

Preparation of Papers for IEEE Sponsored Conferences & Symposia*

Akshaya Thippur¹, Rares Ambrus¹, Kaushik Desai¹, Adria Gallart¹, Malepati Sai Akhil², Gaurav Agrawal², Mayank Jha², Janardhan HR², Nishan Shetty², Prasad NR², John Folkesson¹ and Patric Jensfelt¹

Abstract—

Humans subconsciously exploit various strong correlations amidst different object instances, classes and between different object classes and scene types when analysing indoor environments. Correlations in naive, logical object co-occurrences have been exploited along with the extraction of vision based object-intrinsic descriptive features in previous research. In this paper, we present several alternative learning techniques to model and make estimates of scenes based on a variety of spatial relations - geometric extrinsic features with different amounts of discretization, which capture *how* the objects co-occur; and compare their efficacy in the context of object classification in real-world table-top scenes. We investigate the possibilities of using such techniques to refine the results from a traditional vision-perception system. We also contribute a unique, long-term periodic, large 3D dataset of 20 office table-top scenes, manually annotated with 18 object classes. Apart from our current comparison, we foresee that the dataset will be useful for applications such as generalized learning of spatial models, learning object data structuring based on semantic hierarchy, learning best suited semantic abstractions and grammar for long term autonomy, ground truth for vision-perception systems etc.

I. INTRODUCTION

Objects pervade human environments. If robots are to perform useful service tasks for humans it is crucial that they are able to locate and identify a wide variety of objects in everyday environments. State-of-the-art object recognition/classification typically relies on the extraction features to be matched against models built through machine learning techniques. As - the number of objects a given system is trained to recognise - increases, the uncertainty of individual recognition results tends to increase as greater number of objects increases the chance of overlapping features existence. The reliability of such recognisers is also affected when used on real robots in everyday environments, as objects may be partially occluded by scene clutter or only visible from certain angles, both potentially reducing the visibility of features for their trained models. In this paper we argue that the performance of a robot on an object recognition task can be increased by the addition of *contextual knowledge* about the scene the objects are found in. In particular we

demonstrate how models of the *spatial configuration* of objects, learnt over prior observations of real scenes, can allow a robot to recognise the objects in unseen scenes more reliably.

Our work is performed in the context of developing a mobile service robot for long-term autonomy in indoor human environments, from offices to hospitals. The ability for a robot to run for weeks or months in its task environment opens up a new range of possibilities in terms of capabilities. In particular, any task the robot performs will be done in an environment it may have visited many times before, and we wish to find ways to capture the contextual knowledge gained from previous visits in a way that enables subsequent behaviour to be improved. The use of context to improve object recognition is just one example of this new robotics paradigm. In this paper we focus on the task of *table-top scene understanding*, and more specifically what objects are present on a table-top. Whilst the objects present on a single table may change in position, their overall arrangement has some regularity over time as influenced by the use to which the table is put. For example, if this table is used for computing, then a (relatively static) monitor will be present, with a keyboard in front of it and mouse to one side. A drink, or paper and a pen, may be within an arms length of the keyboard, as may headphones or a cellphone. This arrangement may vary across different tables in the same building, but the overall pattern of arrangements will contain some structure. It is this structure we aim to exploit in order to improve the recognition of table-top objects, e.g. knowing that the object to the right of a keyboard is more likely to be a mouse than a cellphone.

As the absolute positions of objects on a table (or their relative positions with respect to some fixed part of the table) is unlikely to generalise across a range of different tables, we are investigating *relational* models of space, i.e. ways of encoding the position of a target object relative to the position of one or more landmark objects. Using a novel data set of table-top scenes (described in Section IV), in this paper we explore the performance of a variety of representations for relative object position, plus inference techniques for operating on them, on the task of table-top scene understanding (Section ??). In particular we investigate representations that use varying forms of spatial relations, from geometric ones such as distances and angles to more qualitative spatial relations such as *Left* and *Behind* as a means for capturing observations of object configurations over time. The contributions this paper makes are: (1) A novel comparison between mechanisms for representing, learning and inference

*This work was supported by STRANDS

¹KTH Royal Institute of Technology
albert.author@papercept.net

²MSRIT Bangalore b.d.researcher@ieee.org

on object spatial configurations using spatial relations. (2) An evaluation of the use of these mechanisms for augmenting a robot's vision based *perceptual system* (PS). (3) A new large 3D annotated table-top benchmark dataset.

Points discussed:

- 1) **Make this more "Human Oriented"**
- 2) Overview of the problem - why?
 - a) Not only desktops
 - b) How do we treat object classification
- 3) Generalize and Transfer Knowledge
- 4) How are objects correlated, search becomes easier
- 5) Understanding a model for inter object influence for inference
- 6) No benchmarks yet or datasets having 3D spatial information
 - time
 - complete scenes
 - different people
 - different types of people
- 7) need to structure data – metric is tedious, because of large amounts of data
- 8) Contributions:
 - A big 3D data set, annotated – benchmark for results, folding data, classification.
 - suggest SR and QSR
 - something that could augment perception systems

II. RELATED WORK

Spatial relations have been used previously to provide contextual information to vision-related work. [1] used a hierarchy of spatial relations alongside descriptive features to support multiple object detections in a single image. Spatial relations and contextual information are commonly used in activity recognition from video streams. For example, [2] demonstrate the learning of activity phases in airport videos using spatial relations between tracked objects, and [3] use spatial relations to monitor objects and activities in videos of a constrained workflow environment. Recent work has used object co-occurrence to provide context in visual tasks. Examples in 2D include object co-occurrence statistics in class-based image segmentation [4]; and the use of object presence to provide context in activity recognition [5]. However, all this previous work is restricted to 2D images, whereas our approaches work with spatial context in 3D (RGB-D) data. Authors have also worked with spatial context in 3D, including parsing a 2D image of a 3D scene into a simulated 3D field before extracting geometric and contextual features between the objects [6]. Our approaches to encoding 3D spatial context could be applied in these cases, and we use richer, structured models of object relations.

Apart from using the statistics of co-occurrence, a lot of information can be exploited from *how* the objects co-occur in the scene, i.e. the extrinsic, geometric spatial relations between the objects. Recent work in 3D semantic labelling has used such geometric information along with descriptive intrinsic appearance features [7]. They achieve a

high classification accuracy for a large set of object-classes belonging to home and office environments. Scene similarity measurement and classification based on contextual information is conducted by [8]. They also use spatial information for context-based object search using Graph Kernel Methods. The method is further developed to provide synthetic scene examples using spatial relations [9]. In [10] spatial relations between smaller objects, furniture and locations is used for pruning in object search problems in human environments. In [11] a technique is developed for automatic annotation of 3D objects. It uses intrinsic appearance features and geometric features and is employed to build an object and scene classifier using conditional random fields. In [12] the authors utilise both geometric single object features and pair-wise spatial relations between objects to develop an empirical base for scene understanding. Recent studies [13], [12] compute statistics of spatial relations of objects and use it for conditional object recognition for service robotics. Whilst our techniques are comparable to those in the literature, our contribution comes from the explicit comparison of different representations of spatial context (metric vs qualitative) on a novel, long-term learning task. Additionally our qualitative approach relies on relationships which could be provided through other mechanisms than unsupervised machine learning (e.g. through a human tutor describing a spatial scene), and in this way bootstrap the system using expert knowledge.

Our work is evaluated on a new 3D long-term dataset. Other datasets exist: The *B3DO dataset* [14] which contains many single-snapshot instances of indoor human environments having a variety in viewpoints, object-classes, scene-classes and instances. This dataset is in the form of RGB and depth image pairs with manual 2D annotations of object classes, capturing many unique scenes with the sole aim of finding more realistic scenes which are difficult for PSs to perform scene classification. *NYU Depth V1-2* [15] datasets contain different instance examples of object-classes and scene-classes. Each image instance is a combo of synchronous RGB and D images of a different scene-class with semantic annotation provided to every pixel. This dataset is aimed at helping PSs with automatic semantic segmentation and scene classification. The *3D IKEA database* [16] has been collected using robotic maneuvering in different scene-class instances. The aim is to test scene-classification algorithms based on large furniture level objects. The *WRGBD dataset* [17] is aimed to support object classification methods and contains many instances of isolated objects in .pcd format. Annotation is done by assigning every pixel a correct semantic label. None of these datasets contain periodically collected data or easily usable spatial annotations of objects which are key for long-term autonomous scene-learning, based on spatial relations.

The required regularity in instances and time was the main motivation for the construction of this dataset, as currently available datasets either are of individual objects or singular instances of entire rooms

- 1) More on dataset papers

III. MOTIVATING SCENARIO

We are investigating systems that operate for long periods of time in environments populated by humans. As a motivating scenario we will look at security guard in an office building. The robot patrols the working environment and should learn models of what the environment normally looks like and what variations there are. In an implementation of such a system the robot tell when something differs from the ordinary too much and then raise an alarm. Initially we will focus on desktop scenes. We are interested in models for individual desks and as well as general models of desks. Our working hypotheses is that there is some general rules for how desks are organized that we want to be able to extract and later exploit when building the models. We expect that object will change place, within certain semantic bounds. Some objects could be missing at times (coffee mug) which is normal, but some other objects are rarely moved (monitor).

Another aim is to be able to transfer knowledge from one environment to the next. This would allow a robot that just entered a new environment to be functional from the start. Concretely this would correspond to having a reasonable prior which can then be adapted when new observations are available. What information is general and what is environment specific? How do we represent the knowledge in a way that caters for the knowledge transfer, the ability to learn from few samples and adapting existing models? These are some further examples of questions that we want to study.

To study these questions we need data to learn from. The data need to capture both variations across different desks but also over time. None of the datasets available (see Section /refsec:Related Work) meet these requirements which is the motivation for the work behind the dataset that we present in this paper.

Points Discussed:

- 1) This section should stress why it is important to gather to such datasets
- 2) what type of data do we really need? - hint it
- 3) reference to related work section to show that this dataset is unique.

IV. DATASET

Most of the human outdoor and especially indoor environments, are characterized by a supporting surface on which all relevant objects are placed e.g. Land - buildings, living room floor - furniture, dining table - cutlery and dishes etc. It is hence convenient to perceive a structure for the object arrangement in the human environments to be as if on a "nested table-top system". Table-tops are the most commonly present large objects in any floor map. They provide for a favourable prototypical example for analysis of object organization structure. Moreover, it is highly likely that objects on a "table" have more semantic and organizational correlation amidst them, than when compared to those objects not on that same "table" e.g. *Pen* has more semantic/functional correlation with a *Laptop* (on the table) than a *Couch* (not

on that table, here the "table" references to the floor of the room).

With the aim of exploring possibilities to understand, learn and model these organizational formalities amongst objects in human indoor environments, a first, pertinent dataset called *3D Table-Tops Dataset for Long Term Autonomous Learning* (3D-TOTAL), has been composed by periodically capturing observations of entire table-tops of a fixed set of researchers at a computer science research facility. A singular observation of a table-top of one person at a single instance in time, captured in image format, is termed as a *scene*. These scenes have been captured with intervals of few hours, over many days and various instances of object classes. The dataset therefore captures the individual and group variation in object position and pose due to humans and their interaction with the environment. The required regularity in instances and time was the main motivation for the construction of this dataset, as currently available datasets either are of individual objects or singular instances of entire rooms.

It is required for such a long-term autonomous learning to have entire views of the scenes to model the interrelations amidst the member objects. Apart from being variant over object instances across different scenes, the data needed to be periodic over time at a scale of couple of hours so as to capture the individual and group variations in position and pose when there has been regular/irregular human interaction involved. The currently available datasets either are of individual objects or single instances of entire rooms. The required regularity in instances and time was the main motivation for the construction of this dataset (Section II).

A. Dataset Design and Concept

The target research groups to benefit from this dataset are of the kind that develops artificial intelligence for autonomous robots augmenting human activities (especially if they are fatiguing) within indoor human environments. Hence, it becomes essential for robotic learning to pay attention to variances in human environments wrt. scales of time, space and instances in the same space. The dataset has been composed by capturing and manually annotating 3D images of office type table-tops, for a fixed set of people, at fixed times of the day and for a span of many days.

Observing the table-tops of the same set of people at different times of the day gives insight about the daily interactions a human has with his table and over many days gives an understanding of the gradual variances in their table-top setups. If the data is observed for an entire week, including weekends, features in the table-top configurations, that can be used for estimation of the type of the day of the week, can be extracted (e.g. Weekdays, Fridays, Weekends). Table-top models can also be learnt for all the people put together – which gives a gross functional representation of a typical table-top in general for research employees in office environments – or for individual people which helps to gain functional representations of office table-tops for individuals, seniority, gender and so on (Figure 1). Finally, when trying to



Fig. 1: Each column shows a different person's table in top view. The first row displays the morning scenes and the second row displays the evening scenes. The first two columns contain scenes from the same day, whereas the third column shows scenes 12 days apart in time. Notice in Column 1: the slight changes in position of the keyboard, mouse, papers and pen; Column 2: the relatively big changes in position of laptop, mouse, papers, pen, keyboard, lamp. When objects in columns 1,2,3 are compared there is a certain generality in structure (keyboards are always in front of monitors), but also a specificity for each person (occurrence of headphones, position of mouse wrt. keyboard). The clutter in column 3 is significantly more in row 2 than row 1 implying the time since the last time the person tidied-up his/her table. Cluttered tables in columns 2,3 belong to graduate students and a more orderly table in column 1 to a senior researcher, lightly suggesting a characteristic of type of researcher.

find general models over any of these types of data partition, the model learns to be able to generalize over different instances of a fixed set of object classes. In summary: When the dataset is partitioned in different ways with respect to time, people or instances, it richly yields knowledge and hence representations of table-tops in office environments.

As explained in Section III our research intentions are to provide intelligence to an long-term operating, autonomous, human activity augmenting robot for indoor human environments. Thusly, the concept of composing the 3D-TOTAL dataset follows naturally from this research motivation.

B. Dataset Realization

In 3D-TOTAL, 3D images of scenes were captured regularly at 3 times a day for 19 days for 20 people. Each scene has been manually annotated to obtain information about common objects-of-interest generally and regularly present, when observed across the many scenes.

The data was collected using the *SCENECT* software [18] and an *Asus Xtion Pro* RGB-D camera. There is one 3D colour point cloud per scene (.pcd format). Every scene is a reconstructed version of the raw data stream obtained by manual detailed scanning of a table-top with real time visual feedback using SCENECT. The software has in-built – real-time sampling, registration and de-noising algorithms to output the final high resolution point cloud.

The scenes were recorded as periodically as possible and at three fixed time instances of the day: *Morning* (09:00 hrs), *Afternoon* (13:00 hrs) and *Evening* (18:00 hrs). Scenes contain tables of 20 different people collected over 19 days including weekends. A *Scene_ID* is attached to each scene to indicate who the table belongs to and the date and time of the recording. These Scene_IDs help in partitioning the

dataset with respect to time of the day {Morning, Afternoon, Evening}, person {Anna, Bob, Carl, ...}, or day {2013-11-01, 2013-11-06, 2013-11-13, ...}.

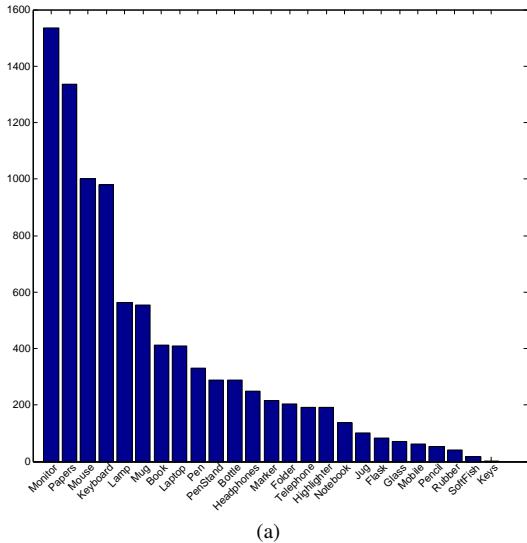
A 3D annotation tool was developed for manually segmenting out objects-of-interest from the point clouds. On average, 12 different objects were labelled, including repeating instances of the same object class, per scene depending upon feasibility and occurrence. The objects belong to the following super set - {Mouse, Keyboard, Monitor, Laptop, Cellphone, Keys, Headphones, Telephone, Pencil, Eraser, Notebook, Papers, Book, Pen, Highlighter, Marker, Folder, Pen-Stand, Lamp, Mug, Flask, Glass, Jug, Bottle}. The information about every scene and object is available in XML and JSON formats. Each scene has a nested list of object data containing {Position, Orientation, Size, Date and Time of recording, Person ID, Point Indices of the point cloud that have been labelled as belonging to the Object}. This manual annotation provides the required ground truth data for long term autonomous learning.

C. Dataset Summary

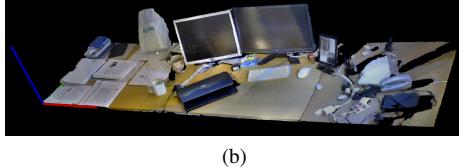
This section gives a brief summary of the dataset to serve as a quick reference.

3D-TOTAL has:

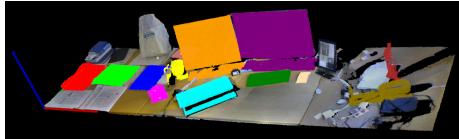
- scenes collected from 20 unique tables, 3 times a day for 19 days and hence in total, approximately 1140 scenes.
- each scene manually annotated with 18 possible object classes and an average of 12 object instances, from these classes, annotated per scene.
- has annotation stored in XML and JSON formats containing scene instance and its object instances' specifications.



(a)



(b)



(c)

Fig. 2: (a) Objects annotated in 3D Long-Term Dataset, sorted in descending order of count of occurrences. X-axis=Object Name, Y-axis=Occurrence Count. (b) Screenshot of one table scene, along with it's annotations in (c).

- occurrences of object instances as depicted in Figure 2a and annotations as exemplified in Figure 2b,2c

V. ANALYSIS

Points Discussed:

- 1) **figure** Scatter plot for variation within person ID
- 2) **figure** Scatter plot for variation wrt object type
- 3) **figure** Scatter plot for variation considering different landmark - trajectors.
- 4) **figure** Scatter plot of 2D footprints of objects
- 5) General conclusions about above figures.
- 6)
- 7)

VI. CONCLUSIONS

VII. FUTURE WORK

Points Discussed:

- 1) Room level data set
- 2) Which QSR for which purpose?

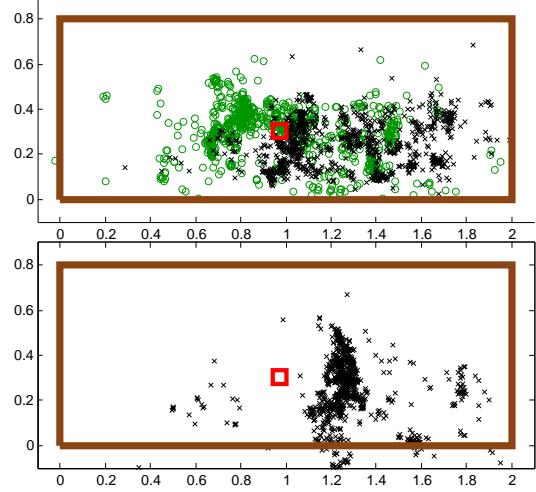


Fig. 3: This figure investigates the correlation in position of keyboard and mouse. Top: The green circles shows the position of the centroid of each keyboard that exist in a scene where there is at least one mouse. The black crosses show the position of all mice in scenes with at least one keyboard. The red square shows the mean position of all keyboards. Bottom: By shifting the position of keyboard and mouse such that the keyboards are at the mean keyboard position we can illustrate the position of the mouse relative to the keyboard. As expected most mice are to the right of the keyboard.

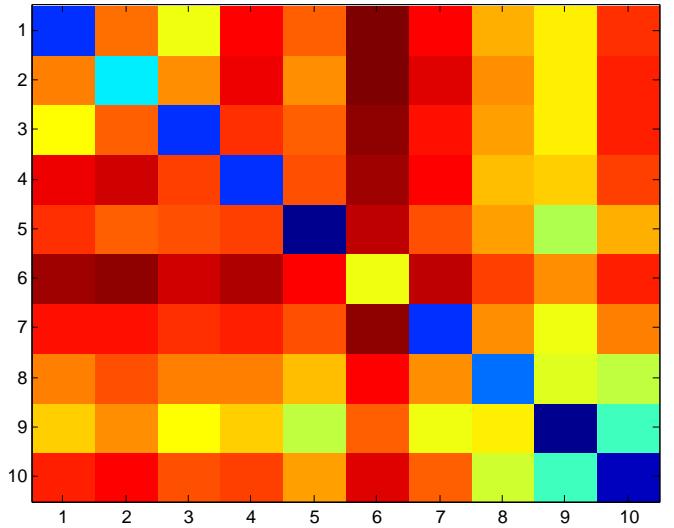


Fig. 4: The figure illustrates the entropy in the distribution of relative positions between a landmark (row) and trajector (column). We show this for a subset of objects from the database. Only objects that occur frequently enough are displayed. The objects in are (in order): 1:keyboard, 2:monitor, 3:mouse, 4:mug, 5:laptop, 6:papers, 7:book, 8:bottle, 9:jug, 10:notebook. Dark red indicates high entropy, i.e., closer to a uniform distribution and dark blue corresponds to low entropy, i.e., a peakier distribution.

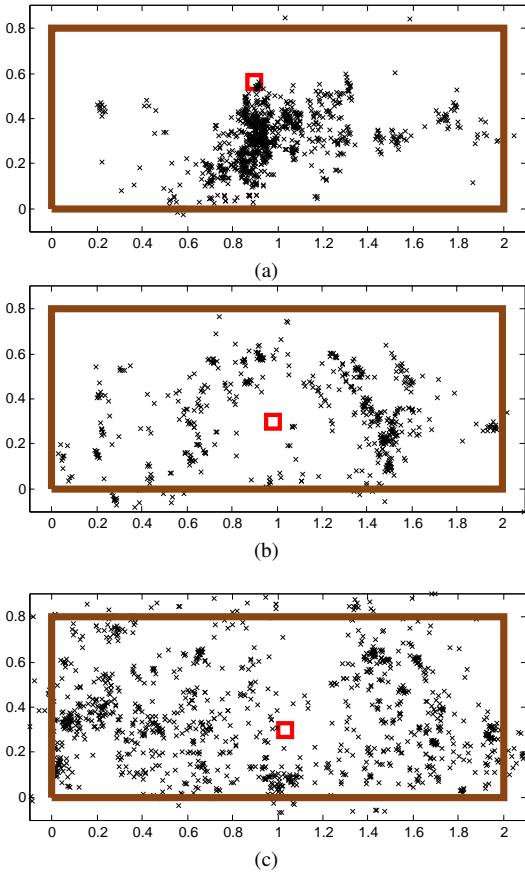


Fig. 5: The figures show the relative position of a) keyboard w.r.t. monitor, b) mug w.r.t. keyboard and c) papers w.r.t. keyboard. Qualitatively keyboards are mostly in front of monitors, mugs are around the keyboard and the position of papers is mostly independent on the position of the keyboard.

VIII. ACKNOWLEDGEMENTS

CVAP-KTH, Accel Partners, STRANDS project

REFERENCES

- [1] M. J. Choi, J. Lim, A. Torralba, and A. Willsky, “Exploiting hierarchical context on a large database of object categories,” in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pp. 129–136, June 2010.
- [2] K. S. R. Dubba, A. G. Cohn, and D. C. Hogg, “Event model learning from complex videos using ilp,” in *Proc. ECAI*, vol. 215 of *Frontiers in Artificial Intelligence and Applications*, pp. 93–98, IOS Press, 2010.
- [3] A. Behera, A. G. Cohn, and D. C. Hogg, “Workflow activity monitoring using dynamics of pair-wise qualitative spatial relations,” in *Advances in Multimedia Modeling*, pp. 196–209, Springer, 2012.
- [4] L. Ladicky, C. Russell, P. Kohli, and P. Torr, “Inference methods for crfs with co-occurrence statistics.” *International Journal of Computer Vision*, vol. 103, no. 2, pp. 213–225, 2013.
- [5] L.-J. Li, H. Su, Y. Lim, and L. Fei-Fei, “Objects as attributes for scene classification,” in *Trends and Topics in Computer Vision* (K. Kutulakos, ed.), vol. 6553 of *Lecture Notes in Computer Science*, pp. 57–69, Springer Berlin Heidelberg, 2012.
- [6] J. Xiao, B. C. Russell, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, “Basic level scene understanding: From labels to structure and beyond,” in *SIGGRAPH Asia 2012 Technical Briefs*, SA ’12, (New York, NY, USA), pp. 36:1–36:4, ACM, 2012.
- [7] H. S. Koppula, A. Anand, T. Joachims, and A. Saxena, “Semantic labeling of 3d point clouds for indoor scenes,” in *Advances in Neural Information Processing Systems*, pp. 244–252, 2011.
- [8] M. Fisher, M. Savva, and P. Hanrahan, “Characterizing structural relationships in scenes using graph kernels,” *ACM Trans. Graph.*, vol. 30, pp. 34:1–34:12, July 2011.
- [9] M. Fisher, D. Ritchie, M. Savva, T. Funkhouser, and P. Hanrahan, “Example-based synthesis of 3d object arrangements,” *ACM Trans. Graph.*, vol. 31, pp. 135:1–135:11, Nov. 2012.
- [10] A. Aydemir, K. Sjoo, J. Folkesson, A. Pronobis, and P. Jensfelt, “Search in the real world: Active visual object search based on spatial relations,” in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pp. 2818–2824, May 2011.
- [11] D. Lin, S. Fidler, and R. Urtasun, “Holistic scene understanding for 3d object detection with rgbd cameras,” *ICCV, December*, 2013.
- [12] A. Kasper, R. Jakel, and R. Dillmann, “Using spatial relations of objects in real world scenes for scene structuring and scene understanding,” in *ICAR 2011: Proceedings of the 15th International Conference on Advanced Robotics*, 2011.
- [13] T. Southey and J. J. Little, “Learning qualitative spatial relations for object classification,” in *IROS 2007 Workshop: From Sensors to Human Spatial Concepts*, 2007.
- [14] A. Janoch, S. Karayev, Y. Jia, J. Barron, M. Fritz, K. Saenko, and T. Darrell, “A category-level 3-d object dataset: Putting the kinect to work,” in *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on*, pp. 1168–1174, Nov 2011.
- [15] P. K. Nathan Silberman, Derek Hoiem, and R. Fergus, “Indoor segmentation and support inference from rgbd images,” in *ECCV*, 2012.
- [16] A. Swadzba and S. Wachsmuth, “A detailed analysis of a new 3d spatial feature vector for indoor scene classification,” *Robotics and Autonomous Systems*, 2012.
- [17] K. Lai, L. Bo, X. Ren, and D. Fox, “A large-scale hierarchical multi-view rgbd object dataset,” in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pp. 1817–1824, May 2011.
- [18] O. Buerkler, “Scenect – faro technologies : <http://www.faro.com/scenect>,” 2012.