

Chapter 3

RGB-D image-based Object Detection: from Traditional Methods to Deep Learning Techniques

RGB-D Image Analysis and Processing

Isaac Ronald Ward, Hamid Laga, and Mohammed Bennamoun

Abstract Object detection from RGB images is a long-standing problem in image processing and computer vision. It has applications in various domains including robotics, surveillance, human-computer interaction, and medical diagnosis. With the availability of low cost 3D scanners, a large number of RGB-D object detection approaches have been proposed in the past years. This chapter provides a comprehensive survey of the recent developments in this field. We structure the chapter into two parts; the focus of the first part is on techniques that are based on hand-crafted features combined with machine learning algorithms. The focus of the second part is on the more recent work, which is based on deep learning. Deep learning techniques, coupled with the availability of large training datasets, have now revolutionized the field of computer vision, including RGB-D object detection, achieving an unprecedented level of performance. We survey the key contributions, summarize the most commonly used pipelines, discuss their benefits and limitations, and highlight some important directions for future research.

3.1 Introduction

Humans are able to efficiently and effortlessly detect objects, estimate their sizes and orientations in the 3D space, and recognize their classes. This capability has long been studied by cognitive scientists. It has, over the past two decades, attracted a lot of interest from the computer vision and machine learning communities mainly be-

Isaac Ronald Ward
University of Western Australia, e-mail: isaac.ward@uwa.edu.au

Hamid Laga
Murdoch University, Australia, and The Phenomics and Bioinformatics Research Centre, University of South Australia, e-mail: H.Laga@murdoch.edu.au

Mohammed Bennamoun
University of Western Australia, e-mail: mohammed.bennamoun@uwa.edu.au

cause of the wide range of applications that can benefit from it. For instance, robots, autonomous vehicles, and surveillance and security systems rely on accurate detection of 3D objects to enable efficient object recognition, grasping, manipulation, obstacle avoidance, scene understanding, and accurate navigation.

Traditionally, object detection algorithms operate on images captured with RGB cameras. However, in the recent years, we have seen the emergence of low-cost 3D sensors, hereinafter referred to as *RGB-D sensors*, that are able to capture depth information in addition to RGB images. Consequently, numerous approaches for object detection from RGB-D images have been proposed. Some of these methods have been specifically designed to detect specific types of objects, e.g. humans, faces and cars. Others are more generic and aim to detect objects that may belong to one of many different classes. This chapter, which focuses on generic object detection from RGB-D images, provides a comprehensive survey of the recent developments in this field. We will first review the traditional methods, which are mainly based on hand-crafted features combined with machine learning techniques. In the second part of the chapter, we will focus on the more recent developments, which are mainly based on deep learning.

The chapter is organized as follows; Section 3.2 formalizes the object detection problem, discusses the main challenges, and outlines a taxonomy of the different types of algorithms. Section 3.3 reviews the traditional methods, which are based on hand-crafted features and traditional machine learning techniques. Section 3.4 focuses on approaches that use deep learning techniques. Section 3.5 discusses some RGB-D-based object detection pipelines and compares their performances on publicly available datasets, using well-defined performance evaluation metrics. Finally, Section 3.6 summarizes the main findings of this survey and discusses some potential challenges for future research.

3.2 Problem statement, challenges, and taxonomy

Object detection from RGB-D images can be formulated as follows; given an RGB-D image, we seek to find the location, size, and orientation of objects of interest, e.g. cars, humans, and chairs. The position and orientation of an object is collectively referred to as the *pose*, where the orientation is provided in the form of Euler angles, quaternion coefficients or some similar encoding. The location can be in the form of a 3D bounding box around the visible and/or non-visible parts of each instance of the objects of interest. It can also be an accurate 2D/3D segmentation, i.e. the complete shape and orientation even if only part of the instance is visible. In general, we are more interested in detecting the whole objects, even if parts of them are not visible due to clutter, self-occlusions, and occlusion with other objects. This is referred to as *amodal* object detection. In this section, we discuss the most important challenges in this field (Section 3.2.1) and then lay down a taxonomy for the state-of-the-art (Section 3.2.2).

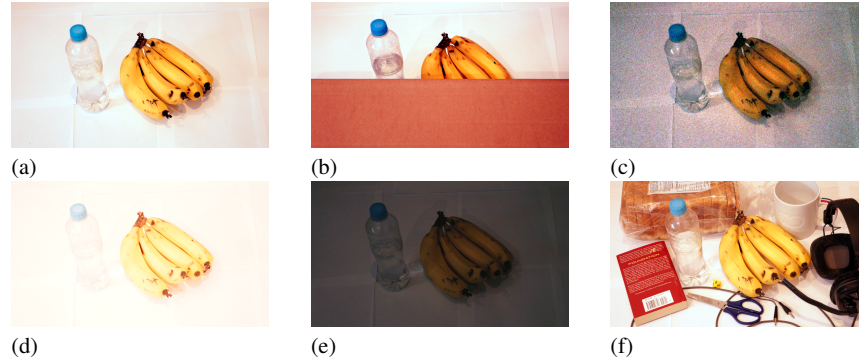


Fig. 3.1: Illustration of some extrinsic challenges in object detection. (a) Objects of interest are clearly separated from each other and from the uniform background. (b) Two objects partially occluded by a cardboard box. (c) Sensor noise that might affect images. (d) An overexposed image. (e) An underexposed image. (f) A cluttered image, which hinders the detection of smaller and occluded objects.

3.2.1 Challenges

Though RGB-D object detection has been extensively investigated, there are a number of challenges that efficient solutions should address. Below, we classify these challenges into whether they are due to intrinsic or extrinsic factors. Extrinsic factors refer to all the *external* factors that might affect object detection (see Fig. 3.1). Extrinsic challenges include:

- **Occlusions and background clutter.** The task of object detection algorithms is to not only localize objects in the 3D world, but also to estimate their physical sizes and poses, even if only parts of them are visible in the RGB-D image. In real-life situations, such occlusions can occur at anytime, especially when dealing with dynamic scenes. Clutter can occur in the case of indoor and outdoor scenes. While biological vision systems excel at detecting objects under such challenging situations, occlusions and background clutter can significantly affect object detection algorithms.
- **Incomplete and sparse data.** Data generated by RGB-D sensors can be incomplete and even sparse in some regions, especially along the z -, i.e. depth, direction. Efficient algorithms should be able to detect the full extent of the object(s) of interest even when significant parts of it are missing.
- **Illumination.** RGB-D object detection pipelines should be robust to changes in lighting conditions. In fact, significant variations in lighting can be encountered in indoor and outdoor environments. For example, autonomously driving drones and domestic indoor robots are required to operate over a full day-night cycle and are likely to encounter extremes in environmental illumination. As such, the appearance of objects can be significantly affected not only in the RGB image

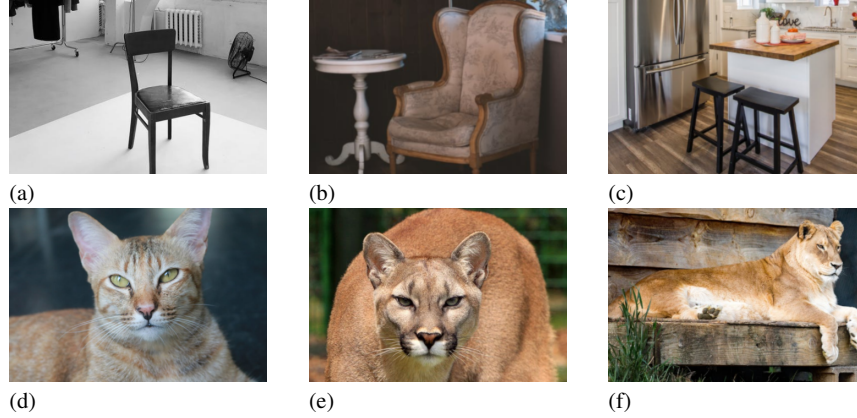


Fig. 3.2: Illustration of some intrinsic challenges in object detection. (a-c) Intra-class variations where objects of the same class (chair) appear significantly different. (d-f) Inter-class similarities where objects belonging to different classes (cat, cougar & lion) appear similar.

but also in the depth map, depending on the type of 3D sensors used for the acquisition.

- **Sensor limitations.** Though sensor limitations classically refer to colour image noise that occurs on imaging sensors, RGB-D images are also prone to other unique sensor limitations. Examples include spatial and depth resolution. The latter limits the size of the objects that can be detected. Depth sensor range limitations are particularly noticeable, e.g. the Microsoft Kinect, which is only sufficiently accurate to a range of approximately 4.5m [79]. This prevents the sensor from adequately providing RGB-D inputs in outdoor contexts where more expensive devices, e.g. laser scanners, may have to be used [56].
- **Computation time.** Many applications, e.g. autonomous driving, require real-time object detection. Despite hardware acceleration, using GPUs, RGB-D-based detection algorithms can be slower when compared to their 2D counterparts. In fact, adding an extra spatial dimension increases, relatively, the size of the data. As such, techniques such as sliding windows and convolution operations, which are very efficient on RGB images, become significantly more expensive in terms of computation time and memory storage.
- **Training data.** Despite the wide-spread use of RGB-D sensors, obtaining large *labelled* RGB-D datasets to train detection algorithms is more challenging when compared to obtaining purely RGB datasets. This is due to the price and complexity of RGB-D sensors. Although low cost sensors are currently available, e.g. the Microsoft Kinect, these are usually more efficient in indoor setups. As such, we witnessed a large proliferation of indoor datasets, whereas outdoor datasets are fewer and typically smaller in size.

Intrinsic factors, on the other hand, refer to factors such as deformations, intra-class variations, and inter-class similarities, which are properties of the objects themselves (see Fig. 3.2):

- **Deformations.** Objects can deform in a rigid and non-rigid way. As such, detection algorithms should be invariant to such shape-preserving deformations.
- **Intra-class variations and inter-class similarities.** Object detection algorithms are often required to distinguish between many objects belonging to many classes. Such objects, especially when imaged under uncontrolled settings, display large intra-class variations. Also, natural and man-made objects from different classes may have strong similarities. Such intra-class variations and inter-class similarities can significantly affect the performance of the detection algorithms, especially if the number of RGB-D images used for training is small.

This chapter discusses how the state-of-the-art algorithms addressed some of these challenges.

3.2.2 Taxonomy

Figure 3.3 illustrates the taxonomy that we will follow for reviewing the state-of-the-art techniques. In particular, both traditional (Section 3.3) and deep-learning based (Section 3.4) techniques operate in a pipeline of two or three stages. **The first** stage takes the input RGB-D image(s) and generates a set of region proposals. **The second** stage then refines that selection using some accurate recognition techniques. It also estimates the accurate locations (i.e. centers) of the detected objects, their sizes, and their pose. This is referred to as the object's bounding box. This is usually sufficient for applications such as object recognition and autonomous navigation. Other applications, e.g. object grasping and manipulation, may require an accurate segmentation of the detected objects. This is usually performed either within the second stage of the pipeline or separately with an additional segmentation module, which only takes as input the region within the detected bounding box.

Note that, in most of the state-of-the-art techniques, the different modules of the pipeline operate in an independent manner. For instance, the region proposal module can use traditional techniques based on hand-crafted features, while the recognition and localization module can use deep learning techniques.

Another important point in our taxonomy is the way the input is represented and fed into the pipeline. For instance, some methods treat the depth map as a one-channel image where each pixel encodes depth. The main advantage of this representation is that depth can be processed in the same way as images, i.e. using 2D operations, and thus there is a significant gain in the computation performance and memory requirements. Other techniques use 3D representations by converting the depth map into either a point cloud or a volumetric grid. These methods, however, require 3D operations and thus can be significantly expensive compared to their 2D counterparts.

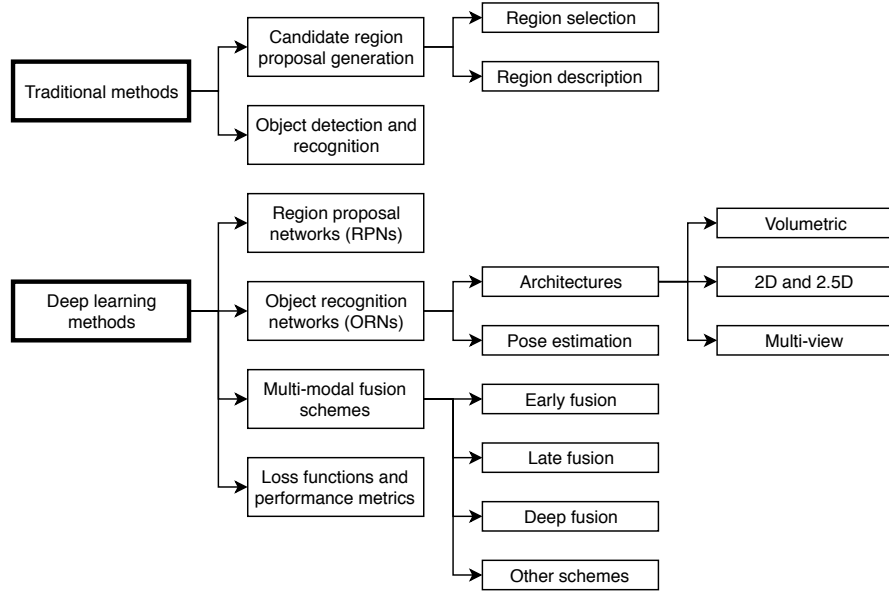


Fig. 3.3: Taxonomy of the state-of-the-art traditional and deep learning methods.

The last important point in our taxonomy is the fusion scheme used to merge multi-modal information. In fact, the RGB and D channels of an RGB-D image carry overlapping as well as complementary information. The RGB image mainly provides information about the colour and texture of objects. The depth map, on the other hand, carries information about the geometry (e.g. size, shape) of objects, although some of this information can also be inferred from the RGB image. Existing state-of-the-art techniques combine this complementary information at different stages of the pipeline (see Fig 3.3).

3.3 Traditional Methods

The first generation of algorithms that aim to detect the location and pose of objects in RGB-D images relies on hand-crafted features¹ combined with machine learning techniques. They operate in two steps: (1) candidate region selection, and (2) detection refinement.

¹ Here we defined hand-crafted or *hand-engineered* features as those which have been calculated over input images using operations which have been explicitly defined by a human designer (i.e. *hand-crafted*), as opposed to learned features which are extracted through optimization procedures in learning pipelines.

3.3.1 Candidate region proposals generation

The first step of an object detection algorithm is to generate a set of candidate regions, also referred to as hypotheses, from image and depth cues. The set will form the potential object candidates and should cover the majority of the true object locations. This step can be seen as a coarse recognition step in which regions are roughly classified into whether they contain the objects of interest or not. While the classification is not required to be accurate, it should achieve a high recall so that object regions will not be missed.

3.3.1.1 Region selection

The initial set of candidate regions can be generated by using (1) bottom-up clustering-and-fitting, (2) sliding windows, or (3) segmentation algorithms.

Bottom-up clustering-and-fitting methods start at the pixel and point level, and iteratively cluster such data into basic geometric primitives such as cuboids. These primitives can be further grouped together to form a set of complex object proposals. In a second stage, an optimal subset is selected using some geometric, structural, and/or semantic cues. Jiang and Xiao [35] constructed a set of cuboid candidates using pairs of superpixels in an RGB-D image. Their method starts by partitioning the depth map, using both colour and surface normals, into superpixels forming piecewise planar patches. Optimal cuboids are then fit to pairs of adjacent planar patches. Cuboids with high fitting accuracy are then considered as potential candidate hypotheses. Representing objects using cuboids has also been used in [8, 39, 51]. Khan et al. [39] used a similar approach to fit cuboids to pairs of adjacent planar patches and also to individual planar patches. The approach also distinguishes between scene bounding cuboids and object cuboids.

These methods tend, in general, to represent a scene with many small components, especially if it contains objects with complex shapes. To overcome this issue, Jiang [34] used approximate convex shapes, arguing that they are better than cuboids in approximating generic objects.

Another approach is to use a **sliding window method** [15, 80, 81] where detection is performed by sliding, or *convolving*, a window over the image. At each window location, the region of the image contained in the window is extracted, and these regions are then classified as potential candidates or not depending on the similarity of that region to the objects of interest. Nakahara et al. [57] extended this process by using multi-scale windows to make the method robust to variations in the size of objects that can be detected. Since the goal of this step is to extract candidate regions, the classifier is not required to be accurate. It is only required to have a high recall to ensure that the selected candidates cover the majority of the true object locations. Thus, these types of methods are generally fast and often generate a large set of candidates.

Finally, **segmentation-based region selection methods** [6, 9] extract candidate regions by segmenting the RGB-D image into meaningful objects and then con-

sidering segmented regions separately. These methods are usually computationally expensive and may suffer in the presence of occlusions and clutter.

3.3.1.2 Region description

Once candidate regions have been identified, the next step is to describe these regions with features that characterize their geometry and appearance. These descriptors can be used to refine the candidate region selection, either by using some supervised recognition techniques, e.g. Support Vector Machines [80], Adaboost [16], and hierarchical cascaded forests [6], or by using unsupervised procedures.

In principle, any type of features which can be computed from the RGB image can be used. Examples include colour statistics, Histogram of Oriented Gradients (HOG) descriptor [13], Scale-Invariant Feature Transform (SIFT) [14], the Chamfer distance [7], and Local Binary Patterns (LBPs) [31]. Some of these descriptors can be used to describe the geometry if computed from the depth map, by treating depth as a grayscale image. Other examples of 3D features include:

- **3D normal features.** These are used to describe the orientation of an object's surface. To compute 3D normals, one can pick n nearest neighbors for each point, and estimate the surface normal at that point using principal component analysis (PCA). This is equivalent to fitting a plane and choosing the normal vector to the surface to be the normal vector to that plane, see [41].
- **Point density features [80].** It is computed by subdividing each 3D cell into $n \times n \times n$ voxels and building a histogram of the number of points in each voxel. Song et al. [80] also applied a 3D Gaussian kernel to assign a weight to each point, canceling the bias of the voxel discretization. After obtaining the histogram inside the cell, Song et al. [80] randomly pick 1000 pairs of entries and compute the difference within each pair, obtaining what is called the stick feature [76]. The stick feature is then concatenated with the original count histogram to form the point density feature. This descriptor captures both the first order (point count) and second order (count difference) statistics of the point cloud [80].
- **Depth statistics.** This can include the first and second order statistics as well as the histogram of depth.
- **Truncated Signed Distance Function (TSDF) [59].** For a region divided into $n \times n \times n$ voxels, the TSDF value of each voxel is defined as the signed distance between the voxel center and the nearest object point on the line of sight from the camera. The distance is clipped to be between -1 and 1 . The sign indicates whether the voxel is in front of or behind a surface, with respect to the camera's line of sight.
- **Global depth contrast (GP-D) [67].** This descriptor measures the saliency of a superpixel by considering its depth contrast with respect to all other superpixels.
- **Local Background Enclosure (LBE) descriptor [22].** This descriptor, which is also used to detect salient objects, is designed based on the observation that salient objects tend to be located in front of their surrounding regions. Thus, the descriptor can be computed by creating patches, via superpixel segmentation [2],

and considering the angular density of the surrounding patches which have a significant depth difference to the center patch (a difference beyond a given threshold). Feng et al. [22] found that LBE features outperform depth-based features such as anisotropic Center-Surround Difference (ACSD) [36], multi-scale depth contrast (LMH-D) [60], and global depth contrast (GP-D) [67], when evaluated on the RGBD1000 [60] and NJUDS2000 [36] RGB-D benchmarks.

- **Cloud of Oriented Gradients (COG) descriptor [70].** It extends the HOG descriptor, which was originally designed for 2D images [10, 24], to 3D data.
- **Histogram of Control Points (HOCP) descriptor [73, 74].** A volumetric descriptor calculated over patches of point clouds featuring occluded objects. The descriptor is derived from the Implicit B-Splines (IBS) feature.

In general, these hand-crafted features are computed at the pixel, super-pixel, point, or patch level. They can also be used to characterize an entire region by aggregating the features computed at different locations on the region using histogram and/or Bag-of-Words techniques, see [41]. For instance, Song et al. [80] aggregate the 3D normal features by dividing, uniformly, the half-sphere into 24 bins along the azimuth and elevation angles. Each bin encodes the frequency of the normal vectors whose orientation falls within that bin. Alternatively, one can use a Bag-of-Words learned from training data to represent each patch using a histogram which encodes the frequency of occurrences of each code-word in the patch [80].

3.3.2 Object detection and recognition

Given a set of candidate objects, one can train, using the features described in Section 3.3.1, a classifier that takes these candidates and classifies them into either true or false detections. Different types of classifiers have been used in the literature. Song and Xiao [80] use Support Vector Machines. Asif et al. [6] use hierarchical cascade forests. In general, any type of classifiers, e.g. AdaBoost, can be used to complete this task. However, in many cases, this approach is not sufficient and may lead to excessive false positives and/or false negatives. Instead, given a set of candidates, other approaches select a subset of shapes that best describes the RGB-D data whilst satisfying some geometrical, structural and semantic constraints. This has been formulated as the problem of optimizing an objective function of the form [34, 35]:

$$E(x) = E_d(x) + E_r(x), \quad (3.1)$$

where x is the indicator vector of the candidate shapes, i.e. $x(i) = 1$ means that the i -th shape is selected for the final subset. Khan et al. [39] extended this formulation to classify RGB-D superpixels as cluttered or non-cluttered, in addition to object detection. As such, another variable y is introduced. It is a binary indicator vector whose i -th entry indicates whether the i -th superpixel is cluttered or not. This has also been formulated as the problem of optimizing an objective function of the form:

$$U(x, y) = E(x) + U(y) + U_c(x, y). \quad (3.2)$$

Here, $E(x)$ is given by Equation (3.1). $U(y)$ also consists of two potentials:

$$U(y) = U_d(y) + U_r(y), \quad (3.3)$$

where the unary term U_d is the data likelihood of a superpixel's label, and the pairwise potential U_r encodes the spatial smoothness between neighboring superpixels. The third term of Equation (3.2) encodes compatibility constraints, e.g. enforcing the consistency of the cuboid labelling and the superpixel labelling.

The data likelihood term E_d . This term measures the cost of matching the candidate shape to the data. For instance, Jiang et al. [35] used cuboids for detection, and defined this term as the cost of matching a cuboid to the candidate shape. On the other hand, Jiang [34] used convex shapes, and defined E_d as a measure of concavity of the candidate shape. Several other papers take a learning approach. For instance, Khan et al. [39] computed seven types of cuboid features (volumetric occupancy, colour consistency, normal consistency, tightness feature, support feature, geometric plausibility feature, and cuboid size feature), and then predicted the local matching quality using machine learning approaches.

The data likelihood term U_d . Khan et al. [39] used a unary potential that captures, on each superpixel, the appearance and texture properties of cluttered and non-cluttered regions. This is done by extracting several cues including image and depth gradient, colour, surface normals, LBP features, and self-similarity features. Then, a Random Forest classifier was trained to predict the probability of a region being a clutter or a non-clutter.

The regularization terms E_r and U_r . The second terms of Equations (3.1) and (3.3) are regularization terms, which incorporate various types of constraints. Jian [34] define this term as:

$$E_r(x) = \alpha N(x) + \beta I(x) - \lambda A(x), \quad (3.4)$$

where α, β , and λ are weights that set the importance of each term. $N(x)$ is the number of selected candidates, $I(x)$ measures the amount of intersection (or overlap) between two neighboring candidate shapes, and $A(x)$ is the amount of area covered by the candidate shapes projected onto the image plane. Jiang and Xiao [35], and later Khan et al. [39], used the same formulation but added a fourth term, $O(x)$, which penalizes occlusions.

For the superpixel pairwise term of Equation (3.3), Khan et al. [39] defined a contrast-sensitive Potts model on spatially neighboring superpixels, which encouraged the smoothness of cluttered and non-cluttered regions.

The compatibility term U_c . Khan et al. [39] introduced the compatibility term to link the labelling of the superpixels to the cuboid selection task. This ensures consistency between the lower level and the higher level of the scene representation. It consists of two terms: a superpixel membership potential, and a superpixel-cuboid

occlusion potential. The former ensures that a superpixel is associated with at least one cuboid if it is not a cluttered region. The latter ensures that a cuboid should not appear in front of a superpixel which is classified as clutter, i.e. a detected cuboid cannot completely occlude a superpixel on the 2D plane which takes a clutter label.

Optimization. The final step of the process is to solve the optimization problem of Equations (3.1) or (3.3). Jian [34] showed that the energy function of Equation (3.1) can be linearized and optimized using efficient algorithms such as the branch-and-bound method. Khan et al. [39], on the other hand, transformed the minimization problem into a Mixed Integer Linear Program (MILP) with linear constraints, which can be solved using the branch-and-bound method.

3.4 Deep Learning Methods

Despite the extensive research, the performance of traditional methods is still far from the performance of the human visual system, especially when it comes to detecting objects in challenging situations, e.g. highly-cluttered scenes and scenes with high occlusions. In fact, while traditional methods perform well in detecting and producing bounding boxes on visible parts, it is often desirable to capture the full extent of the objects regardless of occlusions and clutter. Deep learning-based techniques aim to overcome these limitations. They generally operate following the same pipeline as the traditional techniques (see Section 3.3), i.e. region proposals extraction, object recognition, and 3D bounding box location and pose estimation. However, they replace some or all of these building blocks with deep learning networks. This section reviews the different deep learning architectures that have been proposed to solve these problems. Note that these techniques can be combined with traditional techniques; e.g. one can use traditional techniques for region proposals and deep learning networks for object recognition, bounding box location and pose refinement.

3.4.1 Region proposal networks

In this section, we are interested in amodal detection, i.e. the inference of the full 3D bounding box beyond the visible parts. This critical step in an object detection pipeline is very challenging since different object categories can have very different object sizes in 3D. Region Proposal Networks (RPNs) are central to this task since they reduce the search space considered by the remainder of the object detection pipeline.

We classify existing RPNs into three different categories. **Methods in the first category** perform the detection on the RGB image and then, using the known camera projection matrix, lift the 2D region to a 3D frustum that defines a 3D search

space for the object. In general, any 2D object detector, e.g. [53], can be used for this task. However, using deep CNNs allows the extraction of rich features at varying degrees of complexity and dimensionality, which is beneficial for the purpose of RGB-D object detection tasks. In fact, the success of the recent object detection pipelines can be largely attributed to the automatic feature learning aspect of convolutional neural networks. For instance, Qi et al. [62] used the Feature Pyramid Networks (FPN) [52], which operate on RGB images, to detect region proposals. Lahoud and Ghanem [45], on the other hand, use the 2D Faster R-CNN [69] and VGG-16 [77], pre-trained on the 2D ImageNet database [72], to position 2D bounding boxes around possible objects with high accuracy and efficiency. These methods have been applied for indoor and outdoor scenes captured by RGB-D cameras, and for scenes captured using LIDAR sensors [12].

The rationale behind these methods is that the resolution of data produced by most 3D sensors is still lower than the resolution of RGB images, and that 2D object detectors are mature and quite efficient. RGB-based detection methods, however, do not benefit from the additional information encoded in the depth map.

The second class of methods aim to address these limitations. They treat depth as an image and perform the detection of region proposals on the RGB-D image either by using traditional techniques or by using 2D convolutional networks. For instance, Gupta et al. [29], and later Deng and Latecki [15], computed an improved contour image from an input RGB-D image. An improved contour image is defined as the contour image but augmented with additional features such as the gradient, normals, the geocentric pose, and appearance features such as the soft edge map produced by running the contour detector on the RGB image. They then generalize the multiscale combinatorial grouping (MCG) algorithm [5, 61] to RGB-D images for region proposal and ranking. Note that both the hand-crafted features as well as the region recognition and ranking algorithms can be replaced with deep learning techniques, as in [62].

Chen et al. [12] took the point cloud (produced by LIDAR sensors) and the RGB image, and produced two types of feature maps: the bird's eye view features and the front view features. The bird's eye view representation is encoded by height, intensity, and density. First, the point cloud is projected and discretized into a 2D grid with a fixed resolution. For each cell in the grid, the height feature is computed as the maximum height of the points in that cell. To encode more detailed height information, the point cloud is divided equally into m slices. A height map is computed for each slice, thus obtaining m height maps. The intensity feature is defined as the reflectance value of the point which has the highest height in each cell. Finally, a network that is similar to the region proposal network of [68] was used to generate region proposals from the bird's eye view map.

These methods require the fusion of the RGB and depth data. This can be done by simply concatenating the depth data with the RGB data and using this as input to the RPN. However, depth data encodes geometric information, which is distinct from the spatial and colour information provided by monocular RGB images. As such, Alexandre et al. [3] found that fusing amodal networks with a majority voting scheme produced better results in object recognition tasks, with an improvement of

29% when compared to using simple RGB and depth frame concatenation. Note that, instead of performing fusion at the very early stage, e.g. by concatenating the input modalities, fusion can be performed at a later stage by concatenating features computed from the RGB and D maps, or progressively using the complementarity-aware fusion network of Chen et al. [11], see Section 3.4.3.

The third class of methods take a 3D approach. For instance, Song and Xia [81] projected both the depth map and the RGB image into the 3D space forming a volumetric scene. The 3D scene is then processed with a fully 3D convolutional network, called a *3D Amodal Region Proposal Network*, which generates region proposals in the form of 3D bounding boxes at two different scales. Multi-scale RPNs allow the detection of objects of different sizes. It performs a 3D sliding-window search with varying window sizes, and produces an objectness score for each of the non-empty proposal boxes [4]. Finally, redundant proposals are removed using non-maximum suppression with an IoU threshold of 0.35 in 3D. Also, the approach ranks the proposals based on their objectness score and only selects the top 2000 boxes to be used as input to the object recognition network.

3D detection can be very expensive since it involves 3D convolutional operations. In fact, it can be more than 30 times slower than its 2D counterpart. Also, the solution space is very large since it includes three dimensions for the location and two dimensions for the orientation of the bounding boxes. However, 3D voxel grids produced from depth maps are generally sparse as they only contain information near the shape surfaces. To leverage this sparsity, Engelcke et al. [18] extended the approach of Song et al. [80] by replacing the SVM ensemble with a 3D CNN, which operates on voxelized 3D grids. The key advantage of this approach is that it leverages the sparsity encountered in point clouds to prevent huge computational cost that occurs with 3D CNNs. In this approach, the computational cost is only proportional to the number of occupied grid cells rather than to the total number of cells in the discretized 3D grid as in [84].

3.4.2 Object recognition networks

Once region proposals have been generated, the next step is to classify these regions into whether they correspond to the objects we want to detect or not, and subsequently refine the detection by estimating the accurate location, extent, and pose (position and orientation) of each object's bounding box. The former is a classification problem, which has been well solved using Object Recognition Networks (ORNs) [6, 17, 75]. An ORN takes a candidate region, and assigns to it a class label, which can be binary, i.e. 1 or 0, to indicate whether it is an object of interest or not, or multi-label where the network recognizes the class of the detected objects.

There are several ORN architectures that have been proposed in the literature [12, 15, 28, 55, 81]. In this section, we will discuss some of them based on (1) whether they operate on 2D or 3D (and subsequently whether they are using 2D or

3D convolutions), and (2) how the accurate 3D location and size of the bounding boxes are estimated.

3.4.2.1 Network architectures for object recognition

(1) Volumetric approaches. The first class of methods are volumetric. The idea is to lift the information in the detected regions into 3D volumes and process them using 3D convolutional networks. For instance, Maturana et al. [55] used only the depth information to recognize and accurately detect the objects of interest. Their method first converts the point cloud within each 3D region of interest into a $32 \times 32 \times 32$ occupancy grid, with the z axis approximately aligned with gravity. The point cloud is then fed into a 3D convolutional network, termed *VoxNet*, which outputs the class label of the region.

Song and Xia [81] followed a similar volumetric approach but they jointly learned the object categories and the 3D box regression from both depth and colour information. Their approach operates as follows. For each 3D proposal, the 3D volume from depth is fed to a 3D ConvNet, and the 2D colour patch (the 2D projection of the 3D proposal) is fed to a 2D ConvNet (based on VGG and pre-trained on ImageNet). The two latent representations learned by the two networks are then concatenated and further processed with one fully connected layer. The network then splits into two branches, each composed of one fully connected layer. The first branch is a classification branch as it produces the class label. The second branch estimates the location and size of the 3D amodal bounding box of the detected object. This approach has two important features; **first**, it combines both colour and geometry (through depth) information to perform recognition and regress the amodal bounding box. These two types of information are complementary and thus combining them can improve performance. The **second** important feature is that it does not directly estimate the location and size of the bounding box but instead it estimates the residual. That is, it first takes an initial estimate of the size of the bounding box (using some prior knowledge about the class of shapes of interest). The network is then trained to learn the correction that one needs to apply to the initial estimate in order to obtain an accurate location and size of the bounding box.

(2) 2D and 2.5D approaches. In the context of object detection, 2D approaches operate over the two spatial dimensions in the input (i.e. an RGB image), without exploiting the data encoded in depth images. 2.5D inputs refer to inputs with attached depth images, but importantly, these depth images are treated similarly to how colour images are (without exploiting 3D spatial relationships, i.e. using the depth frames as 2D maps where each pixel encodes the depth value). Finally, 3D approaches use the rich spatial data encoded in volumetric or point cloud representations of data (or any other representation which represents the data over three spatial dimensions).

These 2D and 2.5D approaches are mainly motivated by the performance of the human visual system in detecting and recognizing objects just from partial 2D infor-

mation. In fact, when the majority of an object area on the depth map is not visible, the depth map only carries partial information. However, information encoded in the 2D image is rich, and humans can still perceive the objects and estimate their 3D locations and sizes from such images [15]. 2D and 2.5D approaches try to mimic the human perception and leverage the 2.5D image features directly using current deep learning techniques.

In particular, Deng and Latecki [15] followed the same approach as Song and Xia [81] but operate on 2D maps using 2D convolutional filters. The main idea is to regress the 3D bounding box just from the RGB and depth map of the detected 2D regions of interests. Their approach replaces the 3D ConvNet of Song and Xia [81] with a 2D ConvNet that processes the depth map. Thus, it is computationally more efficient than the approaches which operate on 3D volumes, e.g. [15].

(3) Multi-view approaches. The 2D and 2.5D approaches described above can be extended to operate on multi-view inputs. In fact, many practical systems, e.g. autonomous driving, acquire RGB-D data from multiple view points. Central to multi-view techniques is the fusion mechanism used to aggregate information from different views (see also Section 3.4.3), which can be multiple images and/or depth maps captured from multiple view points. Some of the challenges which need to be addressed include catering for images gathered at varying resolutions.

Chen et al. [12] proposed a Multi-View 3D network (MV3D), a region-based fusion network, which combines features from multiple views. The network jointly classifies region proposals and regresses 3D bounding box orientations. The pipeline operates in two stages: multi-view ROI pooling, and deep fusion. The former is used to obtain feature vectors of the same length, since features from different views/modalities usually have different resolutions. The deep fusion network fuses multi-view features hierarchically to enable more interactions among features of the intermediate layers from different views.

3.4.2.2 Pose estimation

One of the main challenges in amodal object detection from RGB-D images is how to accurately estimate the pose of the bounding box of the detected object, even if parts of the objects are occluded. Early works such as Song and Xia [81] do not estimate orientation but use the major directions of the room in order to orient all proposals. This simple heuristic works fine for indoor scenes, e.g. rooms. However, it cannot be easily extended to outdoor scenes or scenes where no prior knowledge of their structure is known.

Another approach is to perform an exhaustive search of the best orientations over the discretized space of all possible orientations. For example, Maturana et al. [55] performed an exhaustive search over $n = 12$ orientations around the z axis and selected the one with the largest activation. At training time, Maturana et al. [55] augmented the dataset by creating $n = 12$ to $n = 18$ copies of each input instance, each rotated $360^\circ/n$ intervals around the z axis (assuming that the z axis is known).

At testing time, the activations of the output layer over all the n copies are aggregated by pooling. This approach, which can be seen as a voting scheme, has been also used to detect landing zones from LIDAR data [54].

Gupta et al. [28] considered the problem of fitting a complete 3D object model to the detected objects, instead of just estimating the location, orientation, and size of the object’s bounding box. They first detect and segment object instances in the scene and then use a convolutional neural network (CNN) to predict the coarse pose of the object. Gupta et al. [28] then use the detected region (segmentation mask) to create a 3D representation of the object by projecting points from the depth map. The Iterative Closest Point (ICP) algorithm [71] is then used to align 3D CAD models to these 3D points.

Finally, some recent approaches regress pose in the same way as they perform recognition, i.e. using CNNs. This is usually achieved using a region recognition network, which has two branches of fully connected layers; one for recognition and another one for bounding box regression [12, 15]. Existing methods differ in the way the bounding boxes are parameterized. For instance, Cheng et al. [12] represent a bounding box using its eight corners. This is a redundant representation as a cuboid can be described with less information. Deng and Latecki [15] used a seven-entry vector $[x_{cam}, y_{cam}, z_{cam}, l, w, h, \theta]$ where $[x_{cam}, y_{cam}, z_{cam}]$ corresponds to the coordinates of the bounding box’s centroid under the camera coordinate system. $[l, w, h]$ represents its 3D size, and θ is the angle between the principal axis and its orientation vector under the tilt coordinate system. Note that these methods do not directly regress the pose of the bounding box. Instead, starting from an initial estimate provided by the Region Proposal Network, the regression network estimates the offset vector, which is then applied to the initial estimate to obtain the final pose of the bounding box.

3.4.3 Fusion schemes

In the context of RGB-D object detection, we aim to exploit the multiple modalities that are present in RGB-D images, which carry complementary information (colour and depth data). This, however, requires efficient fusion mechanisms. In this section, we discuss some of the strategies that have been used in the literature.

(1) Early fusion. In this scheme, the RGB image and the depth map are concatenated to form a four-channel image [30]. This happens at the earliest point in the network, i.e. before any major computational layers process the image. The concatenated image is then processed with 2D or 3D convolutional filters. This scheme was adopted in [32] for the purpose of saliency detection.

(2) Late fusion. In this scheme, the RGB image and the depth map are processed separately, e.g. using two different networks, to produce various types of features. These features are then fused together, either by concatenation or by further processing using convolutional networks. Eitel et al. [17], for example, used two networks,

one for depth and one for the RGB image, with each network separately trained on ImageNet [40]. The feature maps output by the two networks are then concatenated and presented to a final fusion network, which produces object class predictions. This approach achieved an overall accuracy of $91.3\% \pm 1.4\%$ on the Washington RGB-D Object Dataset, see Table 3.3 for more details regarding pipeline performance.

Note that, Chen et al. [12] showed that early and late fusion approaches perform similarly when tested on the hard category of the KITTI dataset, scoring an average precision of 87.23% and 86.88%.

(3) Deep fusion. Early and late fusion schemes are limited in that they only allow the final joint predictions to operate on early or deep representations, so useful information can be discarded. Chen et al. [12] introduced a deep learning fusion scheme, which fuses features extracted from multiple representations of the input. The fusion pipeline uses element-wise mean pooling operations instead of simple concatenations (as in early or late fusion). Chen et al. [12] showed that this fusion mechanism improved performance by about 1% compared to early and late fusion.

(4) Sharable features learning and complementarity-aware fusion. The fusion methods described above either learn features from colour and depth modalities separately, or simply treat RGB-D as a four-channel data. Wang et al. [83] speculate that different modalities should contain not only some modal-specific patterns but also some shared common patterns. They then propose a multi-modal feature learning framework for RGB-D object recognition. First, two deep CNN layers are constructed, one for colour and another for depth. They are then connected with multimodal layers, which fuse colour and depth information by enforcing a common part to be shared by features of different modalities. This produces features reflecting shared properties as well as modal-specific properties from different modalities.

Cheng et al. [11] proposed a fusion mechanism, termed *complementarity-aware (CA) fusion*, which encourages the determination of complementary information from the different modalities at different abstraction levels. They introduced a CA-Fuse module, which enables cross-modal, cross-level connections and modal/level-wise supervisions, explicitly encouraging the capture of complementary information from the counterpart, thus reducing fusion ambiguity and increasing fusion efficiency.

3.4.4 Loss functions

In general, the region proposal network (RPN) and the object recognition network (ORN) operate separately. The RPN first detects candidate regions. The ORN then refines the detection by discarding regions that do not correspond to the objects of interest. The ORN then further refines the location, size, and orientation of the bounding boxes. As such, most of the state-of-the-art techniques train these networks separately using separate loss functions.

Loss functions inform the network on how poorly it completed its designated task over each training batch using a scalar metric (referred to as the loss, cost, or inverse fitness). The calculation of the loss should incorporate error that the algorithm accumulated during the completion of its task, as the network will change its weights in order to reduce the loss, and thus the error. For example, in object classification networks, the loss might be defined as the mean squared error between the one-hot encoded ground truth labels, and the network's output logits. For object detection networks, the IoU (see Fig. 3.4) of the detected region and the ground truth region may be incorporated. In this way the loss function design is task-dependant, and performance increases have been observed to be contingent on the loss function's design [85]. Typically loss functions are hand-crafted, though weightings between terms can be learned [37]. Naturally, numerous loss functions have been devised to train networks to accomplish various tasks [38].

A common loss function that has been used for classification (in the RPN as well as in the ORN) is the softmax regression loss. Let θ be the parameters of the network, m the number of training samples, n_c the number of classes, and $y_i \in \{1, \dots, n_c\}$ the output of the network for the training sample x_i . The softmax regression loss is defined as:

$$L(\theta) = - \sum_{i=1}^m \sum_{k=1}^{n_c} \mathbf{1}(y_i = k) \log(p(y^i = k | x^i; \theta)). \quad (3.5)$$

Here, $\mathbf{1}(s)$ is equal to 1 if the statement s is true and 0 otherwise. This loss function has been used by Gupta et al. [28] to train their region proposal network, which also provides a coarse estimation of each object's pose.

Song and Xia [81], on the other hand, trained their multiscale region proposal network using a loss function that is a weighted sum of two terms: an objectness term and a box regression term:

$$L(p, p^*, \mathbf{t}, \mathbf{t}^*) = L_{cls}(p, p^*) + \lambda p L_{l_{smooth}}(\mathbf{t}, \mathbf{t}^*), \quad (3.6)$$

where p^* is the predicted probability of this region being an object and p is the ground truth, L_{cls} is the log loss over the two categories (object vs. non-object) [26], and \mathbf{t} is a 6-element vector, which defines the location and scale of the bounding box of the region. The second term, which is the box regression loss term, is defined using the smooth L_1 function as follows:

$$L_{l_{smooth}}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1, \\ |x| - 0.5, & \text{otherwise.} \end{cases} \quad (3.7)$$

Song and Xia [81] also used a similar loss function to train their ORN. The only difference is in the second term, which is set to zero when the ground-truth probability p is zero. This is because there is no notion of a ground-truth bounding box for background RoIs. Hence, the box regression term is ignored. Finally, Maturana et al. [55]

used the multinomial negative log-likelihood plus 0.001 times the L_2 weight norm for regularization.

3.5 Discussion and comparison of some pipelines

In this section, we discuss some pipelines for object detection from RGB-D data and compare their performance on standard benchmarks. We will first review examples of the datasets that have been used for training and testing the techniques (Section 3.5.1), discuss different performance evaluation metrics (Section 3.5.1), and finally compare and discuss the performance of some of the key RGB-D-based object detection pipelines (Section 3.5.3).

3.5.1 Datasets

Many of the state-of-the-art algorithms rely on large datasets to train their models and evaluate their performance. Both traditional machine learning and advanced deep learning approaches require labelled datasets in the form of RGB-D images and their corresponding ground-truth labels. The labels can be in the form of 2D bounding boxes highlighting the object regions in the RGB image and/or the depth map, oriented 3D bounding boxes (3DBBX) delineating the 3D regions of the objects, and/or exact segmentations (in the form of segmentation masks) of the objects of interest.

Table 3.1 summarizes the main datasets and benchmarks that are currently available in the literature. Note that several types of 3D sensors have been used for the acquisition of these datasets. For instance, the SUN RGB-D dataset was constructed using four different sensors: the Intel RealSense 3D Camera, the Asus Xtion LIVE PRO, and Microsoft Kinect v1 and v2. Intel Asus and Microsoft Kinect v1 sensors use infrared (IR) light patterns to generate quantized depth maps, an approach known as *structured light*, whereas Microsoft Kinect v2 uses time-of-flight ranging. These sensors are suitable for indoor scenes since their depth range is limited to a few meters. On the other hand, The KITTI dataset, which includes outdoor scene categories, has been captured using a Velodyne HDL-64E rotation 3D laser scanner.

Note that some datasets, such as the PASCAL3D+, are particularly suitable for testing the robustness of various algorithms to occlusions, since an emphasis was placed on gathering data with occlusions.

Table 3.1: Examples of datasets used for training and evaluating object detection pipelines from RGB-D images.

Name	Description	Size
KITTI 2015 [56]	Cluttered driving scenarios recorded in and around Karlsruhe in Germany.	400 annotated dynamic scenes from the raw KITTI dataset.
KITTI 2012 [23]	As above.	389 image pairs and more than 200,000 3D object annotations.
SUN RGB-D [79]	Indoor houses and universities in North America and Asia.	10,335 images, 800 object categories, and 47 scene categories annotated with 58,657 bounding boxes (3D).
NYUDv2 [58]	Diverse indoor scenes taken from three cities.	1,449 RGB-D images over 464 scenes.
PASCAL3D+ [89]	Vehicular and indoor objects (augments the PASCAL VOC dataset [19]).	12 object categories with 3,000 instances per category.
ObjectNet3D [88]	Indoor and outdoor scenes.	90,127 images sorted into 100 categories. 201,888 objects in these images and 44,147 3D shapes.
RGBD1000 [60]	Indoor and outdoor scenes captured with a Microsoft Kinect.	1,000 RGB-D images.
NJUDS2000 [36]	Indoor and outdoor scenes.	2,000 RGB-D images.
LFSD [50]	Indoor (60) and outdoor (40) scenes captured with a Lytro camera.	100 light fields each composed from raw light field data, a focal slice, an all-focus image and a rough depth map.
Cornell Grasping [1]	Several images and point clouds of typical graspable indoor objects taken at different poses.	1,035 images of 280 different objects.
ModelNet10 [87]	Object aligned 3D CAD models for the 10 most common object categories found in the SUN2012 database [90].	9,798 total instances split amongst 10 object categories, each with their own test/train split.
ModelNet40 [87]	Object aligned 3D CAD models for 40 common household objects.	12,311 total instances split amongst 40 object categories, each with their own test/train split.
Caltech-101 [20]	Single class object centric images with little or no clutter. Most objects are presented in a stereotypical pose.	9,144 images sorted into 101 categories, with 40 to 800 images per category.
Caltech-256 [27]	Single class images with some clutter.	30,607 images sorted into 256 categories with an average of 119 images per category.
Washington RGB-D [46]	Turntable video sequences at varying heights of common household objects.	300 objects organized into 51 categories. Three video sequences per object.

3.5.2 Performance criteria and metrics

Object detection usually involves two tasks; the first is to assess whether the object exists in the RGB-D image (classification) and the second is to exactly localize the object in the image (localization). Various metrics have been proposed to evaluate

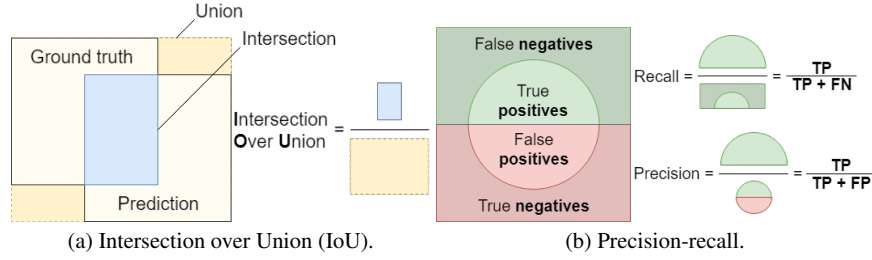


Fig. 3.4: (a) Illustration of the Intersection over Union (IoU) metric in 2D. (b) Illustration of how precision and recall are calculated from a model's test results. Recall is the ratio of correct predictions to total objects in the dataset. It measures how complete the predictions are. Precision is the ratio of correct predictions to total predictions made, i.e. how correct the predictions are.

the performance of these tasks. Below, we discuss the most commonly used ones, see also [42].

Computation time. Object detection algorithms operate in two phases; the training phase and the testing phase. While, in general, algorithms can afford having a slow training phase, the computation time at runtime is a very important performance criterion. Various applications may have different requirements. For instance, time-critical applications such as autonomous driving and surveillance systems should operate in real time. Other applications, e.g. offline indexing of RGB-D images and videos, can afford slower detection times. However, given the large amount of information they generally need to process, real-time detection is desirable. Note that, there is often a trade-off between computation time at runtime and performance.

Intersection over Union (IoU). It measures the overlap between a ground-truth label and a prediction as a proportion of the union of the two regions, see Fig. 3.4. IoU is a useful metric for measuring the predictive power of a 2D/3D object detector. IoU thresholds are applied to sets of detections in order to precisely define what constitutes a positive detection. For example, an $\text{IoU} > 0.5$ might be referred to as a positive detection. Such thresholds are referred to as *overlap criterion*.

Precision-recall curves. Precision and recall are calculated based on a model's test results, see Fig. 3.4. The precision-recall curve is generated by varying the threshold which determines what is counted as a positive detection of the class. The model's precision at varying recall values are then plotted to produce the curve.

Average Precision (AP). It is defined as the average value of precision over the interval from recall $r = 0$ to $r = 1$, which is equivalent to measuring the area under the precision-recall curve (r here is the recall):

$$AP = \int_0^1 \text{precision}(r) dr. \quad (3.8)$$

Mean Average Precision (mAP) score. It is defined as the mean average precision over all classes and/or over all IoU thresholds.

F- and E-Measures. These are two measures which combine precision and recall into a single number to evaluate the retrieval performance. The F-measure is the weighted harmonic mean of precision and recall. It is defined as

$$F_{\alpha} = \frac{(1 + \alpha) \times \text{precision} \times \text{recall}}{\alpha \times \text{precision} + \text{recall}}, \quad (3.9)$$

where α is a weight. When $\alpha = 1$ then

$$F_1 \equiv F = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}. \quad (3.10)$$

The E-Measure is defined as $E = 1 - F$, which is equivalent to

$$E = 2 \left(\frac{1}{\text{precision}} + \frac{1}{\text{recall}} \right)^{-1}. \quad (3.11)$$

Note that the maximum value of the E-measure is 1.0 and the higher the E-measure is, the better is the detection algorithm. The main property of the E-measure is that it quantifies how good are the results retrieved in the top of the ranked list. This is very important since, in general, the user of a search engine is more interested in the first page of the query results than in the later pages.

Localization performance. The localization task is typically evaluated using the Intersection over Union threshold (IoU) as discussed above.

3.5.3 Discussion and performance comparison

Tables 3.2 and 3.3 summarize the performance, on various datasets and using various performance metrics, of some traditional and deep learning-based RGB-D object detection pipelines. Below, we discuss the pipelines whose performances are highlighted in bold in Table 3.3, with a focus on analyzing the operational steps which allow them to offer increased performance.

The Deep Sliding Shapes model for amodal 3D object detection in RGB-D images [81] extends its predecessor [80], which used hand-crafted features and SVMs. The network begins by replacing the 3D exhaustive search with a 3D multi-scale RPN, which produces 3D RoIs. Non-maximum suppression with an IoU constraint of less than 0.35 is enforced on the RPN output to reduce the number of RoIs. Each RoI is then projected to 2D and fed to a VGG-19-based deep feature extractor [77], which produces the class labels as well as the 3D bounding boxes of the detected objects.

Table 3.2: Performance of some traditional RGB-D-based object detection methods across multiple datasets and performance metrics. All values taken directly from the cited publications. *Class* refers to classification accuracy, *Reco* to recognition accuracy and *mAP* to mean Average Precision. All metrics are measured in percentages except F-score, which is in the range $[0, 1]$.

Method	Dataset	Detection rate		Recall	Class.	Reco.	mAP	
		IoU>.7	IoU>.75	IoU>.5			IoU>.25	F-score
Linear Cuboid Matching [35]	NYUDv2	75.0						
Extended CPMC [51]	NYUDv2			42.8	60.5			0.36
Convex Shapes [34]	NYUDv2		78.2					
Sliding Shapes [80]	NYUDv2 ¹						39.6	
DSS [81]	NYUDv2 ¹						72.3	
DSS [81]	NYUDv2 ²						36.3	
Amodal Detection [15]	NYUDv2 ²						40.9	
DSS [81]	SUN RGB-D						29.6	
CoG [70]	SUN RGB-D ₃						47.6	
Separating objects / clutter [39]	RMRC 2013	38.0						
Sliding Shapes [80]	RMRC 2013						62.4	
M-HMP [9]	Caltech Bird-200					30.3		
M-HMP [9]	Caltech-256				58.0			
M-HMP [9]	MIT Scene-67				51.2			
STEM-CaRFs [6]	Washington					97.6		
STEM-CaRFs [6]	Cornell grasping					94.1		
GP [67]	NLPR							0.72
GP [67]	NJUDS400							0.76
LBE [22]	NJUDS2000							0.72
LBE [22]	RGBD1000							0.73
3DHoG [10]	i-LIDS ⁴				92.1			

¹ Considering only 5 categories from the NYUDv2 dataset. ² 3D annotations for the NYUDv2 dataset were improved in [15] and this improved dataset was used to calculate performance. ³ Considering only 10 categories from the SUN RGB-D dataset. ⁴ Scenario 1 of the i-LIDS dataset.

Sun et al. [82] proposed an object detection framework for a mobile manipulation platform. The framework is composed of an RPN, an ORN, and a Scene Recognition Network (SRN). Its main feature is that the convolutional operations of the three modules are shared, subsequently reducing the computational cost. The ORN, which achieved a mean average precision (mAP) of 52.4% on the SUN RGB-D dataset, outperformed Faster R-CNN [69], RGB-D RCNN [29], and DPM [21].

Another example is the object detection pipeline of Qi et al. [62], which produces the full extents of an object's bounding box in 3D from RGB-D images by using four sub-networks, namely:

Table 3.3: Performance of some deep learning-based object detection methods across multiple datasets and performance metrics. All values taken directly from the original publications. *mAP* refers to mean Average Precision, *Class* refers to classification accuracy, *Valid* to validation set, and *Seg* to instance segmentation. All metrics are measured in percentages except F-score, which is in the range $[0, 1]$.

Method	Dataset	Detection mAP			mAP ¹	Seg.	Class.	F-score $\alpha^2 = 0.3$
		IoU>.25	IoU>.5	IoU>.7				
Rich Features [29]	NYUDv2					37.3		
RGB-D DPM [29]	NYUDv2					23.9		
RGB R-CNN [29]	NYUDv2					22.5		
Amodal 3D [15]	NYUDv2	40.9						
DSS [81]	NYUDv2	36.3						
VoxNet [55]	NYUDv2						71.0	
ShapeNet [87]	NYUDv2						58.0	
Frustrum PointNets [62]	SUN RGB-D ²	54.0						
DSS [81]	SUN RGB-D ²	42.1						
COG [70]	SUN RGB-D ²	47.6						
2D-driven [45]	SUN RGB-D ²	45.1						
Faster R-CNN [69]	SUN RGB-D		50.8					
Unified CNN [82]	SUN RGB-D		52.4					
Frustrum PointNets [62]	KITTI (valid) ³			63.7				
MV3D [12]	KITTI (valid) ³			55.1				
MV3D [12]	KITTI (test) ³				79.8			
3D FCN [48]	KITTI (test) ³				68.3			
Vote3Deep [18]	KITTI (test) ³				63.2			
Vote3D [84]	KITTI (test) ³				42.6			
VeloFCN [49]	KITTI (test) ³				46.9			
Complement-Aware [11]	NLPR							0.85
Salient Deep Fusion [66]	NLPR							0.78
LMH [60]	NLPR							0.65
GP [67]	NLPR							0.72
ACSD [36]	NLPR							0.54
PointNet [63]	ModelNet40						89.2	
PointNet++ [65]	ModelNet40						91.9	
VoxNet [65]	ModelNet40						85.9	
3D ShapeNets [86]	ModelNet40						84.7	
Subvolumes [64]	ModelNet40						89.2	
VoxNet [55]	ModelNet40						83.0	
ShapeNet [87]	ModelNet40						77.0	
Faster R-CNN [69] ^{4,5}	PASCAL 2012 ⁶				76.4			
ION [30]	PASCAL 2012 ⁶				74.6			
SSD512 [30]	PASCAL 2012 ⁶				76.8			
Short Conns [32]	PASCAL 2010 ⁶							0.82
MMDL [17]	Washington						91.3	
STEM-CaRFs [6]	Washington						88.1	
CNN Features [75]	Washington						89.4	
VoxNet [55]	ModelNet10						92.0	
ShapeNet [87]	ModelNet10						84.0	
Faster R-CNN on FPN [52]	COCO test-dev				36.2			
AttractionNet [25]	COCO test-dev				35.7			
Complement-Aware [11]	NJUDS							0.86
Short Conns [32]	MSRA-B							0.92
MMDL [17]	Scenes						82.1	

¹ Values in this column use unspecified IoU thresholds. ² Using only 10 categories from the SUN RGB-D dataset. ³ Using the hard, cars subset of the KITTI dataset. ⁴ Results taken from [30]. ⁵ Uses ResNet-101. ⁶ Using the PASCAL Visual Object Classes (VOC) evaluation.

- **A joint 2D RPN/ORN.** It generates 2D region proposals from the RGB image, and classifies them into one of the n_c object categories.
- **A PointNet-based network.** It performs 3D instance segmentation of the point clouds within 3D frustums extended from the proposed regions.
- **A light-weight regression PointNet (T-Net).** It estimates the true center of the complete object and then transforms the coordinates such that the predicted center becomes the origin.
- **A box estimation network.** It predicts, for each object, its amodal bounding box for the entire object even if parts of it are occluded.

The approach simultaneously trains the 3D instance segmentation PointNet, the T-Net, and the amodal box estimation PointNet, using a loss function that is defined as a weighted sum of the losses of the individual subnetworks. Note that Qi et al. [62]’s network architecture is similar to the architecture of the object classification network of [63, 65] but it outputs the object class scores as well as the detected object’s bounding box. The work also shares some similarities to [45], which used hand-crafted features.

3.6 Summary and perspectives

In this chapter, we have reviewed some of the recent advances in object detection from RGB-D images. Initially we focused on traditional methods, which are based on hand-crafted features combined with machine learning techniques. We then shifted our attention to more recent techniques, which are based on deep learning networks. In terms of performance, deep learning-based techniques significantly outperform traditional methods. However, these methods require large datasets for efficient training. We expect that in the near future, this will become less of an issue since RGB-D sensors are becoming cheaper and annotated RGB-D datasets will thus become widely available.

Although they achieve remarkable performance compared to traditional methods, deep learning techniques are still in their infancy. For instance, amodal object detection, i.e. estimating the entire 3D bounding box of an object, especially when parts of the object are occluded, still remains challenging especially in highly cluttered scenes. Self-occluding objects also challenge deep learning-based pipelines, especially when dealing with dynamic objects that deform in a non-rigid way [33, 44]. Similarly, object detection performance for objects at various scales, particularly at small scales, is still relatively low. In fact, the performance comparison of Table 3.3 does not consider the robustness of the methods to scale variation.

Existing techniques focus mainly on the detection of the bounding boxes of the objects of interest. However, many situations, e.g. robust grasping, image editing, and accurate robot navigation, require the accurate detection of object boundaries. Several works have attempted to achieve this using, for example, template matching. This, however, remains an open problem.

Another important avenue for future research is how to incorporate spatial relationships and relationships between semantic classes in deep learning-based RGB-D object detection pipelines. These relationships are important cues for recognition and it has been already shown in many papers that they can significantly boost the performance of traditional techniques [43]. Yet, this knowledge is not efficiently exploited in deep learning techniques.

Finally, there are many aspects of object detection from RGB-D images that have not been covered in this chapter. Examples include saliency detection [11, 47, 66, 78], which aims to detect salient regions in an RGB-D image. Additionally, we have focused in this chapter on generic objects in indoor and outdoor scenes. There is, however, a rich literature on specialized detectors which focus on specific classes of objects, e.g. humans and human body parts such as faces and hands.

Acknowledgements This work is supported by ARC DP 150100294 and ARC DP 150104251.

References

1. Cornell grasping dataset. URL http://pr.cs.cornell.edu/grasping/rect_data/data.php. Accessed: 2018-12-13
2. Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Süsstrunk, S.: Slic superpixels compared to state-of-the-art superpixel methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **34**(11), 2274–2282 (2012). DOI 10.1109/TPAMI.2012.120
3. Alexandre, L.A.: 3d object recognition using convolutional neural networks with transfer learning between input channels. In: IAS (2014)
4. Alexe, B., Deselaers, T., Ferrari, V.: What is an object? In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 73–80 (2010). DOI 10.1109/CVPR.2010.5540226
5. Arbeláez, P., Pont-Tuset, J., Barron, J., Marques, F., Malik, J.: Multiscale combinatorial grouping. In: Computer Vision and Pattern Recognition (2014)
6. Asif, U., Bennamoun, M., Sohel, F.A.: Rgb-d object recognition and grasp detection using hierarchical cascaded forests. *IEEE Transactions on Robotics* **33**(3), 547–564 (2017). DOI 10.1109/TRO.2016.2638453
7. Barrow, H.G., Tenenbaum, J.M., Bolles, R.C., Wolf, H.C.: Parametric correspondence and chamfer matching: Two new techniques for image matching. In: Proceedings of the 5th International Joint Conference on Artificial Intelligence - Volume 2, IJCAI'77, pp. 659–663. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA (1977). URL <http://dl.acm.org/citation.cfm?id=1622943.1622971>
8. Bleyer, M., Rhemann, C., Rother, C.: Extracting 3d scene-consistent object proposals and depth from stereo images. In: European Conference on Computer Vision, pp. 467–481. Springer (2012)
9. Bo, L., Ren, X., Fox, D.: Learning hierarchical sparse features for rgb-(d) object recognition. *The International Journal of Robotics Research* **33**(4), 581–599 (2014)
10. Buch, N.E., Orwell, J., Velastin, S.A.: 3d extended histogram of oriented gradients (3dhog) for classification of road users in urban scenes. In: BMVC (2009)
11. Chen, H., Li, Y.: Progressively complementarity-aware fusion network for rgb-d salient object detection. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

12. Chen, X., Ma, H., Wan, J., Li, B., Xia, T.: Multi-view 3D object detection network for autonomous driving. In: IEEE CVPR, vol. 1, p. 3 (2017)
13. Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. In: 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05), vol. 1, pp. 886–893 vol. 1 (2005). DOI 10.1109/CVPR.2005.177
14. David G. Lowe: Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision (IJCV)* (2004)
15. Deng, Z., Latecki, L.J.: Amodal detection of 3D objects: Inferring 3D bounding boxes from 2D ones in RGB-depth images. In: Conference on Computer Vision and Pattern Recognition (CVPR), vol. 2, p. 2 (2017)
16. E. Schapire, R.: Explaining AdaBoost, pp. 37–52 (2013). DOI 10.1007/978-3-642-41136-6-5
17. Eitel, A., Springenberg, J.T., Spinello, L., Riedmiller, M.A., Burgard, W.: Multimodal deep learning for robust RGB-D object recognition. *CoRR abs/1507.06821* (2015). URL <http://arxiv.org/abs/1507.06821>
18. Engelcke, M., Rao, D., Wang, D.Z., Tong, C.H., Posner, I.: Vote3deep: Fast object detection in 3D point clouds using efficient convolutional neural networks. In: Robotics and Automation (ICRA), 2017 IEEE International Conference on, pp. 1355–1361. IEEE (2017)
19. Everingham, M., Van Gool, L., Williams, C., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. In: 2010 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2010)
20. Fei-Fei, L., Fergus, R., Perona, P.: Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. 2004 Conference on Computer Vision and Pattern Recognition Workshop pp. 178–178 (2004)
21. Felzenszwalb, P.F., Girshick, R.B., McAllester, D., Ramanan, D.: Object detection with discriminatively trained part-based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **32**(9), 1627–1645 (2010)
22. Feng, D., Barnes, N., You, S., McCarthy, C.: Local background enclosure for rgb-d salient object detection. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2343–2350 (2016). DOI 10.1109/CVPR.2016.257
23. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2012)
24. Getto, R., Fellner, D.W.: 3d object retrieval with parametric templates. In: Proceedings of the 2015 Eurographics Workshop on 3D Object Retrieval, 3DOR '15, pp. 47–54. Eurographics Association, Goslar Germany, Germany (2015). DOI 10.2312/3dor.20151054. URL <https://doi.org/10.2312/3dor.20151054>
25. Gidaris, S., Komodakis, N.: Attend refine repeat: Active box proposal generation via in-out localization. *CoRR abs/1606.04446* (2016). URL <http://arxiv.org/abs/1606.04446>
26. Girshick, R.B.: Fast R-CNN. *CoRR abs/1504.08083* (2015). URL <http://arxiv.org/abs/1504.08083>
27. Griffin, G., Holub, A., Perona, P.: Caltech-256 object category dataset. Tech. Rep. 7694, California Institute of Technology (2007). URL <http://authors.library.caltech.edu/7694>
28. Gupta, S., Arbeláez, P., Girshick, R., Malik, J.: Aligning 3d models to rgb-d images of cluttered scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4731–4740 (2015)
29. Gupta, S., Girshick, R.B., Arbelaez, P., Malik, J.: Learning rich features from RGB-D images for object detection and segmentation. *CoRR abs/1407.5736* (2014). URL <http://arxiv.org/abs/1407.5736>
30. Han, J., Zhang, D., Cheng, G., Liu, N., Xu, D.: Advanced deep-learning techniques for salient and category-specific object detection: A survey. *IEEE Signal Processing Magazine* **35**(1), 84–100 (2018). DOI 10.1109/MSP.2017.2749125
31. He, D., Wang, L.: Texture unit, texture spectrum, and texture analysis. *IEEE Transactions on Geoscience and Remote Sensing* **28**(4), 509–512 (1990). DOI 10.1109/TGRS.1990.572934

32. Hou, Q., Cheng, M., Hu, X., Borji, A., Tu, Z., Torr, P.H.S.: Deeply supervised salient object detection with short connections. *CoRR abs/1611.04849* (2016). URL <http://arxiv.org/abs/1611.04849>
33. Jermyn, I.H., Kurtek, S., Laga, H., Srivastava, A.: Elastic shape analysis of three-dimensional objects. *Synthesis Lectures on Computer Vision* **12**(1), 1–185 (2017)
34. Jiang, H.: Finding approximate convex shapes in rgb-d images. In: *European Conference on Computer Vision*, pp. 582–596. Springer (2014)
35. Jiang, H., Xiao, J.: A linear approach to matching cuboids in rgb-d images. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2171–2178 (2013)
36. Ju, R., Ge, L., Geng, W., Ren, T., Wu, G.: Depth saliency based on anisotropic center-surround difference. In: *2014 IEEE International Conference on Image Processing (ICIP)*, pp. 1115–1119 (2014). DOI 10.1109/ICIP.2014.7025222
37. Kendall, A., Cipolla, R.: Geometric loss functions for camera pose regression with deep learning. *CoRR abs/1704.00390* (2017). URL <http://arxiv.org/abs/1704.00390>
38. Khan, S., Rahmani, H., Shah, S.A.A., Bennamoun, M.: *A Guide to Convolutional Neural Networks for Computer Vision*. Morgan and Claypool Publishers (2018)
39. Khan, S.H., He, X., Bennamoun, M., Sohel, F., Togneri, R.: Separating objects and clutter in indoor scenes. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4603–4611 (2015)
40. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1, NIPS’12*, pp. 1097–1105. Curran Associates Inc., USA (2012). URL <http://dl.acm.org/citation.cfm?id=2999134.2999257>
41. Laga, H., Guo, Y., Tabia, H., Fisher, R.B., Bennamoun, M.: *3D Shape Analysis: Fundamentals, Theory, and Applications*. John Wiley & Sons (2018)
42. Laga, H., Guo, Y., Tabia, H., Fisher, R.B., Bennamoun, M.: *3D Shape Analysis: Fundamentals, Theory, and Applications*. Wiley (2019)
43. Laga, H., Mortara, M., Spagnuolo, M.: Geometry and context for semantic correspondences and functionality recognition in man-made 3d shapes. *ACM Transactions on Graphics (TOG)* **32**(5), 150 (2013)
44. Laga, H., Xie, Q., Jermyn, I.H., Srivastava, A.: Numerical inversion of snrf maps for elastic shape analysis of genus-zero surfaces. *IEEE transactions on pattern analysis and machine intelligence* **39**(12), 2451–2464 (2017)
45. Lahoud, J., Ghanem, B.: 2D-Driven 3D Object Detection in RGB-D Images. In: *The IEEE International Conference on Computer Vision (ICCV)* (2017)
46. Lai, K., Bo, L., Ren, X., Fox, D.: Rgb-d object recognition: Features, algorithms, and a large scale benchmark. In: *Consumer Depth Cameras for Computer Vision*, pp. 167–192. Springer (2013)
47. Lei, Z., Chai, W., Zhao, S., Song, H., Li, F.: Saliency detection for rgb-d images using optimization. In: *2017 12th International Conference on Computer Science and Education (ICCSE)*, pp. 440–443 (2017). DOI 10.1109/ICCSE.2017.8085532
48. Li, B.: 3d fully convolutional network for vehicle detection in point cloud. *CoRR abs/1611.08069* (2016). URL <http://arxiv.org/abs/1611.08069>
49. Li, B., Zhang, T., Xia, T.: Vehicle detection from 3D lidar using fully convolutional network. *arXiv preprint arXiv:1608.07916* (2016)
50. Li, N., Ye, J., Ji, Y., Ling, H., Yu, J.: Saliency detection on light field. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **39**(8), 1605–1616 (2017). DOI 10.1109/TPAMI.2016.2610425
51. Lin, D., Fidler, S., Urtasun, R.: Holistic scene understanding for 3D object detection with RGBD cameras. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1417–1424 (2013)
52. Lin, T.Y., Dollár, P., Girshick, R.B., He, K., Hariharan, B., Belongie, S.J.: Feature pyramid networks for object detection. In: *CVPR*, vol. 1, p. 4 (2017)
53. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. *CoRR abs/1411.4038* (2014). URL <http://arxiv.org/abs/1411.4038>

54. Maturana, D., Scherer, S.: 3d convolutional neural networks for landing zone detection from lidar. In: 2015 IEEE International Conference on Robotics and Automation (ICRA), pp. 3471–3478 (2015). DOI 10.1109/ICRA.2015.7139679
55. Maturana, D., Scherer, S.: VoxNet: A 3D Convolutional Neural Network for real-time object recognition. In: Ieee/rsj International Conference on Intelligent Robots and Systems, pp. 922–928 (2015)
56. Menze, M., Geiger, A.: Object scene flow for autonomous vehicles. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2015)
57. Nakahara, H., Yonekawa, H., Sato, S.: An object detector based on multiscale sliding window search using a fully pipelined binarized cnn on an fpga. In: 2017 International Conference on Field Programmable Technology (ICFPT), pp. 168–175 (2017). DOI 10.1109/FPT.2017.8280135
58. Nathan Silberman Derek Hoiem, P.K., Fergus, R.: Indoor segmentation and support inference from rgb-d images. In: ECCV (2012)
59. Newcombe, R.A., Izadi, S., Hilliges, O., Molyneaux, D., Kim, D., Davison, A.J., Kohi, P., Shotton, J., Hodges, S., Fitzgibbon, A.: Kinectfusion: Real-time dense surface mapping and tracking. In: 2011 10th IEEE International Symposium on Mixed and Augmented Reality, pp. 127–136 (2011). DOI 10.1109/ISMAR.2011.6092378
60. Peng, H., Li, B., Xiong, W., Hu, W., Ji, R.: Rgb-d salient object detection: A benchmark and algorithms. In: ECCV (2014)
61. Pont-Tuset, J., Arbeláez, P., Barron, J., Marques, F., Malik, J.: Multiscale combinatorial grouping for image segmentation and object proposal generation. In: arXiv:1503.00848 (2015)
62. Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3d object detection from rgb-d data. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
63. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE **1**(2), 4 (2017)
64. Qi, C.R., Su, H., Nießner, M., Dai, A., Yan, M., Guibas, L.J.: Volumetric and multi-view cnns for object classification on 3d data. *CoRR* **abs/1604.03265** (2016). URL <http://arxiv.org/abs/1604.03265>
65. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: *Advances in Neural Information Processing Systems*, pp. 5099–5108 (2017)
66. Qu, L., He, S., Zhang, J., Tian, J., Tang, Y., Yang, Q.: Rgb-d salient object detection via deep fusion. *IEEE Transactions on Image Processing* **26**(5), 2274–2285 (2017). DOI 10.1109/TIP.2017.2682981
67. Ren, J., Gong, X., Yu, L., Zhou, W., Yang, M.Y.: Exploiting global priors for rgb-d saliency detection. In: 2015 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 25–32 (2015). DOI 10.1109/CVPRW.2015.7301391
68. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: *Advances in neural information processing systems*, pp. 91–99 (2015)
69. Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: towards real-time object detection with region proposal networks. *CoRR* **abs/1506.01497** (2015). URL <http://arxiv.org/abs/1506.01497>
70. Ren, Z., Sudderth, E.B.: Three-dimensional object detection and layout prediction using clouds of oriented gradients. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1525–1533 (2016). DOI 10.1109/CVPR.2016.169
71. Rusinkiewicz, S., Levoy, M.: Efficient variants of the icp algorithm. In: *Proceedings Third International Conference on 3-D Digital Imaging and Modeling*, pp. 145–152 (2001). DOI 10.1109/IM.2001.924423
72. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition

- Challenge. *International Journal of Computer Vision (IJCV)* **115**(3), 211–252 (2015). DOI 10.1007/s11263-015-0816-y
73. Sahin, C., Kouskouridas, R., Kim, T.: Iterative hough forest with histogram of control points for 6 dof object registration from depth images. *CoRR* **abs/1603.02617** (2016). URL <http://arxiv.org/abs/1603.02617>
 74. Sahin, C., Kouskouridas, R., Kim, T.: A learning-based variable size part extraction architecture for 6d object pose recovery in depth. *CoRR* **abs/1701.02166** (2017). URL <http://arxiv.org/abs/1701.02166>
 75. Schwarz, M., Schulz, H., Behnke, S.: Rgb-d object recognition and pose estimation based on pre-trained convolutional neural network features. In: 2015 IEEE International Conference on Robotics and Automation (ICRA), pp. 1329–1335 (2015). DOI 10.1109/ICRA.2015.7139363
 76. Shotton, J., Girshick, R., Fitzgibbon, A., Sharp, T., Cook, M., Finocchio, M., Moore, R., Kohli, P., Criminisi, A., Kipman, A., et al.: Efficient human pose estimation from single depth images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **35**(12), 2821–2840 (2013)
 77. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
 78. Song, H., Liu, Z., Xie, Y., Wu, L., Huang, M.: RGBD co-saliency detection via bagging-based clustering. *IEEE Signal Processing Letters* **23**(12), 1722–1726 (2016)
 79. Song, S., Lichtenberg, S.P., Xiao, J.: Sun rgb-d: A rgb-d scene understanding benchmark suite. In: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 567–576. IEEE (2015)
 80. Song, S., Xiao, J.: Sliding shapes for 3d object detection in depth images. In: European conference on computer vision, pp. 634–651. Springer (2014)
 81. Song, S., Xiao, J.: Deep Sliding Shapes for amodal 3D object detection in RGB-D images. In: CVPR (2016)
 82. Sun, H., Meng, Z., Tao, P.Y., Ang, M.H.: Scene recognition and object detection in a unified convolutional neural network on a mobile manipulator. In: 2018 IEEE International Conference on Robotics and Automation (ICRA), pp. 1–5 (2018). DOI 10.1109/ICRA.2018.8460535
 83. Wang, A., Cai, J., Lu, J., Cham, T.J.: MMSS: Multi-modal sharable and specific feature learning for RGB-D object recognition. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 1125–1133 (2015)
 84. Wang, D.Z., Posner, I.: Voting for voting in online point cloud object detection. In: Robotics: Science and Systems (2015)
 85. Ward, I.R., Jalwana, M.A.A.K., Bennamoun, M.: Improving image-based localization with deep learning: The impact of the loss function **abs/1905.03692** (2019). URL <http://arxiv.org/abs/1905.03692>
 86. Wu, Z., Song, S., Khosla, A., Tang, X., Xiao, J.: 3d shapenets for 2.5d object recognition and next-best-view prediction. *CoRR* (2014)
 87. Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1912–1920 (2015). DOI 10.1109/CVPR.2015.7298801
 88. Xiang, Y., Kim, W., Chen, W., Ji, J., Choy, C., Su, H., Mottaghi, R., Guibas, L., Savarese, S.: ObjectNet3D: A large scale database for 3D object recognition. In: European Conference on Computer Vision, pp. 160–176. Springer (2016)
 89. Xiang, Y., Mottaghi, R., Savarese, S.: Beyond pascal: A benchmark for 3d object detection in the wild. In: Applications of Computer Vision (WACV), 2014 IEEE Winter Conference on, pp. 75–82. IEEE (2014)
 90. Xiao, J., Hays, J., Ehinger, K.A., Oliva, A., Torralba, A.: Sun database: Large-scale scene recognition from abbey to zoo. In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 3485–3492 (2010). DOI 10.1109/CVPR.2010.5539970