

# Point Set Joint Annotation for Indoor Scenes

Siyu Hu<sup>1</sup>

<sup>1</sup> USTC



Figure 1: New EG Logo

---

## Abstract

*The ABSTRACT is to be in fully-justified italicized text, between two horizontal lines, in one-column format, below the author and affiliation information. Use the word “Abstract” as the title, in 9-point Times, boldface type, left-aligned to the text, initially capitalized. The abstract is to be in 9-point, single-spaced type. The abstract may be up to 3 inches (7.62 cm) long. Leave one blank line after the abstract, then add the subject categories according to the ACM Classification Index (see <http://www.acm.org/about/class/1998>).*

Categories and Subject Descriptors (according to ACM CCS): I.3.3 [Computer Graphics]: Picture/Image Generation—Line and curve generation

---

## 1. Introduction

In many researches and applications of indoor scenes the data of segmented and even annotated 3D indoor scenes are required as either data base or training data (e.g. [NXS12] [DSS12] [FRS\*12] [CLW\*14] [FSL\*15]).

One way to build such data base is to interactively compose scenes from 3D shape models resulting in scenes with object segmentation and annotation naturally available, or to manually segment and annotate existing scenes. This procedure can be tedious and time consuming, despite the efforts to improve the interaction experience (e.g. [MSL\*11] [XCF\*13]).

Another way is to automatically generate scenes from 3D shape models according to the input RGB or RGB-D images (e.g. [LZW\*15] [CLW\*14]). In such methods, a retrieval procedure is usually needed and inevitably limit the result to a certain set of 3D models despite the actual 3D model in the input images.

We prefer a approach that helps us build such data set directly from the captured data. One of the major gap between the required data set and available scene capturing framework (e.g. [IKH\*11]) is the general object level segmentation. We want to stress that a gen-

eral object level segmentation problem should not be treated as an equivalence of multilabel classification problem since it is not limited to a certain set of objects. For 3D data, [JGSC15] used some simplified physical prior knowledge (i.e. the block based stability) to help achieving the general objectness segmentation, while the work of [XHS\*15] proposes a practical and rather complete framework to close the gap between the required data set and available scene capturing method. One of the observation in [XHS\*15] is that the motion consistency of rigid object can serve as a strong evidence of general objectness. To exploit this fact, they employ a robot to do proactive push and use the movement tracking to verify and iteratively improve their object level segmentation result. Our work presented in this paper is trying to exploit the same observation from a different approach.

We intend to use the motion consistency that is naturally revealed by human activities along the time. Down to this approach, we are facing the choice of scanning scheme. One way is to record the change of the scene along with the human activities, another is to schedule a daily or even a once every half day sweep to only record the result of human activities but avoid the instant of human mo-

tion. The main challenge brought in by the second scheme is that we may not be able to solve the object correspondence by a local search due to the sparse sampling over time, but the very same challenge exists in the first scheme due to the exclusion caused by human bodies not to mention other additional process (e.g. tracking with severe occlusion) needed for human bodies. With the second scanning scheme, our original intention of building 3D scene data set from capturing naturally leads us to the problem of coupled joint registration and co-segmentation.

In this problem, registration and segmentation are entangled in each other. On one hand the segmentation depends on the registration to connect the point clouds into series of rigid movement so that the objectness segmentation can be done based on the motion consistency, on the other hand, the registration depends on the segmentation to break the problem into a series of rigid joint registration instead of a joint registration with non-coherent point drift (A pair of points is close to each other in one point set but their correspondent pair of points in another point set is far from each other, in other words, the point drift of this pair is non-coherent. This happens when this pair of points actually belong to different objects.) To model the problem, we employ a group of Gaussian mixture models and each of these Gaussian mixture models represents a potential object. This model untangle the registration and segmentation in the way that the segmentation can be done by evaluate the probability of points belongs to the Gaussian mixture models and the registration can be done by evaluate rigid registration against each gaussian mixture models.

In summary our work makes following contributions:

Firstly, as far as we know we are the first work that bring up with the problem of point set joint registration and co-segmentation for indoor scenes.

Secondly, we come up with a Gaussian mixture model based formulation to simultaneously model both the registration and co-segmentation problem.

Thirdly, targeting the disadvantages of our formulation, we design a procedure of interaction and provide a practical tool for point set joint annotation based on it.

## 2. Related Work

In this section we explain how our work is related to the previous works.

### 2.1. Point Set Registration with GMM Representation

There are a series of works that uses gaussian mixture model as representation for point set to formulate the registration problem. [MS10] consider the registration of two point sets as a probability density estimation problem. They force the Gaussian mixture model centroids to move coherently as a group to preserve the topological structure of the point sets. Their method is applicable to both rigid registration and non-rigid registration. As we highlighted in section 1, our problem is different from the non-rigid registration considered in [MS10], the point drift could be non-coherent in our problem. [JV11] summarized the works for point set registration using Gaussian mixture models and present a unified framework for the rigid and nonrigid point set registration problem. These works select one of the point set as the “model”.

Unlike these works, [EKBHP14] treats all the point sets as data: they are all realizations of a Gaussian mixture and the registration is cast into a clustering problem. Comparing to these works, our work is most related to [EKBHP14]. Our formulation can be seen as an extension of the formulation of [EKBHP14] to simultaneously handle joint registration and co-segmentation.

### 2.2. Image segmentation and co-segmentation

[RKB04] is an influential work for interactive image segmentation. It uses two Gaussian mixture model, one for foreground and one for background. To initialize these two Gaussian mixture models, [RKB04] let user place a rectangle that contain the foreground. Our design of interaction is similar to [RKB04] in spirit. [TSS16] jointly recover cosegmentation and dense per-pixel correspondence in two images. Our work solve a similar problem for 3D point sets.

### 2.3. Objectness Segmentation

[JGSC15] [JXG17]

## 3. Method Overview

### 3.1. Problem Statement

Given a set of point clouds which record the same group of rigid indoor objects with different layout. We intend to simultaneously partition the point clouds into objects and align the points of same object to recover layouts for corresponding object. Figure ?? shows an example of input point clouds set.

### 3.2. Formulation

To formulate the relation between the unknown object set and the input point clouds. We come up with a generation model as follows:

$$P(v_{mi}) = \sum_{k=1}^{K_n} p_k N(v_{mi} | \phi_{mn}(x_k), \Sigma_k) \quad (1)$$

which means, The observed point clouds are generated by  $N$  object model. Each object model is represented by a gaussian mixture model with  $K_n$  centroids. Our goal is to maximize the probability of the expected complete-data log-likelihood. The object function can be written as:

$$\Theta = \arg \max_Z \sum_Z P(Z|V, \Theta) \ln P(V, Z; \Theta) \quad (2)$$

in which:

$$\Theta = \{ \{p_k, x_k, \Sigma_k\}_{k=1}^{\sum K_n}, \{\phi_{mn}\}_{m=1, n=1}^{MN} \}$$

is the parameters of the generation model.

$p_n$  is the prior probability that the point is generated by the  $n$ -th object.

$p_k$  is the weight of the  $k$ -th Gaussian.

$x_k$  is the center of the  $k$ -th Gaussian.

$\Sigma_k$  is the standard deviation of the  $k$ -th Gaussian.

There are  $\sum K_n$  Gaussian model in total and among them,  $K_n$  Gaussian models belongs to object  $n$ .

$V$  is the  $M$  input point clouds.

$v_{mi}$  is the  $i$ -th point of the  $m$ -th point cloud.  
 $Z$  is a latent variable set defined as:

$$Z = \{z_{ij} | j = 1 \dots M, i = 1 \dots N_j\}$$

among which if  $z_{ij} = k (k = 1 \dots \sum K_n)$  assign the observation of  $\phi_{mn}(v_{mi})$  to the  $k$ -th component of Gaussian mixture model. Such formulation can be seen as an extension of joint registration formulation in [EKBHP14], upon which we add several gaussian mixture model together to express a group of objects. By solving this new problem we simultaneously solve the object co-segmentation of given observation.

### 3.3. Interaction Design for Annotation

Unfortunately, there are several disadvantages for our optimization method. In this subsection we first introduce our interaction design, showing that it is intuitive for users to annotate the point set this way. We then explain how our formulation and optimization can be turned into a practical annotation tool when equipped with our design of interaction.

## 4. Algorithms and Implementation Details

### 4.1. Expectation Conditional Maximization

Assuming the observed point clouds  $\{V_m\}$  are independent and identically distributed, we can then write the (2) as:

$$\epsilon(\Theta | V, Z) = \sum_{m,i,k} \alpha_{mik} (\log p_k + \log P(\phi_{nm}(v_{mi}) | z_{ji} = k; \Theta)) \quad (3)$$

In which the  $\alpha_{mik} = P(z_{mi} = k | v_{mi}; \Theta)$ ,

---

#### Algorithm 1 Joint Registration and Co-segmentation (JRCS)

---

##### Input:

$\{V_m\}$ : Observed point clouds

$\{\alpha_{mik}^0\}$ : Initial posterior probabilities

##### Output:

$\Theta^q$ : Final parameter set

1.  $q \leftarrow 0$
  2. **repeat**
  3. CM-step-a: Use  $\alpha_{mik}^q, x_k^{q-1}$  to estimate  $\{R_{mn}^q\}$  and  $\{t_{mn}^q\}$
  4. CM-step-b: Use  $\alpha_{mik}^q, \{R_{mn}^q\}$  and  $\{t_{mn}^q\}$  to estimate the Gaussian centers  $x_k^q$
  5. CM-step-c: Use  $\alpha_{mik}^q, \{R_{mn}^q\}$  and  $\{t_{mn}^q\}$  to estimate the covariances  $\Sigma_k^q$
  6. CM-step-d: Use  $\alpha_{mik}^q$  to estimate the priors  $p_k^q$
  7. E-step: Use  $\Theta^{q-1}$  to estimate posterior probabilities.  $\alpha_{mik}^q = P(z_{mi} | v_{mi}; \Theta^{q-1})$
  8.  $q \leftarrow q + 1$
  9. **until** Convergence
  10. **return**  $\Theta^q$
- 

### 4.2. Initialization Techniques

A key advantage motivates our formulation is that the soft correspondence can be initialized more flexibly comparing to the typical initialization techniques such as landmark point pairs in registration.

The result of Clustering:

$$P(B_{mj} \in C_n)$$

### Soft Correspondence Initialization

Then the  $\alpha$  is initialized as:

$$\alpha_{ijk} = P(B_{mj} \in C_n)$$

on the condition that:

$$v_{ij} \in B_{mj} \wedge x_k \in O_n$$

## 5. Experiment and Discussion

### References

- [CLW\*14] CHEN K., LAI Y.-K., WU Y.-X., MARTIN R., HU S.-M.: Automatic semantic modeling of indoor scenes from low-quality rgb-d data using contextual information. *ACM Trans. Graph.* 33, 6 (Nov. 2014), 208:1–208:12. URL: <http://doi.acm.org/10.1145/2661229.2661239>, doi:10.1145/2661229.2661239. 1
- [CR00] CHUI H., RANGARAJAN A.: A new algorithm for non-rigid point matching. In *Computer Vision and Pattern Recognition, 2000. Proceedings. IEEE Conference on* (2000), vol. 2, pp. 44–51 vol.2. doi:10.1109/CVPR.2000.854733.
- [DSS12] DEMA M. A., SARI-SARRAF H.: 3d scene generation by learning from examples. In *Multimedia (ISM), 2012 IEEE International Symposium on* (Dec 2012), pp. 58–64. doi:10.1109/ISM.2012.19. 1
- [EKBHP14] EVANGELIDIS G. D., KOUNADES-BASTIAN D., HORAUD R., PSARAKIS E. Z.: *A Generative Model for the Joint Registration of Multiple Point Sets*. Springer International Publishing, Cham, 2014, pp. 109–122. URL: [http://dx.doi.org/10.1007/978-3-319-10584-0\\_8](http://dx.doi.org/10.1007/978-3-319-10584-0_8), doi:10.1007/978-3-319-10584-0\_8. 2, 3
- [FRS\*12] FISHER M., RITCHIE D., SAVVA M., FUNKHOUSER T., HANRAHAN P.: Example-based synthesis of 3d object arrangements. *ACM Trans. Graph.* 31, 6 (Nov. 2012), 135:1–135:11. URL: <http://doi.acm.org/10.1145/2366145.2366154>, doi:10.1145/2366145.2366154. 1
- [FSL\*15] FISHER M., SAVVA M., LI Y., HANRAHAN P., NIESSNER M.: Activity-centric scene synthesis for functional 3d scene modeling. *ACM Trans. Graph.* 34, 6 (Oct. 2015), 179:1–179:13. URL: <http://doi.acm.org/10.1145/2816795.2818057>, doi:10.1145/2816795.2818057. 1
- [IKH\*11] IZADI S., KIM D., HILLIGES O., MOLYNEAUX D., NEWCOMBE R., KOHLI P., SHOTTON J., HODGES S., FREEMAN D., DAVISON A., FITZGIBBON A.: Kinectfusion: Real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology* (New York, NY, USA, 2011), UIST '11, ACM, pp. 559–568. URL: <http://doi.acm.org/10.1145/2047196.2047270>, doi:10.1145/2047196.2047270. 1
- [JGSC15] JIA Z., GALLAGHER A. C., SAXENA A., CHEN T.: 3d reasoning from blocks to stability. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37, 5 (May 2015), 905–918. doi:10.1109/TPAMI.2014.2359435. 1, 2

- [JV11] JIAN B., VEMURI B. C.: Robust point set registration using gaussian mixture models. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33, 8 (Aug 2011), 1633–1645. doi:10.1109/TPAMI.2010.223. 2
- [JXG17] JAIN S., XIONG B., GRAUMAN K.: Pixel objectness. *arXiv preprint arXiv:1701.05349* (2017). 2
- [LZW\*15] LIU Z., ZHANG Y., WU W., LIU K., SUN Z.: Model-driven indoor scenes modeling from a single image. In *Graphics Interface Conference* (2015). 1
- [MS10] MYRONENKO A., SONG X.: Point set registration: Coherent point drift. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32, 12 (Dec 2010), 2262–2275. doi:10.1109/TPAMI.2010.46. 2
- [MSL\*11] MERRELL P., SCHKUFZA E., LI Z., AGRAWALA M., KOLTUN V.: Interactive furniture layout using interior design guidelines. *ACM Trans. Graph.* 30, 4 (July 2011), 87:1–87:10. URL: <http://doi.acm.org/10.1145/2010324.1964982>, doi:10.1145/2010324.1964982. 1
- [NXS12] NAN L., XIE K., SHARF A.: A search-classify approach for cluttered indoor scene understanding. *ACM Trans. Graph.* 31, 6 (Nov. 2012), 137:1–137:10. URL: <http://doi.acm.org/10.1145/2366145.2366156>, doi:10.1145/2366145.2366156. 1
- [RKB04] ROTHER C., KOLMOGOROV V., BLAKE A.: Grabcut -interactive foreground extraction using iterated graph cuts. URL: <https://www.microsoft.com/en-us/research/publication/grabcut-interactive-foreground-extraction-using-iterated-graph-cuts/>. 2
- [TSS16] TANIAI T., SINHA S. N., SATO Y.: Joint recovery of dense correspondence and cosegmentation in two images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2016). 2
- [XCF\*13] XU K., CHEN K., FU H., SUN W.-L., HU S.-M.: Sketch2scene: Sketch-based co-retrieval and co-placement of 3d models. *ACM Trans. Graph.* 32, 4 (July 2013), 123:1–123:15. URL: <http://doi.acm.org/10.1145/2461912.2461968>, doi:10.1145/2461912.2461968. 1
- [XHS\*15] XU K., HUANG H., SHI Y., LI H., LONG P., CAICHEN J., SUN W., CHEN B.: Autoscanning for coupled scene reconstruction and proactive object analysis. *ACM Trans. Graph.* 34, 6 (Oct. 2015), 177:1–177:14. URL: <http://doi.acm.org/10.1145/2816795.2818075>, doi:10.1145/2816795.2818075. 1