# Geometry-Based Region Proposals for Real-Time Robot Detection of Tabletop Objects

**Alexander Broad[1,2] and Brenna Argall[1,2]**

## Abstract

We present a novel object detection pipeline for localization and recognition in three dimensional environments. Our approach makes use of an RGB-D sensor and combines state-of-the-art techniques from the robotics and computer vision communities to create a robust, real-time detection system. We focus specifically on solving the object detection problem for tabletop scenes, a common environment for assistive manipulators. Our detection pipeline locates objects in a point cloud representation of the scene. These clusters are subsequently used to compute a bounding box around each object in the RGB space. Each defined patch is then fed into a Convolutional Neural Network (CNN) for object recognition. We also demonstrate that our region proposal method can be used to develop novel datasets that are both large and diverse enough to train deep learning models, and easy enough to collect that end-users can develop their own datasets. Lastly, we validate the resulting system through an extensive analysis of the accuracy and run-time of the full pipeline.

## 1 Introduction

As the field of robotics advances, and personal robots that assist users in their home and work environments become more prevalent, it will be necessary to extend a robot's autonomy to include more advanced cognitive reasoning and improve abilities in highly dynamic environments. To do so, the robot will need to have knowledge of many of the same physical attributes of the world that a human does. One such important aspect is the ability to recognize and localize objects in common environments. Through this knowledge, a robot can make informed decisions in achieving tasks like intelligently searching for an object in a novel environment, cleaning a room or retrieving an object for a human partner.

The problem of object detection is not unique to robotics. In computer vision it is used to solve problems like automatic caption generation as demonstrated by Karpathy and Fei-Fei (2015) and automatic tagging of shared social pictures as described in Schroff et al. (2015). However, it is often the case that techniques used in the two communities are distinct from one another. One reason for this is that the desired and available sensor information is frequently different — in computer vision, systems are usually limited to the RGB space while solving problems in robotics generally requires depth as an additional, or primary, modality. The requirement of depth information often necessitates an additional sensor (with a few notable exceptions such as Saxena et al. (2009) and Mur-Artal et al. (2015)), which in turn requires a potentially difficult calibration and sensor fusion problem. For this reason, we frequently see methodologies in the two communities that parallel each other in purpose, such as object recognition, but are divergent in technique. However, due to the rise of RGB-D cameras like the Microsoft Kinect (Zhang (2012)), robotics researchers have access to sensors

that provide both color and depth information in a single device. These cameras can be easily calibrated (up to a level of tolerance) and aligned through a single transformation defined by the static configuration of the two integrated sensors.

In this paper, we propose a novel method for object detection that makes use of both the depth and color modalities of RGB-D sensors to recognize and localize objects in real-time. We focus specifically on tabletop environments as many domestic manipulation tasks take place in this type of configuration. Our approach differs from previous proposed techniques in that it solves the localization and recognition tasks independently — the former through an exploitation of the geometry of the scene, and the latter with state-of-the-art deep learning methods. To the best of our knowledge, this work is the first to combine these ideas. Our method achieves high accuracy in both the position and categorization of the detected objects by using a confluence of ideas from the computer vision and robotics communities.

We begin by discussing related work in Section 2 and then present our approach in detail in Section 3. We also discuss how our region proposal method can be used for data acquisition and developing novel datasets in Section 4. Finally, we describe an experimental validation of our system in Section 5 and the results in Section 6. We then discuss the

[1]Northwestern University, Chicago, IL [2]Rehabilitation Institute of Chicago, Chicago, IL

**Corresponding author:**
Alexander Broad, Department of Electrical Engineering and Computer Science, Northwestern University 633 Clark St, Evanston, IL 60208
Email: alex.broad@u.northwestern.edu

success of our approach and future directions in Section 7 and conclude in Section 8.

## 2 Related Work

From a computational perspective, object detection is a two part problem: (1) Where is the object? and (2) What is the object? The long-standing baseline approach in computer vision is known as a sliding window. In this technique, each patch of size $(m, n)$, from an image of size $(M, N)$, is fed through an object recognition model. The concept is that, while this approach may be inefficient, it maximizes recall by ensuring not to skip any possible object locations. To account for scale the same process is repeated over an image pyramid.

More recently, as real-time and interactive systems have become more popular (Chen and Yuille (2005)) there has been an increased focus on the efficiency of object detection systems. The most common way to improve the run-time of the system is to intelligently reduce the number of candidate regions that are run through the object recognition model. New approaches focus on novel techniques and methods for producing *region proposals*. For example, Girshick (2015) uses multi-layer segmentation to produce region proposals at different positions and scales, Szegedy et al. (2013) train a neural network to predict segmentation masks and Zitnick and Dollár (2014) use edge detections. Hosang et al. (2014) provide a comprehensive comparison of region proposal techniques in the computer vision community. By reducing the number of candidate regions, one is able to perform a more efficient search through position, scale and orientation.

These techniques greatly reduce the computational burden when compared to an exhaustive sliding window approach, however they often still require expensive systems and GPUs to train and run. For example, Fast R-CNN, as proposed by Girshick (2015) has proven very successful, yet when using this approach on a $640\times480$ image with a Core i7 laptop with a mid-tier GPU (nVidia GeForce 860M), the full pipeline takes about 0.75 seconds per image. One reason for this is that the segmentation algorithm produces between 1k and 10k proposals per image depending on the *quality mode* parameter. To increase the speed of this system, Ren et al. (2015) extend Girshick's method to a model called Faster R-CNN, which uses a separate neural network to produce object proposals, decreasing the run-time to about 0.2 seconds per image. Of note, these approaches necessitate significantly more data as the training process requires labeled bounding boxes and an extra background class to reject false positives. Additionally, localizing an object in the 2D plane does not fully solve the problem for robotics applications where localizing the object in three dimensions is equally as important as correctly recognizing the object. Lastly, there is also evidence from Chen et al. (2015) that image-based segmentation approaches are not equally effective on all datasets.

There also is related work from the robotics community in 3D object detection. Early approaches are similar to pre-deep learning methods in the vision community; namely, they focus on developing hand-crafted features in the point cloud space. Examples include local features such as the histogram-based Fast Point Feature Histogram of Rusu et al. (2009a) and the Signature of Histogram of Orientations

of Tombari et al. (2010), as well as global features such as the Viewpoint Feature Histogram of Rusu et al. (2010) which also takes the viewpoint into account. Tang et al. (2012) describe a segmentation approach similar to our own, however the recognition is again done in the point cloud space by comparing features to learned object models. The approaches mentioned here work relatively well, however they rely on hand-crafted features and no single approach has found the type of success or wide-spread adoption of Convolutional Neural Networks in the image space (LeCun et al. (2015)).

More closely related to our own work, researchers have begun to look at other methods for combining the RGB and depth modalities when solving the object detection problem. Song and Xiao (2014) describe a three dimensional version of the sliding window approach in which they fit 3D regions to learned CAD-based object models. Dahan et al. (2012) describe a method for computing region proposals by segmenting the input image using information from both the color and depth channels. Couprie et al. (2014) and Gupta et al. (2014) similarly describe methods for including depth during segmentation and they also augment the standard vision approach by including the depth map as another input channel in training the CNN model. The main difference between these works and our own is that we explicitly generate region proposals in the point-cloud space based on geometric constraints, which produces significantly fewer candidate regions. Pillai and Leonard (2015) present a robotic recognition system that incorporates multi-view object proposals and efficient feature encoding methods to solve a similar problem. In particular the researchers develop a SLAM-aware system that incorporates a detection model to improve robotic object recognition. However, again, this work is distinct from our own as it performs recognition only in the point cloud space.

Lastly, Chen et al. (2015) describe a 3D object proposal method that is particularly focused on autonomous driving (e.g. detected objects include cars, pedestrians and cyclists). In this work, researchers similarly use known geometric features to reduce possible candidate regions and then propose an energy minimization formulation to compute region proposals. This approach places a greater emphasis on finding the *best* bounding box for each object, which requires additional prior information such as object size priors, point densities and free space information. As we perform localization in the point cloud space, the fit of the bounding box to a hand-labeled source is not nearly as important. While tight bounding boxes may be important in applications like self driving cars (as this information may be necessary for both high level planning and low-level dynamics considerations), the increased fidelity requires additional computation and therefore the system runs at an average of 1.2 seconds per image (with N=2000 proposals) at runtime. We instead focus on domestic environments (like Rusu et al. (2009b) and Stückler et al. (2013)), which allows us to retain the desired recognition accuracy at significantly improved speed by relaxing the requirements on our region proposal to *any* bounding box suitable for object recognition.
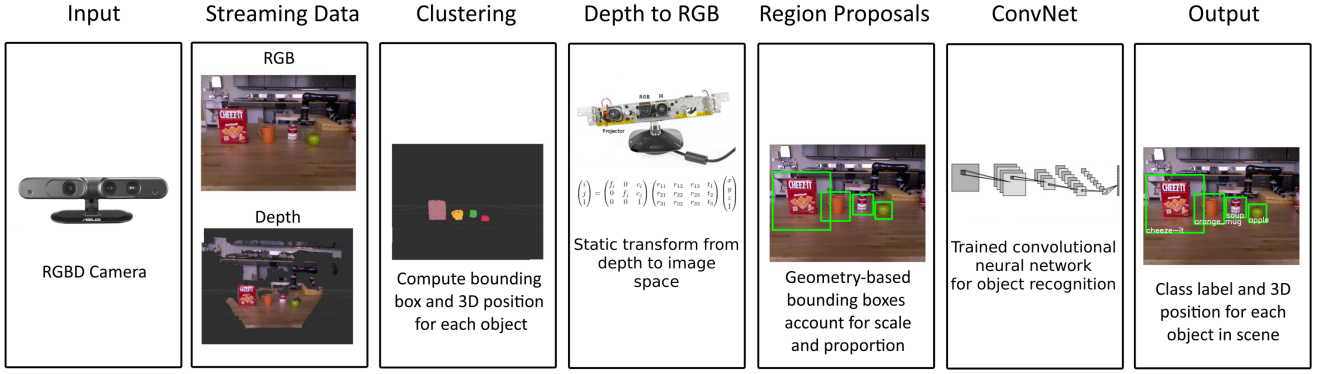
| Input | Streaming Data | Clustering | Depth to RGB | Region Proposals | ConvNet | Output |
|---|---|---|---|---|---|---|

**Figure 1.** High-level step-through of our object detection pipeline. From left to right: We use an RGB-D camera to capture color and depth information about our scene. We then exploit known geometric properties to compute a bounding box and 3D position for each visible object. The bounding boxes are translated to RGB space using the static transform defined by the position of the two sensors in the camera. These patches are fed through a trained CNN for object recognition. The output is a class label and 3D position for each object in the scene.

## 3  Our Approach

Our approach leverages both RGB and depth sensing modalities in a single object detection pipeline. In this work we focus specifically on detection in a tabletop setting, a common environment for assistive manipulators and particularly useful to researchers. We take inspiration from the computer vision community and develop a novel region proposal method, however, our technique is rooted in robotics perception and makes use of three dimensional point cloud data. To do so, we exploit the geometry of the environment to produce a minimum set of region proposals described in Section 3.1. We then translate our candidate regions from three dimensional bounding boxes into their two dimensional representation in the image plane, described in Section 3.2. This image patch is then fed into a CNN for classification, described in Section 3.3. The complete approach is outline in Figure 1.

### 3.1  Object Localization

Our object localization method is detailed in Algorithm 1. This algorithm simultaneously computes a bounding box, $\mathbb{B}$, and three dimensional position, $\mathbb{P}$, of each object in the scene. The localization method capitalizes on known geometric properties of the table to reduce the computational burden and produce highly reproducible results.

The input to the algorithm is a point cloud, $C$, which we then downsample (Line 2) to ensure coverage and speed. We downsample the input with a voxel filter which reduces the number of voxels necessary to represent the scene by replacing each set of $V$ voxels with a single voxel located at their centroid. The downsampling parameter, $\alpha$, defines the fraction of voxels in our final representation compared to the input point cloud — in our experiments this was set to 0.1. An optional step to further reduce the computational burden is to also run the point cloud through pass-through filters parameterized by the geometry of the tabletop (Line 3). These filters removes voxels in the scene that are outside the physical bounds of the table. Our experimental results in Section 6 use this step. We can further remove any points belonging to the tabletop itself by using the random sample consensus (RANSAC) method (Lines 4-5) to find

the dominant plane, $T$, in the point cloud scene, $C$. We can then filter these voxels from the remaining scene to remove any voxels belonging to the table. A Euclidean clustering algorithm is run on the remaining points to discover continuous objects in the scene (Line 6). We can then compute a bounding box $\mathbb{B}$ around each object in the set of clusters $o \in \mathbb{O}$ by finding the upper left, $U$, (Line 9) and lower right, $L$, (Line 10) corners of the cluster $o$. We also compute the three dimensional position $\mathbb{P}$ of an object by computing its centroid (Line 12).

---

**Algorithm 1** Geometric Region Proposal

1: **Given** Point Cloud $C$, *optional:* table dimensions
2: $C \leftarrow downsample(C, \alpha)$
3: *optional:* $C \leftarrow passthrough(C,$ table dimensions$)$
4: $T_{inliers} \leftarrow RANSAC(C)$
5: $T_{outliers} \leftarrow C - T_{inliers}$
6: $\mathbb{O} \leftarrow Cluster(T_{outliers})$
7: **Init** $\mathbb{B} \leftarrow \emptyset, \mathbb{P} \leftarrow \emptyset$
8: **for** $o \in \mathbb{O}$ **do**
9:     $U \leftarrow (x_{min}, y_{max}, z_{max})$
10:     $L \leftarrow (x_{max}, y_{min}, z_{max})$
11:     $\mathbb{B} \cup (U, L)$
12:     $\mathbb{P} \cup centroid(o)$
13: **return** $\mathbb{B}, \mathbb{P}$

---

Similar to other model-free segmentation approaches, a benefit of our method is that the region proposal algorithm does not rely on learning a model from a large dataset. Instead, we make use of the geometry of the scene to develop candidate object locations. Therefore, as only the recognition portion of the pipeline happens in the image-space, our model does not require training data that includes additional meta-data such as object bounding boxes. This approach has the added benefit of significantly decreasing the number of region proposals when compared to other methods. That is, for each object in the scene we propose only a *single* region by using the physical properties of the object to account for both position and scale. Our method is particularly well-suited for our problem domain as a vast number of manipulation objects are easily clustered due to their shape

and size. Additionally, even in cluttered environments, the depth dimensionality helps separate objects that otherwise look nearby in RGB space. Another benefit to computing the region proposals in the depth modality is that the localization of the object is very accurate due to the resolution and precision of the RGB-D sensor (Khoshelham and Elberink (2012)). For example, in our experiments the table was one meter wide, indicating a maximum error of $\sim 6mm.$*

## 3.2   Translation between Depth and RGB space

The next step in our object detection pipeline is to classify each proposal region. We choose to perform the classification in the image space due to the demonstrated accuracy and expressivity of deep learning methods. Therefore, we must translate the bounding box from the depth frame into the image frame. The coordinates of the bounding boxes in these two modalities are not directly aligned due to a physical offset in the sensor, however we can compute the transformation between them as described in Karan (2015).

To transform the bounding box in depth space to its representation in RGB space, we can begin by representing the RGB-D sensor as a pinhole camera. Under this assumption, each point in the depth space $(x, y, z) \in \mathbb{R}^3$ and each point in the image space $(i, j) \in \mathbb{R}^2$ is mapped into its homogeneous coordinate definitions, $(x, y, z, 1)$ and $(i, j, 1)$ respectively. We can then define a projective relationship between the two representations based on the intrinsic $(f, c)$ and extrinsic $(r, t)$ parameters of the camera as seen in Equation 1.

$$\begin{pmatrix} i \\ j \\ 1 \end{pmatrix} = \begin{pmatrix} f_i & 0 & c_i \\ 0 & f_j & c_j \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}$$
(1)

In this equation the first matrix represents the camera's intrinsic parameters and describes a transformation between the optical center of the camera and a given point in the image frame. Specifically, $f_i, f_j$ represent the focal length in pixel space and $c_i, c_j$ represent the physical offset between the origins of each frame in pixel space. The second matrix represents the camera's extrinsic parameters consisting of a rotation matrix, $R$, and the position of the origin in the world frame, $T$, which describes a transformation between the position and orientation of the depth- and RGB-cameras. These are defined by the sensor hardware and are often both readable and tunable using the associated driver (the values used in our implementation are included in the open source code). Through the use of Equation 1, we can therefore translate each bounding box in the point cloud to a bounding box in the image space.

Until this point, the bounding box we have computed tightly constrains each object in the scene, however, for the recognition portion of our pipeline it is useful to have a border around the object itself. This is because most image based recognition networks are trained with patches that include a border around the object of interest. For this reason, we slightly expand the bounding box associated with each object. The size of the border can be tuned (in our work, we expanded the border by $40\%$), however the same parameters should be used during data collection and at runtime.

## 3.3   Object Recognition

The final step in our object detection pipeline is recognition, which we solve using a convolutional neural network. Each region proposal is extracted from the full image, scaled to the input size needed for our trained model and classified.

The specific network architecture chosen for the classification portion of the pipeline is easily replaced and adjustable to stay in line with the state-of-the-art in deep learning. In our experiments (see Section 5), we evaluate a small model that we train from scratch as well as three state-of-the-art architectures initialized with weights learned on the ImageNet (Krizhevsky et al. (2012)) dataset and finetuned on our dataset. Importantly, using larger CNN models does not have a particularly large effect on the runtime of our system as we only evaluate the recognition model once per region proposal. Since there are no widely circulated networks weights trained on a dataset that encompass all of the objects we are interested in we are not able to evaluate our pipeline using a recognition network learned on a large dataset without some finetuning.

## 4   Dataset Acquisition

A secondary application of our region proposal method is data acquisition. Developing new datasets suitable for training deep learning models is normally a heavily human-time intensive process (Deng et al. (2009)). This is particularly important for robotics applications, where there is a dearth of pre-trained recognition models. Using our object localization method, researchers can quickly and easily create labeled data for objects not commonly found in circulated datasets. Our approach is described in Algorithm 2. This algorithm works by employing the use of Algorithm 1 on the set of object classes of interest. By placing an instance of a known object class in the view of the RGB-D sensor (Alg. 2, Line 4), we can store the streaming output of Algorithm 1 along with the user-provided label (Alg. 2, Lines 5-7) in a supervised learning dataset. This process is repeated for the full set of objects that a user is interested in at multiple locations throughout the scene. As Algorithm 1 is very fast, it is possible to store a large quantity of data very quickly.
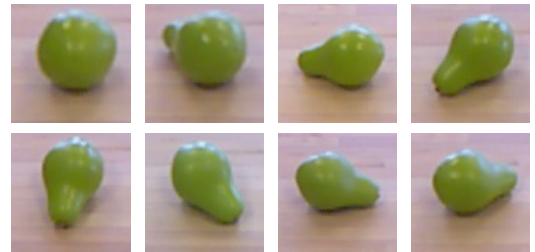


**Figure 2.** Example data captured using our object localization procedure with multiple orientations.

---

*The error in depth measurements using an RGB-D sensors is calculated through triangulation. The error increases with the distance squared as described in Zhang (2012).

**Algorithm 2** Dataset Acquisition
___
1: **Given** object class labels, $\mathbb{Y}_{labels}$
2: **Init** $\mathbb{X}_{data} \leftarrow \emptyset$, $\mathbb{Y}_{data} \leftarrow \emptyset$
3: **for** $y_{label} \in \mathbb{Y}_{labels}$ **do**
4:     place object of type $y_{label}$ in the scene
5:     $\mathbb{B}, \mathbb{P} \leftarrow$ Algorithm 1(object point cloud)
6:     $b_{rgb} \leftarrow$ Convert $\mathbb{B}$ to RGB space
7:     $x \leftarrow$ image patch defined by $b_{rgb}$
8:     $\mathbb{X}_{data} \leftarrow \{\mathbb{X}_{data} \cup x\}$
9:     $\mathbb{Y}_{data} \leftarrow \{\mathbb{Y}_{data} \cup y_{label}\}$
10: **return** $\mathbb{X}_{data}, \mathbb{Y}_{data}$
___

While capturing example images, it helps the model to generalize if the position and orientation of object(s) are altered thereby providing multiple views of each class. It can also help to alter aspects like lighting conditions and out-of-plane rotations. An example of the types of data collected via this method can be seen in Figure 2. Source code for the dataset acquisition process is a part of the released ROS package (Extension A).

## 5 Experimental Design

To evaluate the efficacy of our system, we analyze the accuracy of our approach in localizing and identifying objects in a variety of realistic tabletop scenarios. We begin by demonstrating the dataset building capabilities of our system (Section 5.1), which allows us to train our own CNN model from scratch and finetune three well known architectures whose weights were pre-trained on the ImageNet dataset (Section 5.2). We then evaluate our object detection pipeline on 40 realistic tabletop scenes. By evaluating our pipeline with four different recognition models, we demonstrate the ease with which one can update the underlying classification model to stay in line with the state-of-the-art. We compare the accuracy of the different models to examine the effect of different CNN architectures on the efficacy of our pipeline (Section 5.3).

### 5.1 Object Dataset

We began by creating an object dataset consisting of 19 object classes and 22 total object instances. The objects were almost exclusively taken from the YCB object dataset (Calli et al. (2015)). The specific objects chosen were relevant to our target domain: namely, common household items that a user of a robotic arm may wish to interact with. The full object set can be seen in Figure 3. The dataset we collected consists of a total of 2640 images split evenly by class.

### 5.2 Model Architectures

We trained and tested four different CNN architectures. The first is a small model that we trained from scratch. The other three are well tested architectures that we initialized with weights learned from the ImageNet dataset and finetuned on our own dataset. To train each model we used an $80/20\%$ train/test split of our dataset. During training of each network we also used different forms of data augmentation including random rotations, width and height shifts, shear mapping, and horizontal and vertical flipping. All models are implemented using the Keras library (Chollet (2015)).

*5.2.1 Small Model* The first network architecture that we evaluated is a small 6 layer convolutional neural network. The input layer is connected to a sequence of 3 convolutional layers with 3x3 filters. Each layer uses the ReLu non-linear activation function and is followed by 2x2 max pooling. This sequence is then followed by 3 additional layers of non-convolutional filters. During training, dropout is applied after the first two of the fully connected layers. The output layer is a learned soft-max classifier.

*5.2.2 VGG-16* The second network architecture that we evaluated is the VGG-16 network developed by Simonyan and Zisserman (2015). This network has 16 layers and was the first published work to use very small 3x3 convolutional filters. One of the key insights of this work was that sequences of small convolutional filters are capable of representing higher-order features otherwise captured by larger (more computationally expensive) receptive fields, like 7x7 or 9x9 filters. This network has previously been demonstrated to work well on many object recognition datasets.

*5.2.3 Inception Network* The third network architecture that we evaluated is the Inception v3 network developed by Szegedy et al. (2015). This network has 22 layers and is another popular architecture that that places a specific focus on reducing the necessary computation at test time to improve the model's efficiency. In this work, the authors propose the parallel computation of pooling, 1x1 and 3x3 filters which are then combined into a concatenated vector space (known as an inception module). By using 1x1 convolutions to reduce the filter space before computing the relatively more expensive 3x3 convolutions, the authors are able to reduce the computational complexity of this operation while improving the overall performance of the model.

*5.2.4 Residual Network* The fourth, and last, network architecture that we evaluated is the Resdiual network proposed by He et al. (2015). The standard implementation of this network is 152 layers deep and the result of this work is a network structure that allows one to train much deeper networks. In particular, the main insight of this work is the



**Figure 3.** Object set used to test detection pipeline. 22 object instances and 19 object classes in total. YCB Food: mustard, soup (x2), pringles, ground coffee, spam, jello (x2), apple (x2), pear, banana. YCB Kitchen: mug, bowl, bleach. YCB Shape: marbles, rubiks cube, soccer ball, softball, toy. Other: cup.
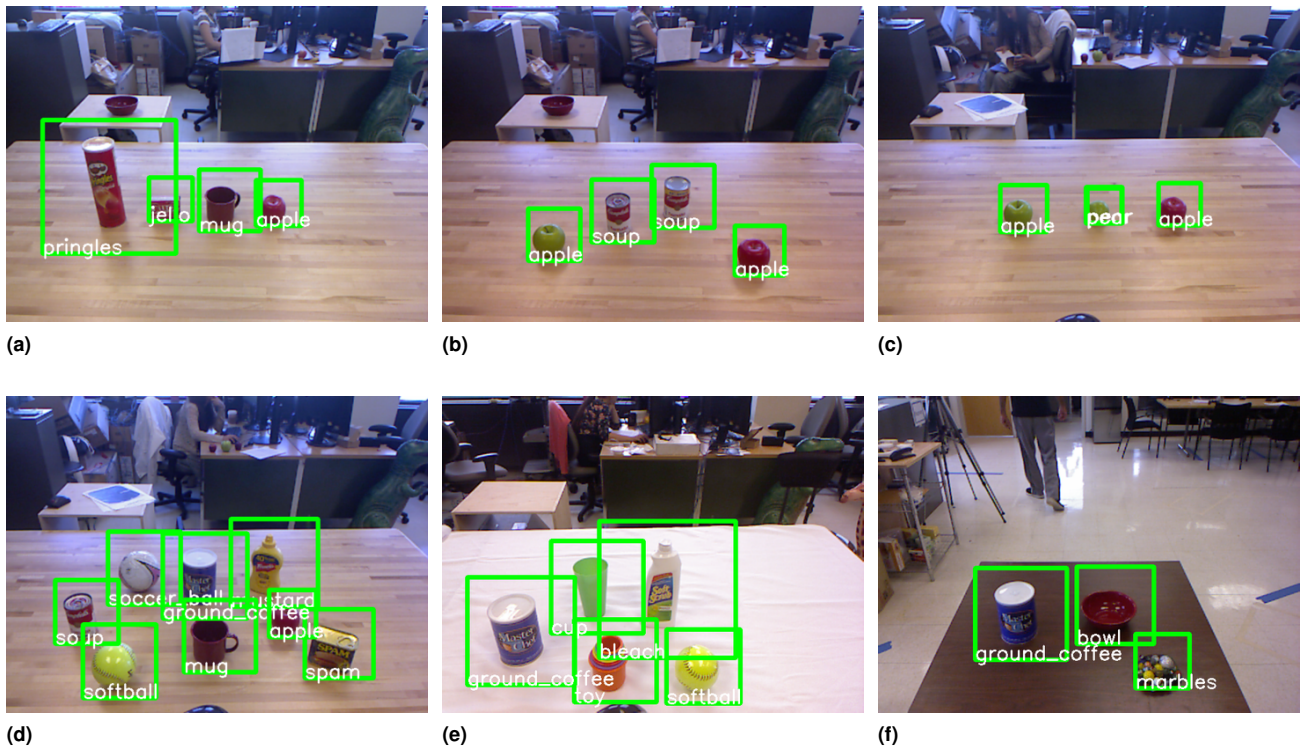
**Figure 4.** Six example scenes from our experimental testing set with various object configurations and environments. These scenes demonstrate the efficacy and highlight the generalizability of our approach. (a) Four objects, all a similar color. (b) Two distinct instances of two different classes. (c) Two distinct but very similar classes — note, the long edge of the pear is hidden. (d) Eight objects in various configurations. (e) Five objects on a white tablecloth. (f) Three objects on a dark brown table. The scenes in (e) and (f) are tested on tables distinct from the table used in the original data collection.

idea of *skip connections*. That is, instead of connecting the output of each layer to each following layer sequentially, one connects the output of each layer to the layer *after* the next layer. The concept behind this architecture is to encourage the network to learn *residual updates* from one layer to the next. In their paper, He et al. (2015) demonstrate that naively adding more and more layers to a network does not necessarily improve the performance of the network, while using residual connections dramatically improves the results.

All four models were successfully trained in less than 50 epochs with a batch size of 32. By observing the validation loss during training, it was clear that the pre-trained networks learned significantly quicker than the model trained from scratch. However, while all networks can run efficiently during test-time on a mid-tier GPU (nVidia 860m), only the smaller network can be trained on this GPU. Finetuning the larger networks requires a more powerful computer and GPU — during training we used an nVidia Titan X. We then transferred the network weights for these models to the less powerful computer for the experiments.

## 5.3 Evaluation Scenes

To analyze the accuracy of our pipeline, we developed 40 realistic tabletop scenes with varying numbers of objects, object configurations, clutter and backgrounds. The number of objects in a scene ranges from three to eight. In all evaluation scenes, the objects themselves are the same physical objects used to collect the training data, however

they are collected separately and we randomize each object's position and orientation in the evaluation scenes. Of the 40 different scenes, we include 20 in the same environment as the initial data collection, with random object configurations (Scenes 1-20). We then evaluate five in the same environment with the addition of a white tablecloth to hide the original table top surface (Scenes 21-25) and five on a new table with a much darker tabletop (Scenes 26-30). These two sets of tabletop scenes demonstrate the generalizability of the recognition models to novel environments. The final 10 scenes are made up of five scenes collected with a moving camera where the camera beginning on the left side of the environment and moving towards the right side (Scenes 31-35) and five with a moving camera with the camera beginning high up in the environment and slowly moving down (Scenes 36-40). This final set of experiments is particularly relevant to mobile robotics where the platform may be moving. The full set of 40 test scenes can be seen with their descriptions in Appendix B. Six example scenes can be seen in Figure 4. Our experimental scenes are similar to those released by Lai et al. (2011), however, we develop a larger number of scenes for testing and specifically focus on scenes that demonstrate particular capabilities of our approach (e.g. larger number of objects in a scene, invariance to object features like color, more cluttered environments and a variety of backgrounds).

**Average Recognition Accuracy**

| Model | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Small | 1.00 | 1.00 | *0.75* | *0.96* | *0.95* | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | *0.33* | *0.83* | 1.00 | 1.00 | 1.00 |
| VGG16 | 1.00 | *0.75* | *0.99* | 1.00 | *0.67* | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| Inception | 1.00 | 1.00 | 1.00 | 1.00 | *0.46* | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 0.91 | *0.85* | 1.00 | 1.00 |
| Residual | 1.00 | 1.00 | 1.00 | 1.00 | *0.66* | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | *0.95* | *0.67* | 1.00 | 1.00 |

| Model | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Small | *0.67* | 1.00 | 1.00 | 1.00 | 1.00 | *0.75* | 1.00 | *0.96* | *0.80* | *0.75* | *0.33* | *0.33* | 1.00 | *0.38* | *0.67* |
| VGG16 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | *0.95* | *0.74* | *0.80* | 1.00 | 1.00 | 1.00 | *0.79* | 1.00 | 1.00 |
| Inception | 1.00 | 1.00 | *0.67* | 1.00 | 1.00 | *0.98* | *0.96* | *0.37* | *0.80* | *0.95* | *0.83* | *0.71* | 1.00 | *0.95* | *0.93* |
| Residual | 1.00 | 1.00 | *0.75* | 1.00 | 1.00 | 1.00 | 1.00 | *0.75* | *0.80* | 1.00 | 1.00 | 1.00 | *0.67* | 1.00 | 1.00 |

| Model | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Small | 1.00 | *0.99* | 1.00 | *0.94* | *0.87* | 1.00 | *0.88* | 1.00 | *0.99* | 1.00 | *0.88* |
| VGG16 | *0.97* | 1.00 | 1.00 | 1.00 | *0.99* | *0.99* | 1.00 | 1.00 | 1.00 | 1.00 | *0.97* |
| Inception | *0.95* | 1.00 | 1.00 | *0.71* | *0.99* | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | *0.93* |
| Residual | 1.00 | 1.00 | 1.00 | *0.99* | *0.97* | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | *0.95* |

**Table 1.** Recognition Accuracy. The average recognition of all four models on each of the 40 scenes we tested in our experiments. The value computed for each scene is an average of the accuracy over all objects in the first 100 frames captured at 5Hz using an RGB-D camera. Accuracy falling short of 100% is italicized.

## 6 Results

To evaluate the accuracy of our pipeline on the scenes described in Section 5.3, we collect 100 continuous pointcloud frames captured at 5Hz from the RGB-D camera and compare the output of our pipeline to hand-labeled ground truth. Our experiments were run on a Core i7 laptop with a mid-tier mobile GPU (nVidia GeForce 860M). We compute the average accuracy of the pipeline by taking the mean of the accuracy over those 100 frames. That is, for the pipeline to achieve 100% on a given scene, it needs to correctly predict the correct object class for *each object* in *each frame*. Averaging the success of our approach over all the frames is particularly important for the scenes in which the camera is moving as there is greater variation between frames than when the camera and objects are static. Example results can be seen in Figure 4 and a full breakdown is presented in Table 1.

The results demonstrate that all four tested architectures are able to perform well on the 40 experimental test scenes. Of note is the fact that these scenes were varied in the amount of clutter, in the orientation and position of the objects, in the presentation of which objects, and in some cases included novel tabletops (Appendix B, Figures 8a-8j) and camera motion (Appendix B, Figures 8k-8t). In order of increasing performance, we find that the small model achieved an average accuracy of 88% on the full test set, the Inception network achieved an average accuracy of 93%, the Residual network achieved an average accuracy of 95% and the VGG-16 network achieved an average accuracy of 97%. While we do observe a clear difference between the small model and the three well known models, recall that the small model is trained from scratch on a mid-tier GPU while the known models are first pre-trained on the ImageNet dataset and then further trained on a more powerful GPU. We note that even in robotics applications that require on-board computation, it is often possible to train a model off-line on a more capable computer so long as it can later run on the on-board computer at test time.

To asses the suitability of our approach for use by robots in making online predictions, we also evaluate the running time of our system. The full pipeline runs at an average of 12Hz, which is suitable for robotic manipulation tasks in our target domain (i.e. households).

## 7 Discussion

The presented paradigm appears to be a promising direction for practical robotic perception systems. The marriage of state-of-the-art techniques from robotics and computer vision helps produce a fast and accurate object detection framework that can be easily incorporated into any tabletop manipulation task. It is a real-time system that does not require top-of-the-line hardware and produces competitive results. This methodology is additionally useful for creating novel datasets which suggests that this approach could be useful for other researchers and advanced users alike.

Our approach differs from related work in three main ways. The first is that unlike image-only based approaches, we use the 3D geometric features of a scene to compute the region proposals that we then feed to our recognition model for classification. This ensures that we only send a single image patch per object in the scene which is extremely efficient when compared to state-of-the-art image-based approaches. The second way in which our work differs, is that unlike point-cloud based approaches that both localize and recognize the objects in the depth space, we localize points in the depth space, but recognize the objects in the image space using convolutional neural networks. Deep learning based approaches have proven extremely effective in the computer vision community and our work demonstrates their applicability to robotics as well. In particular, we observe 97% average precision over all scenes by the best performing model (Table 1). Lastly, unlike with image-based approaches which locate objects only in 2D, our

method produces the precise 3D position of the object with minimal additional computation.

In the remainder of the discussion we will compare our approach to image-based methods (Section 7.1), examine the generalizability of our system (Section 7.2), note specific cases of failure (Section 7.3) and discuss future directions (Section 7.4).

## 7.1 Comparison to Image-Based Methods

Ideally, our analysis would have included a direct comparison between our geometric region proposal method and state-of-the-art image-based approaches from the computer vision community (e.g. R-CNN (Girshick et al. (2014)), Fast R-CNN (Girshick (2015)) and Faster R-CNN (Ren et al. (2015))). However, there are a number of key differences between our approach and this body of work that (1) make such a direct comparison challenging and (2) highlight some of the gains of our approach (for use within robotics in particular). These differences include the required training data, available 3D information and execution speed, which are discussed next.

*7.1.1 Training Data Requirements* One notable distinction is a significant difference in the *type* and *amount* of required training data.

In our object detection pipeline, objects are localized autonomously by exploiting known geometric properties of the scene. All that the human provides is a class label. By comparison, Fast R-CNN and Faster R-CNN require training data that includes images labeled not only with the object class, but also with the corresponding bounding box for each object in the scene. Each bounding box is drawn by hand—a significant human-time intensive process. For this reason, it is unlikely that this type of data will become widely available for all objects of interest in our target domain (the home) in the near future, which prohibits the training or finetuning any of these models on novel datasets (that do not include bounding box labels).

In addition to the required localization data, the image-only based approaches rely much more heavily on the training data including a *background* class, which is used to recognize false positives nominated by the region proposal method. Unlike Fast R-CNN and Faster R-CNN, the original R-CNN approach (a significantly slower approach) does not require training a network that localizes objects in a scene. However, likewise, it still utilizes a region proposal method that proposes (on average) significantly more object locations than there actually are objects (R-CNN and Fast R-CNN generate ∼2000 object proposals and Faster R-CNN generates ≤ 300). To solve this problem, each of the aforementioned approaches utilize a background class which can reject the false positives. Again, this requires collecting more training data. Even worse, it increases the likelihood of a false positive at the end of the pipeline, whereas we saw *zero* false positives (related to the existence of an object) in our experiments.

*7.1.2 3D Localization* Another notable distinction is that approaches from the computer vision community only solve the localization problem in 2D. This is insufficient for robotics applications as the robot itself exists in 3D and must interact with other objects in the same space. Additionally,

not only do these computer vision models lack depth information, but that information is not directly available in the datasets used to train these models (such as ImageNet), which increases the difficulty of incorporating depth into these models.

*7.1.3 Speed at Test Time* We furthermore compare our runtime to that of competitive approaches in speed and accuracy, the image-based R-CNN, Fast R-CNN and Faster R-CNN methods. Comparisons were run on the same computer, with the same size images (640×480), and using test set images from the same dataset used to create a model's training set. Under these conditions, R-CNN was able to run at an average of 0.4Hz (using Selective Search (Uijlings et al. (2013)) for region proposals), Fast R-CNN was able to run at an average of 1.33Hz and Faster R-CNN was able to run at an average of 5Hz. Our approach ran at an average of 12Hz, and thus demonstrates a speed up factor of 30x, 9.2x and 2.4x, respectively.

## 7.2 Generalizability

To asses generalizability, we tested our pipeline on 40 realistic tabletop scenes covering a large variety of possible object classes and configurations. The scenes vary in which objects are present, where they are located, their orientation and the amount of clutter (see Appendix B). We also examine the effect of a moving camera and testing on tabletops that are unique from the one used during data collection.

We highlight a few noteworthy examples. The scene shown in Figure 4a demonstrates that the recognition system is capable of distinguishing between objects of *similar color*. The scene shown in Figure 4b demonstrates that the system is able to generalize to different instances of physical objects of the *same class*, and the scene shown in Figure 4c demonstrates that the system is able to differentiate between classes that are *extremely similar*. In particular, in this last image, notice that the pear is intentionally oriented away from the camera hiding the top half of the pear, likely the largest visually differentiating factor between a pear and an apple. In Figure 5 we see that the recognition system is able to account for *out-of-plane orientation* changes (note the orientation of the wood block in both images).

We also demonstrate that the system performs well when there are *many objects* in the scene (Figure 4d) and when the tabletop is distinct in appearance from the tabletop used to capture the training data (Figures 4e and 4f as well as Appendix B, Figures 8a-8j). The system moreover performs well when the *camera is moving* during the data capture,
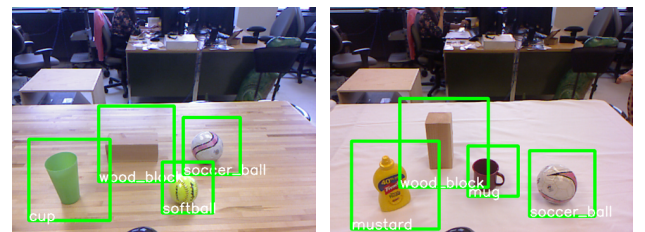


**Figure 5.** In both scenes the pipeline is able to correctly localize and recognize all objects. In particular, by observing the wood block in both scenes we see that the recognition network is capable of accounting for out-of-plane orientation changes.

demonstrating that this approach is viable for mobile robots (Appendix B, Figures 8k-8t). In the first five of these scenes, the camera was held at a constant height and was moved about one meter from the left side of the scene to the right. In the second five of these scenes, the camera was moved about one meter from the top of the scene to the bottom.

### 7.3 Failures

In the cases of recognition failure, it was often fairly clear why there was a misclassification. For example, one of the more common mis-labelings was the toy and rubiks cube which have similar color schemes (Figure 6). However, it is clear from the results that these small errors can be easily overcome by improved modeling techniques. For example, Scenes 16, 25 and 35 each include either the toy or rubiks cube and our results show that the small model is unable to correctly identify these objects. However, the three models pre-trained on the (larger) ImageNet dataset are able to correctly classify all objects in each of these scenes (Table 1).



**Figure 6.** Visually similar objects. Left: Rubiks cube. Right: Toy.

Additionally, when reviewing the accuracy of these models over the entire test set, we see a large improvement when we move from the small model trained from scratch to the pre-trained models that we finetuned. In particular it appears clear that finetuning a pre-trained network helps with generalization to different environments. For example, we see a significantly larger effect of the background color on the small model (where the average accuracy of the model drops to 70%) than any of the pre-trained networks (e.g. the accuracy of the VGG16 only drops to 93%). However, the trade-off is that the three larger pre-trained networks needed to be finetuned on the more powerful and more expensive Titan X GPU.

One point of failure is how the recognition model handles objects that were not in the training set. Our current approach will choose the most likely class label as defined by the probabilities that we get from the softmax output layer of each network. However, this choice is not well suited under the open set world assumption where we expect to see novel objects that were not included in the training data. Instead a robot needs to be aware of when it comes across a novel object. To solve this problem, we can incorporate statistical techniques for detecting class outliers and incorporating novel objects as described in Bendale and Boult (2015).

### 7.4 Future Directions

We expect that highly cluttered environments will require improved segmentation approaches in the point cloud representation. In particular, it is likely that depth-only segmentation will fail in scenes where objects actually sit on top of one another, for example, imagine objects sitting on a bookshelf. A potential area of further work in this research would be to combine the current depth-based segmentation with image-based segmentation approaches to incorporate color information as well into the segmentation procedure.

An additional area of possible improvement and refinement in our system would be to include the depth information from the RGB-D camera as a fourth channel in our CNN architecture. Similar to the benefit of using color information in the segmentation process, the recognition portion of our pipeline could be improved by incorporating depth information into the model. However, at least for now, this would reduce our ability to finetune pre-trained networks (a cheap and efficient way to use features learned from larger datasets) as currently these networks only include color information.

## 8 Conclusion

In this paper, we describe and demonstrate a simple, and fast, object detection pipeline for tabletop manipulation tasks using robot vision. We validate the efficacy of our approach with a thorough study to test the speed, accuracy and generalizability of our method. We found our system to be capable of running in real-time (12Hz) on limited-capability hardware. The described system owes its speed and computational efficiency to the minimal set of regions proposed through unsupervised methods of analysis in the point cloud space. We also demonstrate that our method makes it easy to collect novel datasets which can be used to train recognition models from scratch or used to finetune models pre-trained on larger datasets. The modular design of the pipeline makes it easy to incorporate new recognition models in order to stay in-line with the state-of-the-art. In our experiments, we found that incorporating the state-of-the-art CNN architectures allowed us to achieve a 97% detection accuracy on our varied experimental dataset. Our approach owes it accuracy and generalizability in the recognition space to Convolutional Neural Networks.

In future work we hope to improve the capabilities of our approach by incorporating color information into our segmentation approach and incorporating depth information into our recognition models. We also plan to demonstrate how this approach can be used under the open-set assumption. We believe that integrating these changes will also allow us to expand the classes of scenes in which this approach can be validated.

The code for both the full object detection pipeline as well as the dataset acquisition portion can be found at https://github.com/asbroad/geom_rcnn and are included with an open-source MIT license in Extension A.

## Acknowledgments

## References

Abhijit Bendale and Terrance Boult. Towards Open World Recognition. In *Proceedings of the IEEE Conference on*

*Computer Vision and Pattern Recognition (CVPR)*, pages 1893–1902, 2015.

Berk Calli, Aaron Walsman, Arjun Singh, Siddhartha Srinivasa, Pieter Abbeel, and Aaron Dollar. Benchmarking in Manipulation Research: The YCB Object and Model Set and Benchmarking Protocols. In *IEEE Robotics and Automation Magazine*, August 2015.

Xiangrong Chen and Alan L Yuille. A Time-Efficient Cascade for Real-Time Object Detection: With Applications for the Visually Impaired. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 28–28, 2005.

Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Andrew G Berneshawi, Huimin Ma, Sanja Fidler, and Raquel Urtasun. 3D Object Proposals for Accurate Object Class Detection. In *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, pages 424–432, 2015.

Franois Chollet. keras. https://github.com/fchollet/keras, 2015.

Camille Couprie, Clement Farabet, Laurent Najman, and Yann LeCun. Convolutional Nets and Watershed Cuts for Real-Time Semantic Labeling of RGBD Videos. *Journal of Machine Learning Research*, 15:3489–3511, 2014.

Meir Johnathan Dahan, Nir Chen, Ariel Shamir, and Daniel Cohen-Or. Combining Color and Depth for Enhanced Image Segmentation and Retargeting. *The Visual Computer*, 28(12): 1181–1193, 2012.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009.

Ross Girshick. Fast R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1440–1448, 2015.

Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 580–587, 2014.

Saurabh Gupta, Ross Girshick, Pablo Arbeláez, and Jitendra Malik. Learning Rich Features from RGB-D Images for Object Detection and Segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 345–360, 2014.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

Jan Hosang, Rodrigo Benenson, and Bernt Schiele. How Good are Detection Proposals, Really? In *Proceedings of the British Machine Vision Conference (BMVC)*, 2014.

Huaizu Jiang and Erik Learned-Miller. Face detection with the faster R-CNN. *arXiv:1606.03473*, 2016.

Branko Karan. Calibration of Kinect-type RGB-D sensors for robotic applications. *FME Transactions*, 43(1):47–54, 2015.

Andrej Karpathy and Li Fei-Fei. Deep Visual-Semantic Alignments for Generating Image Descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3128–3137, 2015.

Kourosh Khoshelham and Sander Oude Elberink. Accuracy and Resolution of Kinect Depth Data for Indoor Mapping Applications. *Sensors*, 12(2):1437–1454, 2012.

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet Classification with Deep Convolutional Neural Networks. In *Proceedings of Advances in Neural Information Processing Systems*, pages 1097–1105, 2012.

Kevin Lai, Liefeng Bo, Xiaofeng Ren, and Dieter Fox. A Large-Scale Hierarchical Multi-View RGB-D Object Dataset. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 1817–1824, 2011.

Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.

Raul Mur-Artal, JMM Montiel, and Juan D Tardos. ORB-SLAM: A Versatile and Accurate Monocular SLAM System. *IEEE Transactions on Robotics*, 31(5):1147–1163, 2015.

Sudeep Pillai and John Leonard. Monocular SLAM Supported Object Recognition. In *Proceedings of Robotics: Science and Systems (RSS)*, July 2015.

Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, pages 91–99, 2015.

Radu Bogdan Rusu, Nico Blodow, and Michael Beetz. Fast Point Feature Histograms (FPFH) for 3D Registration. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 3212–3217, 2009a.

Radu Bogdan Rusu, Nico Blodow, Zoltan Csaba Marton, and Michael Beetz. Close-Range Scene Segmentation and Reconstruction of 3D Point Cloud Maps for Mobile Manipulation in Domestic Environments. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–6, 2009b.

Radu Bogdan Rusu, Gary Bradski, Romain Thibaux, and John Hsu. Fast 3D Recognition and Pose using the Viewpoint Feature Histogram. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2155–2162, 2010.

Ashutosh Saxena, Min Sun, and Andrew Y. Ng. Make3D: Learning 3D Scene Structure from a Single Still Image. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 31(5):824–840, 2009.

Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A Unified Embedding for Face Recognition and Clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 815–823, 2015.

Karen Simonyan and Andrew. Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2015.

Shuran Song and Jianxiong Xiao. Sliding Shapes for 3D Object Detection in Depth Images. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 634–651. 2014.

Jörg Stückler, Ricarda Steffens, Dirk Holz, and Sven Behnke. Efficient 3D Object Perception and Grasp Planning for Mobile Manipulation in Domestic Environments. *Robotics and Autonomous Systems*, 61(10):1106–1115, 2013.

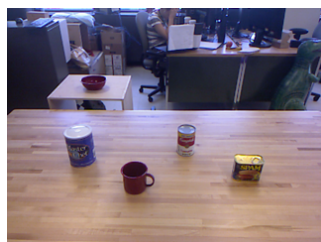Christian Szegedy, Alexander Toshev, and Dumitru Erhan. Deep Neural Networks for Object Detection. In *Proceedings of*

*Advances in Neural Information Processing Systems (NIPS)*, pages 2553–2561, 2013.

Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. *arXiv:1512.00567*, 2015.

Jie Tang, Stephen Miller, Arjun Singh, and Pieter Abbeel. A Textured Object Recognition Pipeline for Color and Depth Image Data. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 3467–3474, 2012.

Federico Tombari, Samuele Salti, and Luigi Di Stefano. Unique Signatures of Histograms for Local Surface Description. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 356–369. 2010.

Jasper RR Uijlings, Koen EA van de Sande, Theo Gevers, and Arnold WM Smeulders. Selective Search for Object Recognition. *International Journal of Computer Vision*, 104 (2):154–171, 2013.

Liliang Zhang, Liang Lin, Xiaodan Liang, and Kaiming He. Is Faster R-CNN Doing Well for Pedestrian Detection? In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 443–457, 2016.

Zhengyou Zhang. Microsoft Kinect Sensor and its Effect. *IEEE MultiMedia*, 19(2):4–10, 2012.

C Lawrence Zitnick and Piotr Dollár. Edge Boxes: Locating Object Proposals from Edges. In *Proceedings of European Conference on Computer Vision (ECCV)*. 2014.

## A   Index to Multimedia Extensions

The multimedia extensions to this article are at: `www.ijrr.org`.

| Extension | Media Type | Description |
|---|---|---|
| 1 | Code | An efficient framework for 3D object detection in Robotics applications. It is specifically designed for detecting objects in table-top scenes (or similar environments with objects sitting on a dominant plane). The output of the system is a class label and 3D position for each object in the scene. |

## B   Evaluation Data

**(a)** Scene 1. Four objects: ground coffee, mug, soup, spam.

**(b)** Scene 2. Four objects: apple, ground coffee, mustard, toy.

**(c)** Scene 3. Four objects: bleach, softball, marbles, pringles.

**(d)** Scene 4. Four objects: two different apples and soups.
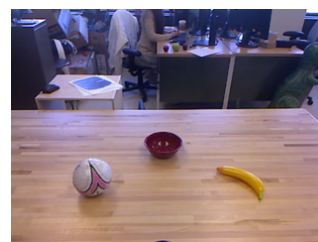
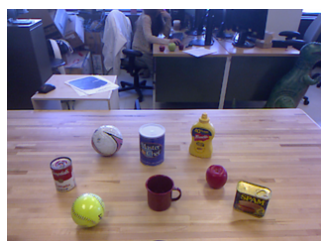**(e)** Scene 5. Three objects: ground coffee, mustard, soup. Two objects are occluded.

**(f)** Scene 6. Four objects: pringles, jellow, mug, apple. All of which have similar colors.

**(g)** Scene 7. Three objects: two apples and a pear. The long edge of the pear is hidden.

**(h)** Scene 8. Three objects: soccer ball, bowl, banana.

**(i)** Scene 9. Eight objects: soup, soccer ball, ground coffee, mustard, apple, spam, mug, softball.

**(j)** Scene 10. Three objects: wood block, cup, pear.

**(k)** Scene 11. Three objects: jello, bleach, rubiks cube.

**(l)** Scene 12. Five objects: bowl, cup, jello, banana, marbles.

**(m)** Scene 13. Three objects: banana, wood block, jello.

**(n)** Scene 14. Four objects: bowl, jello, cup, pear.

**(o)** Scene 15. Four objects: cup, wood block, softball, soccer ball.

**(p)** Scene 16. Three objects: rubiks cube, toy, marbles. All objects have a similar combination of colors.

**(q)** Scene 17. Seven objects: mustard, spam, soup, apple, jello, banana, pear.

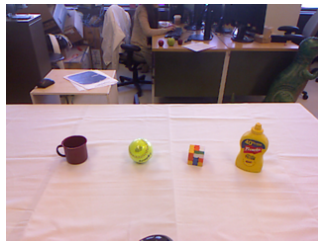**(r)** Scene 18. Three objects: wood block, softball, mug.

**(s)** Scene 19. Three objects: mustard, banana, pear.
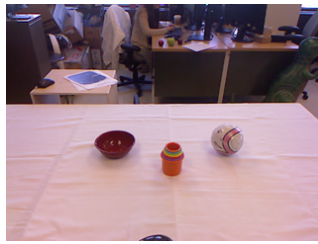
**(t)** Scene 20. Three objects: soccer ball, spam, jello.

**Figure 7.** Example test scenes (1-20)

**(a)** Scene 21. Four objects: mug, softball, rubiks cube, mustard. The table is covered in a white tablecloth.

**(b)** Scene 22. Three objects: bowl, toy, soccer ball. The table is covered in a white tablecloth.

**(c)** Scene 23. Four objects: mustard, wood block, mug, soccer ball. The table is covered in a white tablecloth.

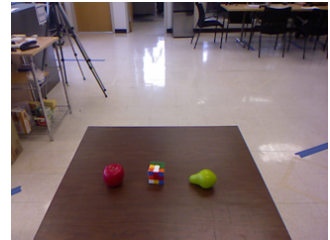**(d)** Scene 24. Five objects: ground coffee, cup, toy, bleach, softball. The table is covered in a white tablecloth.

**(e)** Scene 25. Four objects: mug, toy, soccer ball, banana. The table is covered in a white tablecloth.

**(f)** Scene 26. Three objects: mustard, ground coffee, softball. The tabletop is dark brown.

**(g)** Scene 27. Three objects: apple, soup, pear. The tabletop is dark brown.

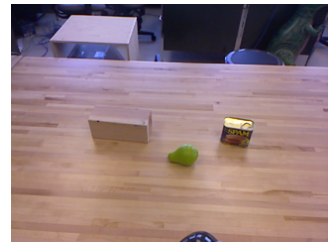**(h)** Scene 28. Three objects: pringles, banana, cup. The tabletop is dark brown.

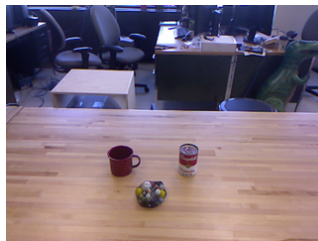**(i)** Scene 29. Three objects: pringles banana, cup. The tabletop is dark brown.

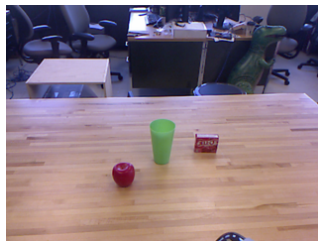**(j)** Scene 30. Three objects: ground coffee, bowl, marbles. The tabletop is dark brown.

**(k)** Scene 31. Three objects: mustard, softball, bowl. The camera was moving left to right.

**(l)** Scene 32. Three objects: wood block, pear, spam. The camera was moving left to right.
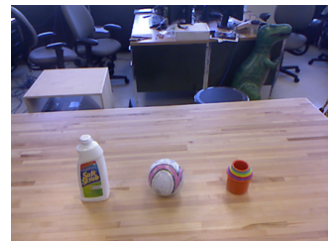
**(m)** Scene 33. Three objects: mug, marbles, soup. The camera was moving left to right.

**(n)** Scene 34. Three objects: apple, cup, jello. The camera was moving left to right.

**(o)** Scene 35. Three objects: pringles rubiks cube, banana. The camera was moving left to right.
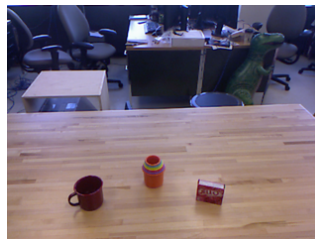
**(p)** Scene 36. Three objects: bleach, soccer ball, rubiks cube. The camera was moving from high to low.

**(q)** Scene 37. Three objects: mug, ground coffee, apple. The camera was moving from high to low.

**(r)** Scene 38. Three objects: cup, ground coffee, soccer ball. The camera was moving from high to low.

**(s)** Scene 39. Three objects: mug, toy, jello. The camera was moving from high to low.

**(t)** Scene 40. Three objects: mustard, softball, bowl. The camera was moving from high to low.

**Figure 8.** Example test scenes (21-40)