some boundary regions. [Chang and Zwicker 2011] presents a global registration algorithm to reconstruct *articulated 3D models* from dynamic range scan sequences. The surface motion is modeled by a reduced deformable model. Joints and skinning weights are solved in the system to register point clouds in different poses. (xuejin: We may also consider the furniture objects in indoor environments as articulated models, whose shapes under different poses can be deformed through connectors like hinge, slide, and so on. )  A new formulation of the ICP algorithm is proposed using sparse inducing norms [Bouaziz et al. 2013]. While it achieves superior registration result on the data

with outliers and missing region, only rigid alignment is handled. A proactive capturing is employed by asking the user to move the objects to capture both interior and exterior of a scene [Yan et al. 2014]. The correspondence between adjacent frames is built first then segmentation. Xu et al. [2015] employ a robot to move objects during scene reconstruction so that the ambiguities in object structures can be solved from the dynamic data. All these techniques take the advantage of the differences caused by motion to extract valuable and semantic information of object structure. We do not only take the motion information in object modeling, but also take its advantage of implicitly encoding object behavior in a scene.

**Functionality Analysis using Context.**   Many techniques of functional analysis of one category of objects using 3D model collections have been proposed [Huang et al. 2014; Su et al. 2014]. Besides of functionality analysis of a single object, interaction between objects has drawn more and more attentions. The Icon descriptor is proposed to represent the functionality of 3D objects with its context [Hu et al. 2015]. Human action is involved to establish the correlations between the geometry and functionality of a region [Savva et al. 2014]. With the association of object arrangements and human activities, novel 3D scenes can be synthesized towards specific functions [Fisher et al. 2015]. Inspired by these methods, we put our effort on behavior recovery, but from more challenging data. The point clouds can be easily to capture using RGBD cameras. However, the noisy and incompleteness brings tremendous challenges for extracting accurate object correspondences.

**Data-Driven Furniture Layout.**   The general way producing the layout of furniture objects is to model a set of design rules and then to optimize an energy function given constraints by individuals. [Merrell et al. 2011] formulates a group of layout guidelines in a density function according to professional manuals on furniture layout. When the user specifies the room shape and an initial arrangement of the set of furniture to be placed in the room, this system generates a number of layout suggestions by a hardware-accelerated Monte Carlo sampler. Instead of manually define the layout guidelines, the hierarchical and spatial relationships of the furniture objects can be learned from a set of examples [Yu et al. 2011]. Assembling these relationships and other ergonomic factors into a cost function, multiple arrangements can be yielded quickly by simulated annealing using a Metropolis-Hastings state search step. In these methods, manual labours are required in modeling the design rules and providing an initial layout. Fisher et al. [**?**] trains a probabilistic model for indoor scenes from a small number of examples. A variety of indoor scenes can be automatically synthesized from a few of user specified examples. Indoor scenes bring more difficulties for scene analysis because there are always many cluttered objects in different scales, shapes, and functions. A focal-driven analysis and organization framework is presented for heterogeneous collections of indoor scenes [Xu et al. 2014]. They develop a co-analysis algorithm which interleaves frequent patten mining and subspace clustering. The interrelations between objects play important role during furniture arrangement in these systems. However, the 3D scene models takes many efforts to collect for training. In comparison, our system provides an efficient framework to generate 3d model examples for many further applications.

## 3  Overview

The objective of our algorithm is to recover object behaviors from a set of point clouds scanned at different times for an indoor scene. To achieve this goal, objects and correspondences between objects and point clouds must be extracted. More specifically, our problem is to transfer the noisy and incomplete point-level data into object-level
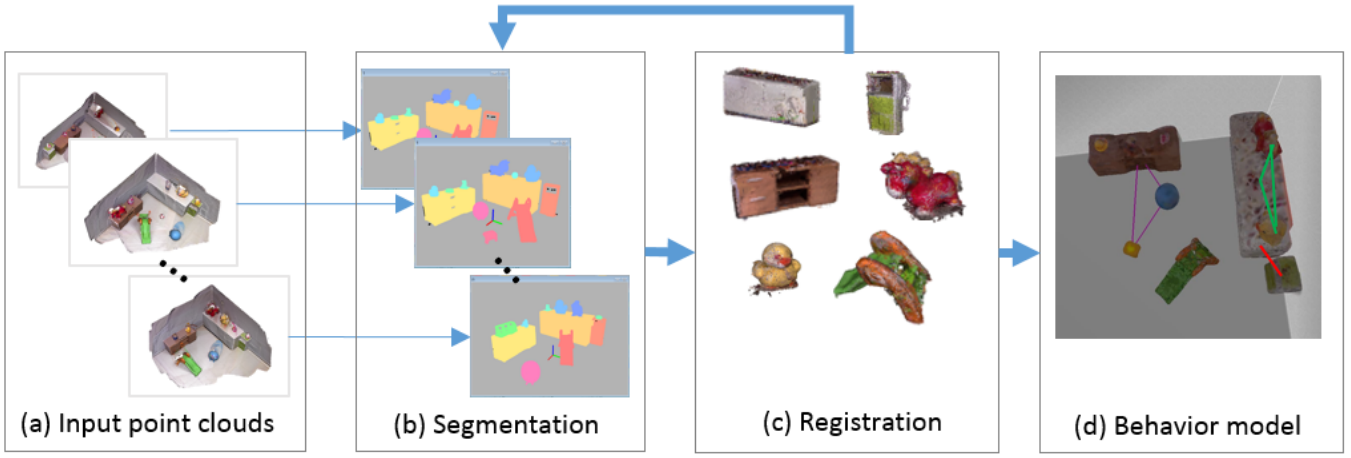
**Figure 2:** *System overview.*

models and correspondences to understand the object behaviors in the scene.

Our system consists of two main steps, *point cloud segmentation* and *object registration*, as Figure 2 shows. The **input** is a set of point clouds scanned using a RGBD camera (a). We scan each scene at different times during a month using a Microsoft Kinect V1. A set of dense point clouds are generated using a real-time fusion system [Nießner et al. 2013]. It is burdensome to scan every detailed structure in a cluttered indoor scene. As a result, each point cloud is incomplete and noisy due to object occlusions.

We first segment each frame simply using region growing (b). There are very likely many wrong boundaries in the generated patches. Then we cluster all the patches from all frames into $k$ clusters using $k$means. [Jia et al. 2015] have demonstrated the power of features based on bounding box in the segmentation of indoor objects. Therefore, we design the descriptor of each patch as the cascades of length, width and height of its bounding box, mean and standard deviation of the distance of each point to its bounding box, percentage of closest points to the faces of its bounding box. The feature dimension is then reduced using PCA. (xuejin: Currently, this step is manually done.)

In each cluster, the patches are registered using a joint registration method [Evangelidis et al. 2014] to produce the object model of this cluster (c). There are inevitably wrong registration due to wrong clustering. Therefore, we project the generated object models back into each frame using the estimated transformation. Resegmentation is performed using the model consistency and neighborhood information in each frame, as described in Sec. 4. By iteratively register resegmented patches and re-segment frames, our system converges to a set of well-registered 3D object models and accurate correspondence between object model and all point clouds. Then we learn the behavior model (d) based the object correspondences from all input point clouds (Sec. 5), then apply the behavior model into many applications (Sec. 6).

# 4 Iteratively Co-Segmentation and Joint Registration

## 4.1 Region Grow

## 4.2 Object Clustering

## 4.3 Joint Registration

## 4.4 Global Consistent Graph Cut

Figure 3 shows the segmentation is progressively refined.

# 5 Behavior Model

Our behavior model in a dynamic scene consists of three parts. One is the behavior of single object, which is represented by the statistics of the object motion (rotations/translations) interpreted from its orientations and positions in frames. The second is the behavior of a group of object, which means the correlations of the statistics of the motions of objects. Their supporting/proximity relationship or arrangements of a group of objects keep consistent. c. The relationship variations caused by motion of the objects in the set of point clouds. In other words, the behavior of the object relationship. For example, the support relationship of vase can change from one desktop to another desktop or so.

## 5.1 Pairwise Object Relations

In a cluttered indoor scene, there are many pairs of objects that usually co-occur in the scene with some typical geometric relations in a long time period. In our system, we consider two types of pairwise relations, which are commonly used to describe object relations in many previous methods [Fisher et al. 2012; Nathan Silberman and Fergus 2012; Xu et al. 2013; Xu et al. 2014].

**Support**   $SR(\mathbf{O}_a, \mathbf{O}_b)$ indicates that $\mathbf{O}_a$ is supported by $\mathbf{O}_b$. We use a Gaussian function to indicate the distribution of $\mathbf{O}_a$ on the supporting surface of $\mathbf{O}_b$; In each frame, we check any pair of objects connected by a surface. Typically, smaller object is supported by the larger one, the upper object is supported by the bottom object. We transform the object been supported $\mathbf{O}_a$ into the local coordinates of its support object $\mathbf{O}_b$ to get its location $\mathbf{x}_a$ and orientation
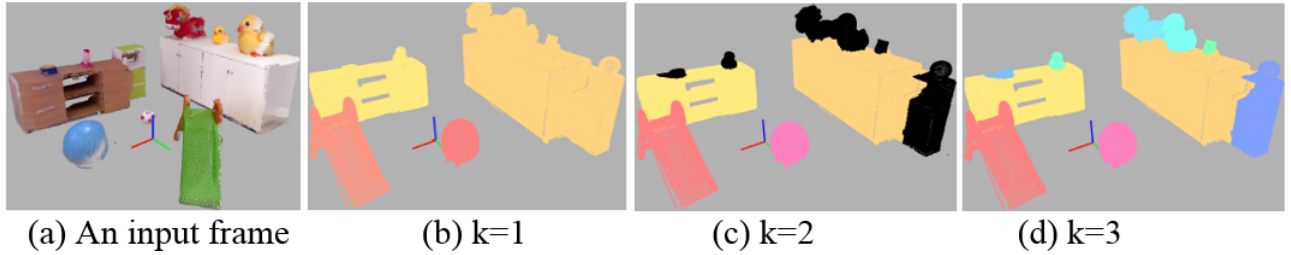
**Figure 3:** *The segmentation of each frame is progressively refined based on registered object models. From left to right: an input point cloud (a), segmentation updates at three iterations.* *(xuejin: show corresponding object model at each iteration.)*

$\theta_a$. According to all the examples of this pair in the input frames, we compute a Gaussian function

$$SR(\mathbf{O}_a, \mathbf{O}_b) = \mathcal{N}_{position}(\mu_x, \Sigma_x)\mathcal{N}_{orientation}(\mu_o, \sigma_o) \qquad (1)$$

Each relation has a confidence or frequency $w$ in all the frames.

**Proximity** $PR(\mathbf{O}_a, \mathbf{O}_b)$ indicates $\mathbf{O}_a$ is close to $\mathbf{O}_b$ in with a Gaussian function $\mathcal{N}(\mu, \sigma)$ including position $x, y$ in the local coordinate system of $\mathbf{O}_a$ and the orientation $\theta_b$ of $\mathbf{O}$ in the local coordinate system of $\mathbf{O}_b$, defined as

$$PR(\mathbf{O}_a, \mathbf{O}_b) = \mathcal{N}_{position}(\mu_x, \Sigma_x)\mathcal{N}_{orientation}(\mu_o, \sigma_o). \qquad (2)$$

If the background floor and wall are taken into account for pairwise relations, the relation between a single object and the background is actually the behavior of this object itself.

## 5.2 Group Behavior

Based on all the extracted pairwise relation, a directed graph $G(V, E)$ is constructed to encode the mutual relations of all the objects in the scene, as shown in Figure 4. Each node is an object model in the scene. Each directed edge represents a pairwise geometric relation between two objects. The information carried in each edge includes:

1. **Relation type**: supporting or proximity. Each edge is directed from a small object to an larger object indicating that the small object is usually associated with a larger object.

2. **Edge weight**: is the frequency of that two objects co-occur in the scene. We define it

$$w(i \rightarrow j) = \frac{c(i \rightarrow j)}{K} \qquad (3)$$

   where $c(i \rightarrow j)$ is the total counts of that two objects co-occur and follow the relation defined by the edge in all the frames. $K$ is the number of all frames.

3. **Relation description**: Associated with each edge, the statistics of the spatial relations defined in Sec. 5.1 is also stored.

A **structural group** $G_s$ is defined as a *complete graph* in which every pair of nodes in connected by a unique edge, to describe the mutual relations between nodes in a subset. From the directed graph $G(V, E)$, a set of structural groups can be extracted. There are many potential structural groups in $G$. The smallest structural group is an edge connecting two nodes. However, we only focus on reliable structural groups for further application, such as layout re-arrangement, unusual event detection, etc. The reliability of a structural group $G_s(V_s, E_s)$ with $k$ objects is defined as defined as

$$\rho_{sg} = \Big( \prod_{e_i \in E_s} w(e_i) \Big)^{\frac{1}{k}} \qquad (4)$$
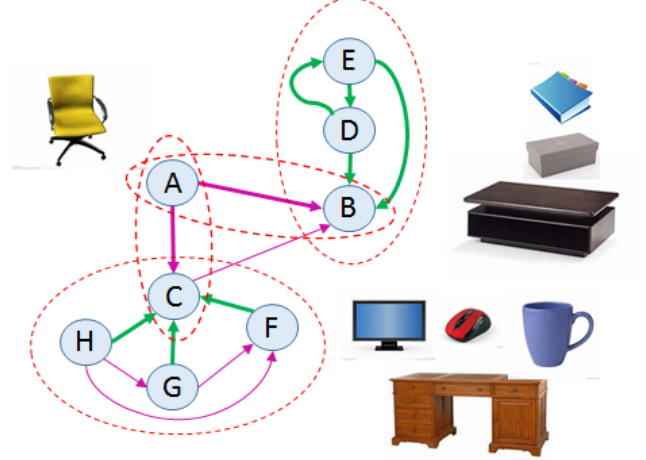


**Figure 4:** *Construction a directed geometric graph from pairwise relations. Green edges indicate supporting relations. Purple edges represent proximity relations. The edge width indicates the frequency of each relation. Each dashed ellipse represents a structural group.* *The pair (B,C) should be very weak because they are far away from other other, and can be removed. Otherwise, (A,B,C) is another structural group.*

# 6 Applications

# 7 Results and Discussions

# 8 Conclusions

In this paper, we present

# References

BOUAZIZ, S., TAGLIASACCHI, A., AND PAULY, M. 2013. Sparse iterative closest point. *Computer Graphics Forum (Symposium on Geometry Processing) 32*, 5, 1–11.

CHANG, W., AND ZWICKER, M. 2011. Global registration of dynamic range scans for articulated model reconstruction. *ACM Transactions on Graphics 30*, 3.

EVANGELIDIS, G., KOUNADES-BASTIAN, D., HORAUD, R., AND E.Z., P. 2014. A generative model for the joint registration of multiple point sets. In *European Conference on Computer Vision (ECCV)*.

FISHER, M., RITCHIE, D., SAVVA, M., FUNKHOUSER, T., AND HANRAHAN, P. 2012. Example-based synthesis of 3D object arrangements. *ACM Trans. Graph. 31*, 6 (Nov.), 135:1–135:11.

FISHER, M., SAVVA, M., LI, Y., HANRAHAN, P., AND NIESSNER, M. 2015. Activity-centric scene synthesis for functional 3d scene modeling. *ACM Transactions on Graphics (TOG) 34*, 6.

HENRY, P., KRAININ, M., HERBST, E., REN, X., AND FOX, D. 2012. RGB-D mapping: Using kinect-style depth cameras for dense 3D modeling of indoor environments. *International Journal of Robotics Research (IJRR) 31*, 5 (April), 647–663.

HU, R., ZHU, C., VAN KAICK, O., LIU, L., SHAMIR, A., AND ZHANG, H. 2015. Interaction context (icon): Towards a geometric functionality descriptor. *ACM Trans. Graph. 34*, 4, 83:1–83:12.

HUANG, Q., WANG, F., AND GUIBAS, L. 2014. Functional map networks for analyzing and exploring large shape collections. *ACM Trans. Graph. 33*, 4 (July), 36:1–36:11.

IZADI, S., KIM, D., HILLIGES, O., MOLYNEAUX, D., NEW-COMBE, R., KOHLI, P., SHOTTON, J., HODGES, S., FREE-MAN, D., DAVISON, A., AND FITZGIBBON, A. 2011. Kinect-fusion: Real-time 3d reconstruction and interaction using a moving depth camera. In *ACM Symposium on User Interface Software and Technology*.

JIA, Z., GALLAGHER, A. C., SAXENA, A., AND CHEN, T. 2015. 3D reasoning from blocks to stability. *IEEE Transactions on Pattern Analysis and Machine Intelligence 37*, 5, 905–918.

KIM, Y. M., MITRA, N. J., YAN, D.-M., AND GUIBAS, L. 2012. Acquiring 3d indoor environments with variability and repetition. *ACM Transactions on Graphics 31*, 6, 138:1–138:11.

KOPPULA, H., ANAND, A., JOACHIMS, T., AND SAXENA, A. 2011. Semantic labeling of 3D point clouds for indoor scenes. In *Conference on Neural Information Processing Systems (NIPS)*.

LIU, Z., TANG, S., XU, W., BU, S., HAN, J., AND ZHOU, K. 2014. Automatic 3D indoor scene updating with rgbd cameras. *Computer Graphics Forum (Pacific Graphics) 33*, 7.

MERRELL, P., SCHKUFZA, E., LI, Z., AGRAWALA, M., AND KOLTUN, V. 2011. Interactive furniture layout using interior design guidelines. *ACM Trans. Graph. (Siggraph'11)*.

NAN, L., XIE, K., AND SHARF, A. 2012. A search-classify approach for cluttered indoor scene understanding. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia 2012) 31*, 6.

NATHAN SILBERMAN, DEREK HOIEM, P. K., AND FERGUS, R. 2012. Indoor segmentation and support inference from rgbd images. In *ECCV*.

NIESSNER, M., ZOLLHÖFER, M., IZADI, S., AND STAMMINGER, M. 2013. Real-time 3d reconstruction at scale using voxel hashing. *ACM Trans. Graph. 32*, 6, 169:1–169:11.

SAVVA, M., CHANG, A. X., HANRAHAN, P., FISHER, M., AND NIESSNER, M. 2014. Scenegrok: Inferring action maps in 3d environments. *ACM Transactions on Graphics (TOG) 33*, 6.

SHAO, T., XU, W., ZHOU, K., WANG, J., LI, D., AND GUO, B. 2012. An interactive approach to semantic modeling of indoor scenes with an rgbd camera. *ACM Trans. Graph.*, 136–136.

SU, H., HUANG, Q., MITRA, N. J., LI, Y., AND GUIBAS, L. 2014. Estimating image depth using shape collections. *ACM Trans. Graph. 33*, 4 (July), 37:1–37:11.

WAND, M., JENKE, P., HUANG, Q., BOKELOH, M., GUIBAS, L., AND SCHILLING, A. 2007. Reconstruction of deforming geometry from time-varying point clouds. In *Proceedings of the Fifth Eurographics Symposium on Geometry Processing*, Eurographics Association, Aire-la-Ville, Switzerland, Switzerland, SGP '07, 49–58.

XIAO, J., AND FURUKAWA, Y. 2012. Reconstructing the world's museums. In *Proceedings of the 12th European Conference on Computer Vision*, ECCV '12.

XU, K., CHEN, K., FU, H., SUN, W.-L., AND HU, S.-M. 2013. Sketch2scene: Sketch-based co-retrieval and co-placement of 3d models. *ACM Transactions on Graphics 32*, 4, 123:1–123:12.

XU, K., MA, R., ZHANG, H., ZHU, C., SHAMIR, A., COHEN-OR, D., AND HUANG, H. 2014. Organizing heterogeneous scene collection through contextual focal points. *ACM Transactions on Graphics, (Proc. of SIGGRAPH 2014) 33*, 4, to appear.

XU, K., HUANG, H., SHI, Y., LI, H., LONG, P., CAICHEN, J., SUN, W., AND CHEN, B. 2015. Autoscanning for coupled scene reconstruction and proactive object analysis. *ACM Transactions on Graphics (Proc. of SIGGRAPH Asia 2015) 34*, 6, to appear.

YAN, F., SHARF, A., LIN, W., HUANG, H., AND CHEN, B. 2014. Proactive 3d scanning of inaccessible parts. *ACM Transactions on Graphics(Proc. of SIGGRAPH 2014) 33*, 4.

YU, L.-F., YEUNG, S. K., TANG, C.-K., TERZOPOULOS, D., CHAN, T. F., AND OSHER, S. 2011. Make it home: automatic optimization of furniture arrangement. *ACM Trans. Graph. 30*, 4, 86.