# Deep Learning for Generic Object Detection: A Survey

**Li Liu** [1,2] · **Wanli Ouyang** [3] · **Xiaogang Wang** [4] ·
**Paul Fieguth** [5] · **Jie Chen** [2] · **Xinwang Liu** [1] · **Matti Pietikäinen** [2]

**Abstract** Object detection, one of the most fundamental and challenging problems in computer vision, seeks to locate object instances from a large number of predefined categories in natural images. Deep learning techniques have emerged as a powerful strategy for learning feature representations directly from data and have led to remarkable breakthroughs in the field of generic object detection. Given this period of rapid evolution, the goal of this paper is to provide a comprehensive survey of the recent achievements in this field brought about by deep learning techniques. More than 300 research contributions are included in this survey, covering many aspects of generic object detection: detection frameworks, object feature representation, object proposal generation, context modeling, training strategies, and evaluation metrics. We finish the survey by identifying promising directions for future research.

**Keywords** Object detection · deep learning · convolutional neural networks · object recognition

**Fig. 1** Most frequent keywords in ICCV and CVPR conference papers from 2016 to 2018. The size of each word is proportional to the frequency of that keyword. We can see that object detection has received significant attention in recent years.

## 1 Introduction

As a longstanding, fundamental and challenging problem in computer vision, object detection (illustrated in Fig. 1) has been an active area of research for several decades [76]. The goal of object detection is to determine whether there are any instances of objects from given categories (such as humans, cars, bicycles, dogs or cats) in an image and, if present, to return the spatial location and extent of each object instance (*e.g.,* via a bounding box [68, 234]). As the cornerstone of image understanding and computer vision, object detection forms the basis for solving complex or high level vi-

✉ Li Liu (li.liu@oulu.fi)
Wanli Ouyang (wanli.ouyang@sydney.edu.au)
Xiaogang Wang (xgwang@ee.cuhk.edu.hk)
Paul Fieguth (pfieguth@uwaterloo.ca)
Jie Chen (jie.chen@oulu.fi)
Xinwang Liu (xinwangliu@nudt.edu.cn)
Matti Pietikäinen (matti.pietikainen@oulu.fi)
1 National University of Defense Technology, China
2 University of Oulu, Finland
3 University of Sydney, Australia
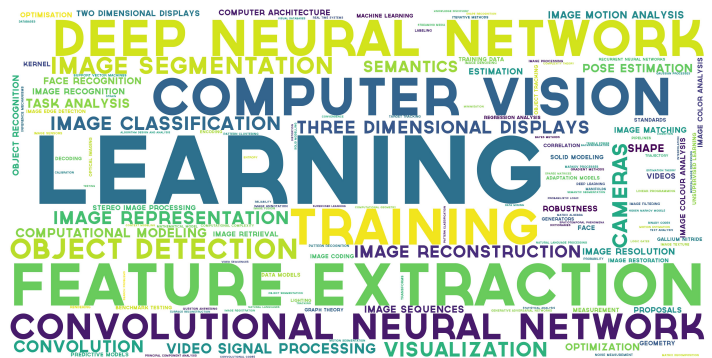4 Chinese University of Hong Kong, China
5 University of Waterloo, Canada

sion tasks such as segmentation, scene understanding, object tracking, image captioning, event detection, and activity recognition. Object detection supports a wide range of applications, including robot vision, consumer electronics, security, autonomous driving, human computer interaction, content based image retrieval, intelligent video surveillance, and augmented reality.

Recently, deep learning techniques [105, 149] have emerged as powerful methods for learning feature representations automatically from data. In particular, these techniques have provided major improvements in object detection, as illustrated in Fig. 3.

As illustrated in Fig. 2, object detection can be grouped into one of two types [91, 310]: detection of specific instances versus the detection of broad categories. The first type aims to detect instances of a particular object (such as Donald Trump's face, the Eiffel Tower, or a neighbor's dog), essentially a matching problem. The goal of the second type is to detect (usually previously unseen) instances of some predefined object categories (for example humans, cars, bicycles, and dogs). Historically, much of the effort in the field of object detection has focused on the detection of a single category (typically faces and pedestrians) or a few specific categories. In contrast, over the past several years, the research community has started moving towards the more challenging goal of building general purpose object detection systems where the breadth of object detection ability rivals that of humans.
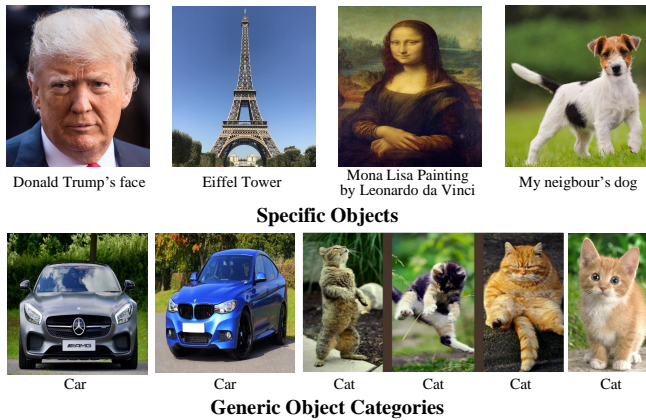
**Specific Objects**

Donald Trump's face — Eiffel Tower — Mona Lisa Painting by Leonardo da Vinci — My neigbour's dog

**Generic Object Categories**

Car — Car — Cat — Cat — Cat — Cat

**Fig. 2** Object detection includes localizing instances of a *particular* object (top), as well as generalizing to detecting object *categories* in general (bottom). This survey focuses on recent advances for the latter problem of generic object detection.
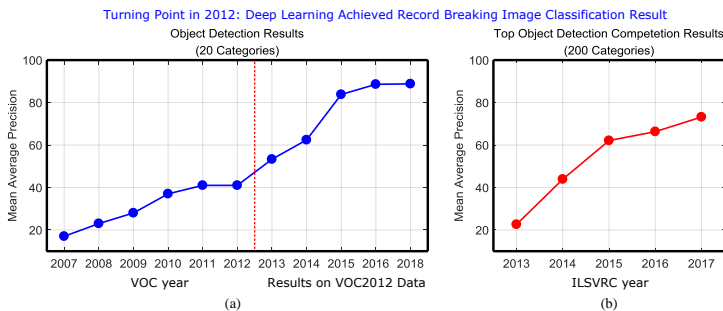


**Fig. 3** An overview of recent object detection performance: We can observe a significant improvement in performance (measured as mean average precision) since the arrival of deep learning in 2012. (a) Detection results of winning entries in the VOC2007-2012 competitions, and (b) Top object detection competition results in ILSVRC2013-2017 (results in both panels use only the provided training data).

In 2012, Krizhevsky *et al.* [140] proposed a Deep Convolutional Neural Network (DCNN) called AlexNet which achieved record breaking image classification accuracy in the Large Scale Visual Recognition Challenge (ILSVRC) [234]. Since that time, the research focus in most aspects of computer vision has been specifically on deep learning methods, indeed including the domain of generic object detection [85, 99, 84, 239, 230]. Although tremendous progress has been achieved, illustrated in Fig. 3, we are unaware of comprehensive surveys of this subject over the past five years. Given the exceptionally rapid rate of progress, this article attempts to track recent advances and summarize their achievements in order to gain a clearer picture of the current panorama in generic object detection.

## 1.1 Comparison with Previous Reviews

Many notable object detection surveys have been published, as summarized in Table 1. These include many excellent surveys on the problem of *specific* object detection, such as pedestrian detection [66, 79, 59], face detection [294, 301], vehicle detection [258] and text detection [295]. There are comparatively few recent surveys focusing directly on the problem of generic object detection, except for the work by Zhang *et al.* [310] who conducted a survey on the topic of object class detection. However, the research

reviewed in [91], [5] and [310] is mostly pre-2012, and therefore prior to the recent striking success and dominance of deep learning and related methods.

Deep learning allows computational models to learn fantastically complex, subtle, and abstract representations, driving significant progress in a broad range of problems such as visual recognition, object detection, speech recognition, natural language processing, medical image analysis, drug discovery and genomics. Among different types of deep neural networks, DCNNs [148, 140, 149] have brought about breakthroughs in processing images, video, speech and audio. To be sure, there have been many published surveys on deep learning, including that of Bengio *et al.* [13], LeCun *et al.* [149], Litjens *et al.* [170], Gu *et al.* [92], and more recently in tutorials at ICCV and CVPR.

In contrast, although many deep learning based methods have been proposed for object detection, we are unaware of any comprehensive recent survey. A thorough review and summary of existing work is essential for further progress in object detection, particularly for researchers wishing to enter the field. Since our focus is on *generic* object detection, the extensive work on DCNNs for *specific* object detection, such as face detection [154, 306, 116], pedestrian detection [307, 109], vehicle detection [322] and traffic sign detection [329] will not be considered.

## 1.2 Scope

The number of papers on generic object detection based on deep learning is breathtaking. There are so many, in fact, that compiling any comprehensive review of the state of the art is beyond the scope of any reasonable length paper. As a result, it is necessary to establish selection criteria, in such a way that we have limited our focus to top journal and conference papers. Due to these limitations, we sincerely apologize to those authors whose works are not included in this paper. For surveys of work on related topics, readers are referred to the articles in Table 1. This survey focuses on major progress of the last five years, and we restrict our attention to still pictures, leaving the important subject of video object detection as a topic for separate consideration in the future.

The main goal of this paper is to offer a comprehensive survey of deep learning based generic object detection techniques, and to present some degree of taxonomy, a high level perspective and organization, primarily on the basis of popular datasets, evaluation metrics, context modeling, and detection proposal methods. The intention is that our categorization be helpful for readers to have an accessible understanding of similarities and differences between a wide variety of strategies. The proposed taxonomy gives researchers a framework to understand current research and to identify open challenges for future research.

The remainder of this paper is organized as follows. Related background and the progress made during the last two decades are summarized in Section 2. A brief introduction to deep learning is given in Section 3. Popular datasets and evaluation criteria are summarized in Section 4. We describe the milestone object detection frameworks in Section 5. From Section 6 to Section 9, fundamental sub-problems and the relevant issues involved in designing object detectors are discussed. Finally, in Section 10, we conclude
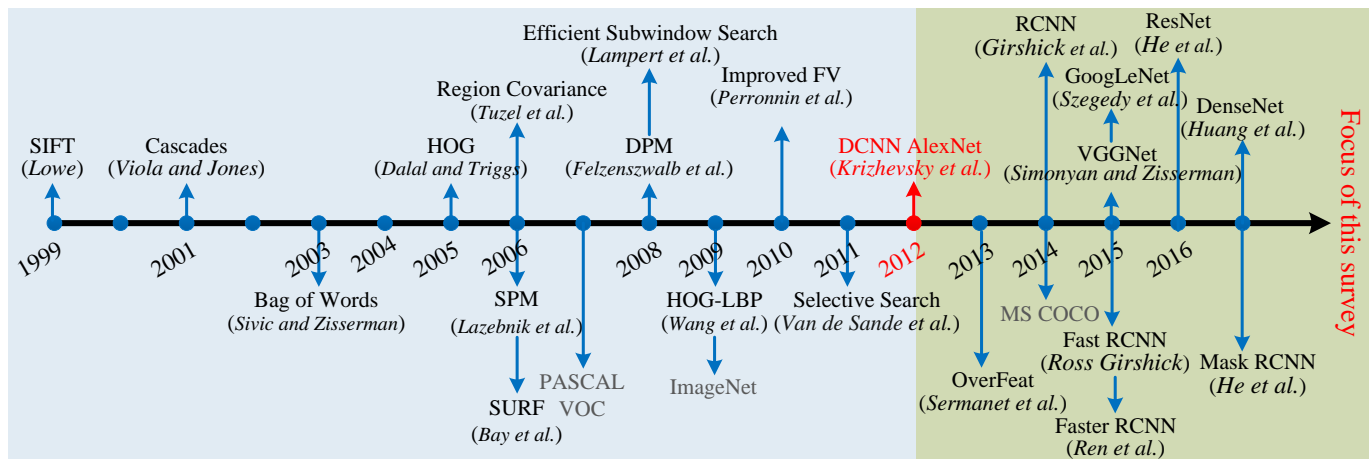
**Fig. 4** Milestones of object detection and recognition, including feature representations [47, 52, 101, 140, 147, 178, 179, 212, 248, 252, 263, 276, 279], detection frameworks [74, 85, 239, 271, 276], and datasets [68, 166, 234]. The time period up to 2012 is dominated by handcrafted features, a transition took place in 2012 with the development of DCNNs for image classification by Krizhevsky *et al.* [140], with methods after 2012 dominated by related deep networks. Mostof the listed methods are highly cited and won a major ICCV or CVPR prize. See Section 2.3 for details.

**Table 1** Summary of related object detection surveys since 2000.

| No. | Survey Title | Ref. | Year | Venue | Content |
|---|---|---|---|---|---|
| 1 | Monocular Pedestrian Detection: Survey and Experiments | [66] | 2009 | PAMI | An evaluation of three pedestrian detectors |
| 2 | Survey of Pedestrian Detection for Advanced Driver Assistance Systems | [79] | 2010 | PAMI | A survey of pedestrian detection for advanced driver assistance systems |
| 3 | Pedestrian Detection: An Evaluation of the State of The Art | [59] | 2012 | PAMI | A thorough and detailed evaluation of detectors in monocular images |
| 4 | Detecting Faces in Images: A Survey | [294] | 2002 | PAMI | First survey of face detection from a single image |
| 5 | A Survey on Face Detection in the Wild: Past, Present and Future | [301] | 2015 | CVIU | A survey of face detection in the wild since 2000 |
| 6 | On Road Vehicle Detection: A Review | [258] | 2006 | PAMI | A review of vision based on-road vehicle detection systems |
| 7 | Text Detection and Recognition in Imagery: A Survey | [295] | 2015 | PAMI | A survey of text detection and recognition in color imagery |
| 8 | Toward Category Level Object Recognition | [215] | 2007 | Book | Representative papers on object categorization, detection, and segmentation |
| 9 | The Evolution of Object Categorization and the Challenge of Image Abstraction | [56] | 2009 | Book | A trace of the evolution of object categorization over four decades |
| 10 | Context based Object Categorization: A Critical Survey | [78] | 2010 | CVIU | A review of contextual information for object categorization |
| 11 | 50 Years of Object Recognition: Directions Forward | [5] | 2013 | CVIU | A review of the evolution of object recognition systems over five decades |
| 12 | Visual Object Recognition | [91] | 2011 | Tutorial | Instance and category object recognition techniques |
| 13 | Object Class Detection: A Survey | [310] | 2013 | ACM CS | Survey of generic object detection methods before 2011 |
| 14 | Feature Representation for Statistical Learning based Object Detection: A Review | [160] | 2015 | PR | Feature representation methods in statistical learning based object detection, including handcrafted and deep learning based features |
| 15 | Salient Object Detection: A Survey | [19] | 2014 | arXiv | A survey for salient object detection |
| 16 | Representation Learning: A Review and New Perspectives | [13] | 2013 | PAMI | Unsupervised feature learning and deep learning, probabilistic models, autoencoders, manifold learning, and deep networks |
| 17 | Deep Learning | [149] | 2015 | Nature | An introduction to deep learning and applications |
| 18 | A Survey on Deep Learning in Medical Image Analysis | [170] | 2017 | MIA | A survey of deep learning for image classification, object detection, segmentation and registration in medical image analysis |
| 19 | Recent Advances in Convolutional Neural Networks | [92] | 2017 | PR | A broad survey of the recent advances in CNN and its applications in computer vision, speech and natural language processing |
| 20 | Tutorial: Tools for Efficient Object Detection | – | 2015 | ICCV15 | A short course for object detection only covering recent milestones |
| 21 | Tutorial: Deep Learning for Objects and Scenes | – | 2017 | CVPR17 | A high level summary of recent work on deep learning for visual recognition of objects and scenes |
| 22 | Tutorial: Instance Level Recognition | – | 2017 | ICCV17 | A short course of recent advances on instance level recognition, including object detection, instance segmentation and human pose prediction |
| 23 | Tutorial: Visual Recognition and Beyond | – | 2018 | CVPR18 | A tutorial on methods and principles behind image classification, object detection, instance segmentation, and semantic segmentation. |
| **24** | **Deep Learning for Generic Object Detection** | **Ours** | **2019** | **VISI** | **A comprehensive survey of deep learning for generic object detection** |

the paper with an overall discussion of object detection, state-of-the- art performance, and future research directions.
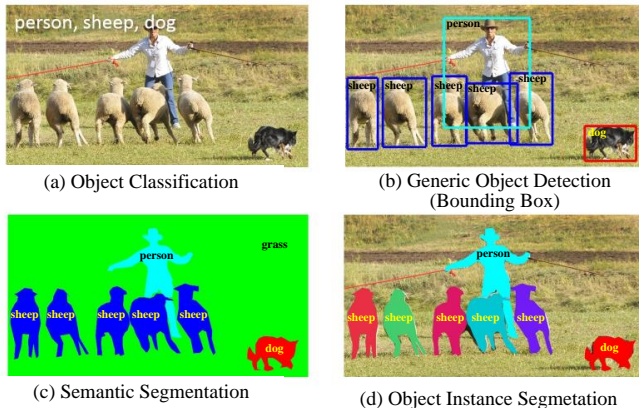
**Fig. 5** Recognition problems related to generic object detection: (a) Image level object classification, (b) Bounding box level generic object detection, (c) Pixel-wise semantic segmentation, (d) Instance level semantic segmentation.

## 2 Generic Object Detection

### 2.1 The Problem

*Generic object detection*, also called generic object category detection, object class detection, or object category detection [310], is defined as follows. Given an image, determine whether or not there are instances of objects from predefined categories (usually *many* categories, *e.g.,* 200 categories in the ILSVRC object detection challenge) and, if present, to return the spatial location and extent of each instance. A greater emphasis is placed on detecting a broad range of natural categories, as opposed to specific object category detection where only a narrower predefined category of interest (*e.g.,* faces, pedestrians, or cars) may be present. Although thousands of objects occupy the visual world in which we live, currently the research community is primarily interested in the localization of highly structured objects (*e.g.,* cars, faces, bicycles and airplanes) and articulated objects (*e.g.,* humans, cows and horses) rather than unstructured scenes (such as sky, grass and cloud).

The spatial location and extent of an object can be defined coarsely using a bounding box (an axis-aligned rectangle tightly bounding the object) [68, 234], a precise pixelwise segmentation mask [310], or a closed boundary [166, 235], as illustrated in Fig. 5. To the best of our knowledge, for the evaluation of generic object detection algorithms, it is bounding boxes which are most widely used in the current literature [68, 234], and therefore this is also the approach we adopt in this survey. However, as the research community moves towards deeper scene understanding (from image level object classification to single object localization, to generic object detection, and to pixelwise object segmentation), it is anticipated that future challenges will be at the pixel level [166].

There are many problems closely related to that of generic object detection[1]. The goal of *object classification* or *object categorization* (Fig. 5 (a)) is to assess the presence of objects from a given set of object classes in an image; *i.e.,* assigning one or more object class labels to a given image, determining the presence without the need of location. The additional requirement to locate the in-
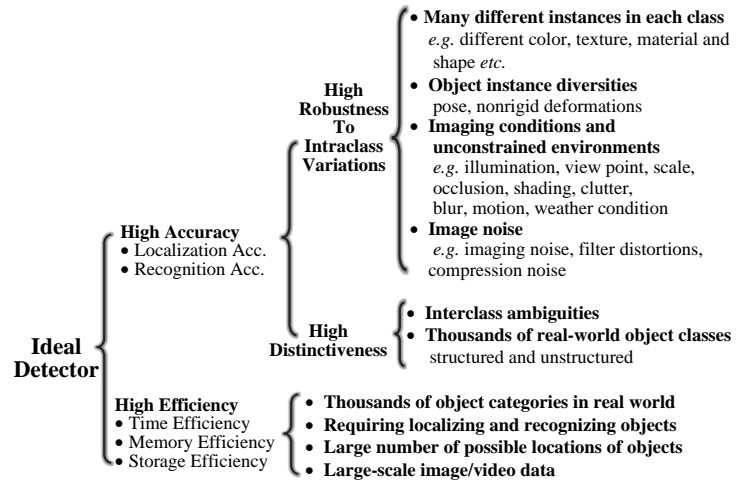
---

[1] To the best of our knowledge, there is no universal agreement in the literature on the definitions of various vision subtasks. Terms such as detection, localization, recognition, classification, categorization, verification, identification, annotation, labeling, and understanding are often differently defined [5].



**Fig. 6** Taxonomy of challenges in generic object detection.



**Fig. 7** Changes in appearance of the same class with variations in imaging conditions (a-h). There is an astonishing variation in what is meant to be a single object class (i). In contrast, the four images in (j) appear very similar, but in fact are from four *d*ifferent object classes. Most images are from ImageNet [234] and MS COCO [166].

stances in an image makes detection a more challenging task than classification. The *object recognition* problem denotes the more general problem of identifying/localizing all the objects present in an image, subsuming the problems of object detection and classification [68, 234, 198, 5]. Generic object detection is closely related to *semantic image segmentation* (Fig. 5 (c)), which aims to assign each pixel in an image to a semantic class label. *Object instance segmentation* (Fig. 5 (d)) aims to distinguish different instances of the same object class, as opposed to semantic segmentation which does not.

### 2.2 Main Challenges

The ideal of generic object detection is to develop a general-purpose algorithm that achieves two competing goals of *high quality/accuracy* and *high efficiency* (Fig. 6). As illustrated in Fig. 7, high quality

detection must accurately localize and recognize objects in images or video frames, such that the large variety of object categories in the real world can be distinguished (*i.e.,* high distinctiveness), and that object instances from the same category, subject to intra-class appearance variations, can be localized and recognized (*i.e.,* high robustness). High efficiency requires that the entire detection task runs in real time with acceptable memory and storage demands.

### 2.2.1 Accuracy related challenges

Challenges in detection accuracy stem from 1) the vast range of intra-class variations and 2) the huge number of object categories.

Intra-class variations can be divided into two types: intrinsic factors and imaging conditions. In terms of intrinsic factors, each object category can have many different object instances, possibly varying in one or more of color, texture, material, shape, and size, such as the "chair" category shown in Fig. 7 (*i*). Even in a more narrowly defined class, such as human or horse, object instances can appear in different poses, subject to nonrigid deformations or with the addition of clothing.

Imaging condition variations are caused by the dramatic impacts unconstrained environments can have on object appearance, such as lighting (dawn, day, dusk, indoors), physical location, weather conditions, cameras, backgrounds, illuminations, occlusion, and viewing distances. All of these conditions produce significant variations in object appearance, such as illumination, pose, scale, occlusion, clutter, shading, blur and motion, with examples illustrated in Fig. 7 (*a-h*). Further challenges may be added by digitization artifacts, noise corruption, poor resolution, and filtering distortions.

In addition to *intra*class variations, the large number of object categories, on the order of $10^4 - 10^5$, demands great discrimination power from the detector to distinguish between subtly different *inter*class variations, as illustrated in Fig. 7 (j). In practice, current detectors focus mainly on structured object categories, such as the 20, 200 and 91 object classes in PASCAL VOC [68], ILSVRC [234] and MS COCO [166] respectively. Clearly, the number of object categories under consideration in existing benchmark datasets is much smaller than can be recognized by humans.

### 2.2.2 Efficiency and scalability related challenges

The prevalence of social media networks and mobile/wearable devices has led to increasing demands for analyzing visual data. However, mobile/wearable devices have limited computational capabilities and storage space, making efficient object detection critical.

The efficiency challenges stem from the need to localize and recognize, computational complexity growing with the (possibly large) number of object categories, and with the (possibly very large) number of locations and scales within a single image, such as the examples in Fig. 7 (c, d).

A further challenge is that of scalability: A detector should be able to handle previously unseen objects, unknown situations, and high data rates. As the number of images and the number of categories continue to grow, it may become impossible to annotate them manually, forcing a reliance on weakly supervised strategies.

### 2.3 Progress in the Past Two Decades

Early research on object recognition was based on template matching techniques and simple part-based models [76], focusing on specific objects whose spatial layouts are roughly rigid, such as faces. Before 1990 the leading paradigm of object recognition was based on geometric representations [190, 215], with the focus later moving away from geometry and prior models towards the use of statistical classifiers (such as Neural Networks [233], SVM [201] and Adaboost [276, 290]) based on appearance features [191, 236]. This successful family of object detectors set the stage for most subsequent research in this field.

The milestones of object detection in more recent years are presented in Fig. 4, in which two main eras (SIFT *vs.* DCNN) are highlighted. The appearance features moved from global representations [192, 260, 267] to local representations that are designed to be invariant to changes in translation, scale, rotation, illumination, viewpoint and occlusion. Handcrafted local invariant features gained tremendous popularity, starting from the Scale Invariant Feature Transform (SIFT) feature [178], and the progress on various visual recognition tasks was based substantially on the use of local descriptors [187] such as Haar-like features [276], SIFT [179], Shape Contexts [12], Histogram of Gradients (HOG) [52] Local Binary Patterns (LBP) [196], and region covariances [268]. These local features are usually aggregated by simple concatenation or feature pooling encoders such as the Bag of Visual Words approach, introduced by Sivic and Zisserman [252] and Csurka *et al.* [47], Spatial Pyramid Matching (SPM) of BoW models [147], and Fisher Vectors [212].

For years, the multistage hand tuned pipelines of handcrafted local descriptors and discriminative classifiers dominated a variety of domains in computer vision, including object detection, until the significant turning point in 2012 when DCNNs [140] achieved their record-breaking results in image classification.

The use of CNNs for detection and localization [233] can be traced back to the 1990s, with a modest number of hidden layers used for object detection [272, 233, 238], successful in restricted domains such as face detection. However, more recently, deeper CNNs have led to record-breaking improvements in the detection of more general object categories, a shift which came about when the successful application of DCNNs in image classification [140] was transferred to object detection, resulting in the milestone Region-based CNN (RCNN) detector of Girshick *et al.* [85].

The successes of deep detectors rely heavily on vast training data and large networks with millions or even billions of parameters. The availability of GPUs with very high computational capability and large-scale detection datasets (such as ImageNet [54, 234] and MS COCO [166]) play a key role in their success. Large datasets have allowed researchers to target more realistic and complex problems from images with large intra-class variations and inter-class similarities [166, 234]. However, accurate annotations are labor intensive to obtain, so detectors must consider methods that can relieve annotation difficulties or can learn with smaller training datasets.

The research community has started moving towards the challenging goal of building general purpose object detection systems whose ability to detect many object categories matches that of humans. This is a major challenge: according to cognitive scientists,
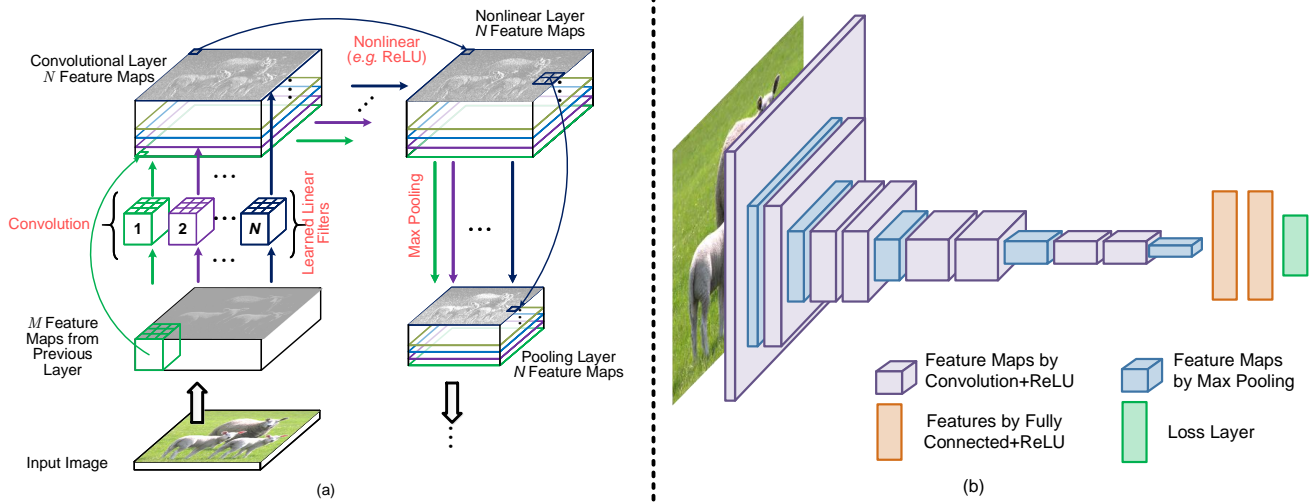
**Fig. 8** (a) Illustration of three operations that are repeatedly applied by a typical CNN: Convolution with a number of linear filters; Nonlinearities (*e.g.* ReLU); and Local pooling (*e.g.* Max Pooling). The $M$ feature maps from a previous layer are convolved with $N$ different filters (here shown as size $3 \times 3 \times M$), using a stride of 1. The resulting $N$ feature maps are then passed through a nonlinear function (*e.g.* ReLU), and pooled (*e.g.* taking a maximum over $2 \times 2$ regions) to give $N$ feature maps at a reduced resolution. (b) Illustration of the architecture of VGGNet [248], a typical CNN with 11 weight layers. An image with 3 color channels is presented as the input. The network has 8 convolutional layers, 3 fully connected layers, 5 max pooling layers and a softmax classification layer. The last three fully connected layers take features from the top convolutional layer as input in vector form. The final layer is a $C$-way softmax function, $C$ being the number of classes. The whole network can be learned from labeled training data by optimizing an objective function (*e.g.* mean squared error or cross entropy loss) via Stochastic Gradient Descent.

human beings can identify around 3,000 entry level categories and 30,000 visual categories overall, and the number of categories distinguishable with domain expertise may be to the order of $10^5$ [15]. Despite the remarkable progress of the past years, designing an accurate, robust, efficient detection and recognition system that approaches human-level performance on $10^4 - 10^5$ categories is undoubtedly an unresolved problem.

## 3 A Brief Introduction to Deep Learning

Deep learning has revolutionized a wide range of machine learning tasks, from image classification and video processing to speech recognition and natural language understanding. Given this tremendously rapid evolution, there exist many recent survey papers on deep learning [13, 89, 92, 149, 170, 216, 287, 297, 313, 320, 325]. These surveys have reviewed deep learning techniques from different perspectives [13, 89, 92, 149, 216, 287, 320], or with applications to medical image analysis [170], natural language processing [297], speech recognition systems [313], and remote sensing [325].

Convolutional Neural Networks (CNNs), the most representative models of deep learning, are able to exploit the basic properties underlying natural signals: translation invariance, local connectivity, and compositional hierarchies [149]. A typical CNN, illustrated in Fig. 8, has a hierarchical structure and is composed of a number of layers to learn representations of data with multiple levels of abstraction [149]. We begin with a convolution

$$\boldsymbol{x}^{l-1} * \boldsymbol{w}^l \qquad (1)$$

between an input feature map $\boldsymbol{x}^{l-1}$ at a feature map from previous layer $l - 1$, convolved with a 2D convolutional kernel (or filter or weights) $\boldsymbol{w}^l$. This convolution appears over a sequence of layers,

subject to a nonlinear operation $\sigma$, such that

$$\boldsymbol{x}^l_j = \sigma\left(\sum_{i=1}^{N^{l-1}} \boldsymbol{x}^{l-1}_i * \boldsymbol{w}^l_{i,j} + b^l_j\right), \qquad (2)$$

with a convolution now between the $N^{l-1}$ input feature maps $\boldsymbol{x}^{l-1}_i$ and the corresponding kernel $\boldsymbol{w}^l_{i,j}$, plus a bias term $b^l_j$. The elementwise nonlinear function $\sigma(\cdot)$ is typically a rectified linear unit (ReLU) for each element,

$$\sigma(x) = \max\{x, 0\}. \qquad (3)$$

Finally, pooling corresponds to the downsampling/upsampling of feature maps. These three operations (convolution, nonlinearity, pooling) are illustrated in Fig. 8 (a); CNNs having a large number of layers, a "deep" network, are referred to as Deep CNNs (DCNNs), with a typical DCNN architecture illustrated in Fig. 8 (b).

Most layers of a CNN consist of a number of feature maps, within which each pixel acts like a neuron. Each neuron in a convolutional layer is connected to feature maps of the previous layer through a set of weights $\boldsymbol{w}_{i,j}$ (essentially a set of 2D filters). As can be seen in Fig. 8 (b), where the early CNN layers are typically composed of convolutional and pooling layers, the later layers are normally fully connected. From earlier to later layers, the input image is repeatedly convolved, and with each layer, the receptive field or region of support increases. In general, the initial CNN layers extract low-level features (*e.g.,* edges), with later layers extracting more general features of increasing complexity [303, 13, 149, 199].

DCNNs have a number of outstanding advantages: a hierarchical structure to learn representations of data with multiple levels of abstraction, the capacity to learn very complex functions, and learning feature representations directly and automatically from data with minimal domain knowledge. What has particularly made

**Table 2** Most frequent object classes for each detection challenge. The size of each word is proportional to the frequency of that class in the training dataset.


(a) PASCAL VOC (20 Classes)


(b) MS COCO (80 Classes)


(c) ILSVRC (200 Classes)


(d) Open Images Detection Challenge (500 Classes)

DCNNs successful has been the availability of large scale labeled datasets and of GPUs with very high computational capability.

Despite the great successes, known deficiencies remain. In particular, there is an extreme need for labeled training data and a requirement of expensive computing resources, and considerable skill and experience are still needed to select appropriate learning parameters and network architectures. Trained networks are poorly interpretable, there is a lack of robustness to degradations, and many DCNNs have shown serious vulnerability to attacks [88], all of which currently limit the use of DCNNs in real-world applications.

## 4 Datasets and Performance Evaluation

### 4.1 Datasets

Datasets have played a key role throughout the history of object recognition research, not only as a common ground for measuring and comparing the performance of competing algorithms, but also pushing the field towards increasingly complex and challenging problems. In particular, recently, deep learning techniques have brought tremendous success to many visual recognition problems,

and it is the large amounts of annotated data which play a key role in their success. Access to large numbers of images on the Internet makes it possible to build comprehensive datasets in order to capture a vast richness and diversity of objects, enabling unprecedented performance in object recognition.

For generic object detection, there are four famous datasets: PASCAL VOC [68, 69], ImageNet [54], MS COCO [166] and Open Images [143]. The attributes of these datasets are summarized in Table 3, and selected sample images are shown in Fig. 9. There are three steps to creating large-scale annotated datasets: determining the set of target object categories, collecting a diverse set of candidate images to represent the selected categories on the Internet, and annotating the collected images, typically by designing crowdsourcing strategies. Recognizing space limitations, we refer interested readers to the original papers [68, 69, 166, 234, 143] for detailed descriptions of these datasets in terms of construction and properties.

The four datasets form the backbone of their respective detection challenges. Each challenge consists of a publicly available dataset of images together with ground truth annotation and standardized evaluation software, and an annual competition and corresponding workshop. Statistics for the number of images and object instances in the training, validation and testing datasets[2] for the detection challenges are given in Table 4. The most frequent object classes in VOC, COCO, ILSVRC and Open Images detection datasets are visualized in Table 2.

**PASCAL VOC** [68, 69] is a multi-year effort devoted to the creation and maintenance of a series of benchmark datasets for classification and object detection, creating the precedent for standardized evaluation of recognition algorithms in the form of annual competitions. Starting from only four categories in 2005, the dataset has increased to 20 categories that are common in everyday life. Since 2009, the number of images has grown every year, but with all previous images retained to allow test results to be compared from year to year. Due the availability of larger datasets like ImageNet, MS COCO and Open Images, PASCAL VOC has gradually fallen out of fashion.

**ILSVRC**, the ImageNet Large Scale Visual Recognition Challenge [234], is derived from ImageNet [54], scaling up PASCAL VOC's goal of standardized training and evaluation of detection algorithms by more than an order of magnitude in the number of object classes and images. ImageNet1000, a subset of ImageNet images with 1000 different object categories and a total of 1.2 million images, has been fixed to provide a standardized benchmark for the ILSVRC image classification challenge.

**MS COCO** is a response to the criticism of ImageNet that objects in its dataset tend to be large and well centered, making the ImageNet dataset atypical of real-world scenarios. To push for richer image understanding, researchers created the MS COCO database [166] containing complex everyday scenes with common objects in their natural context, closer to real life, where objects are labeled using fully-segmented instances to provide more accurate detector evaluation. The COCO object detection challenge [166] features two object detection tasks: using either bounding

---

[2] The annotations on the test set are not publicly released, except for PASCAL VOC2007.

**Table 3** Popular databases for object recognition. Example images from PASCAL VOC, ImageNet, MS COCO and Open Images are shown in Fig. 9.

| Dataset Name | Total Images | Categories | Images Per Category | Objects Per Image | Image Size | Started Year | Highlights |
|---|---|---|---|---|---|---|---|
| PASCAL VOC (2012) [69] | 11,540 | 20 | 303 ∼ 4087 | 2.4 | 470 × 380 | 2005 | Covers only 20 categories that are common in everyday life; Large number of training images; Close to real-world applications; Significantly larger intraclass variations; Objects in scene context; Multiple objects in one image; Contains many difficult samples. |
| ImageNet [234] | 14 millions+ | 21,841 | – | 1.5 | 500 × 400 | 2009 | Large number of object categories; More instances and more categories of objects per image; More challenging than PASCAL VOC; Backbone of the ILSVRC challenge; Images are object-centric. |
| MS COCO [166] | 328,000+ | 91 | – | 7.3 | 640 × 480 | 2014 | Even closer to real world scenarios; Each image contains more instances of objects and richer object annotation information; Contains object segmentation notation data that is not available in the ImageNet dataset. |
| Places [319] | 10 millions+ | 434 | – | – | 256 × 256 | 2014 | The largest labeled dataset for scene recognition; Four subsets Places365 Standard, Places365 Challenge, Places 205 and Places88 as benchmarks. |
| Open Images [143] | 9 millions+ | 6000+ | – | 8.3 | varied | 2017 | Annotated with image level labels, object bounding boxes and visual relationships; Open Images V5 supports large scale object detection, object instance segmentation and visual relationship detection. |



(a) PASCAL VOC

(b) ILSVRC

(c) MS COCO

(d) Open Images Detection

**Fig. 9** Some example images with object annotations from PASCAL VOC, ILSVRC, MS COCO and Open Images. See Table 3 for a summary of these datasets.

box output or object instance segmentation output. COCO introduced three new challenges:

1. It contains objects at a wide range of scales, including a high percentage of small objects [249];
2. Objects are less iconic and amid clutter or heavy occlusion;
3. The evaluation metric (see Table 5) encourages more accurate object localization.

Just like ImageNet in its time, MS COCO has become the standard for object detection today.

**OICOD** (the Open Image Challenge Object Detection) is derived from Open Images V4 (now V5 in 2019) [143], currently the largest publicly available object detection dataset. OICOD is different from previous large scale object detection datasets like ILSVRC and MS COCO, not merely in terms of the significantly increased number of classes, images, bounding box annotations and instance segmentation mask annotations, but also regarding the annotation process. In ILSVRC and MS COCO, instances of all classes in the dataset are exhaustively annotated, whereas for Open Images V4 a classifier was applied to each image and only those labels with sufficiently high scores were sent for human verification. Therefore in OICOD only the object instances of human-confirmed positive labels are annotated.

## 4.2 Evaluation Criteria

There are three criteria for evaluating the performance of detection algorithms: detection speed in Frames Per Second (FPS), preci-

sion, and recall. The most commonly used metric is *Average Precision* (AP), derived from precision and recall. AP is usually evaluated in a category specific manner, *i.e.*, computed for each object category separately. To compare performance over all object categories, the *mean AP* (mAP) averaged over all object categories is adopted as the final measure of performance[3]. More details on these metrics can be found in [68, 69, 234, 108].

The standard outputs of a detector applied to a testing image **I** are the predicted detections $\{(b_j, c_j, p_j)\}_j$, indexed by object $j$, of Bounding Box (BB) $b_j$, predicted category $c_j$, and confidence $p_j$. A predicted detection $(b, c, p)$ is regarded as a True Positive (TP) if

- The predicted category $c$ equals the ground truth label $c_g$.
- The overlap ratio IOU (Intersection Over Union) [68, 234]

$$\text{IOU}(b, b^g) = \frac{area(b \cap b^g)}{area(b \cup b^g)}, \tag{4}$$

between the predicted BB $b$ and the ground truth $b^g$ is not smaller than a predefined threshold $\varepsilon$, where $\cap$ and $cup$ denote intersection and union, respectively. A typical value of $\varepsilon$ is 0.5.

---

[3] In object detection challenges, such as PASCAL VOC and ILSVRC, the winning entry of each object category is that with the highest AP score, and the winner of the challenge is the team that wins on the most object categories. The mAP is also used as the measure of a team's performance, and is justified since the ranking of teams by mAP was always the same as the ranking by the number of object categories won [234].

**Table 4** Statistics of commonly used object detection datasets. Object statistics for VOC challenges list the non-difficult objects used in the evaluation (all annotated objects). For the COCO challenge, prior to 2017, the test set had four splits (*Dev*, *Standard*, *Reserve*, and *Challenge*), with each having about 20K images. Starting in 2017, the test set has only the *Dev* and *Challenge* splits, with the other two splits removed. Starting in 2017, the train and val sets are arranged differently, and the test set is divided into two roughly equally sized splits of about 20,000 images each: Test Dev and Test Challenge. Note that the 2017 Test Dev/Challenge splits contain the same images as the 2015 Test Dev/Challenge splits, so results across the years are directly comparable.

| Challenge | Object Classes | Number of Images | | | Number of Annotated Objects | | Summary (Train+Val) | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Train | Val | Test | Train | Val | Images | Boxes | Boxes/Image |
| PASCAL VOC Object Detection Challenge | | | | | | | | | |
| VOC07 | 20 | 2,501 | 2,510 | 4,952 | 6,301(7,844) | 6,307(7,818) | 5,011 | 12,608 | 2.5 |
| VOC08 | 20 | 2,111 | 2,221 | 4,133 | 5,082(6,337) | 5,281(6,347) | 4,332 | 10,364 | 2.4 |
| VOC09 | 20 | 3,473 | 3,581 | 6,650 | 8,505(9,760) | 8,713(9,779) | 7,054 | 17,218 | 2.3 |
| VOC10 | 20 | 4,998 | 5,105 | 9,637 | 11,577(13,339) | 11,797(13,352) | 10,103 | 23,374 | 2.4 |
| VOC11 | 20 | 5,717 | 5,823 | 10,994 | 13,609(15,774) | 13,841(15,787) | 11,540 | 27,450 | 2.4 |
| VOC12 | 20 | 5,717 | 5,823 | 10,991 | 13,609(15,774) | 13,841(15,787) | 11,540 | 27,450 | 2.4 |
| ILSVRC Object Detection Challenge | | | | | | | | | |
| ILSVRC13 | 200 | 395,909 | 20,121 | 40,152 | 345,854 | 55,502 | 416,030 | 401,356 | 1.0 |
| ILSVRC14 | 200 | 456,567 | 20,121 | 40,152 | 478,807 | 55,502 | 476,668 | 534,309 | 1.1 |
| ILSVRC15 | 200 | 456,567 | 20,121 | 51,294 | 478,807 | 55,502 | 476,668 | 534,309 | 1.1 |
| ILSVRC16 | 200 | 456,567 | 20,121 | 60,000 | 478,807 | 55,502 | 476,668 | 534,309 | 1.1 |
| ILSVRC17 | 200 | 456,567 | 20,121 | 65,500 | 478,807 | 55,502 | 476,668 | 534,309 | 1.1 |
| MS COCO Object Detection Challenge | | | | | | | | | |
| MS COCO15 | 80 | 82,783 | 40,504 | 81,434 | 604,907 | 291,875 | 123,287 | 896,782 | 7.3 |
| MS COCO16 | 80 | 82,783 | 40,504 | 81,434 | 604,907 | 291,875 | 123,287 | 896,782 | 7.3 |
| MS COCO17 | 80 | 118,287 | 5,000 | 40,670 | 860,001 | 36,781 | 123,287 | 896,782 | 7.3 |
| MS COCO18 | 80 | 118,287 | 5,000 | 40,670 | 860,001 | 36,781 | 123,287 | 896,782 | 7.3 |
| Open Images Challenge Object Detection (OICOD) (Based on Open Images V4 [143]) | | | | | | | | | |
| OICOD18 | 500 | 1,643,042 | 100,000 | 99,999 | 11,498,734 | 696,410 | 1,743,042 | 12,195,144 | 7.0 |

**Input**: $\{(b_j, p_j)\}_{j=1}^M$: $M$ predictions for image $\mathbf{I}$ for object class $c$, ranked by the confidence $p_j$ in decreasing order;
$\quad\quad \mathcal{B} = \{b_k^g\}_{k=1}^K$: ground truth BBs on image $\mathbf{I}$ for object class $c$;
**Output**: $\boldsymbol{a} \in \mathbb{R}^M$: a binary vector indicating each $(b_j, p_j)$ to be a TP or FP.
Initialize $\boldsymbol{a} = 0$;
**for** $j = 1, ..., M$ **do**
$\quad$ Set $\mathcal{A} = \varnothing$ and $t = 0$;
$\quad$ **foreach** *unmatched object* $b_k^g$ *in* $\mathcal{B}$ **do**
$\quad\quad$ **if** $IOU(b_j, b_k^g) \geq \varepsilon$ *and* $IOU(b_j, b_k^g) > t$ **then**
$\quad\quad\quad$ $\mathcal{A} = \{b_k^g\}$;
$\quad\quad\quad$ $t = IOU(b_j, b_k^g)$;
$\quad\quad$ **end**
$\quad$ **end**
$\quad$ **if** $\mathcal{A} \neq \varnothing$ **then**
$\quad\quad$ Set $\boldsymbol{a}(i) = 1$ since object prediction $(b_j, p_j)$ is a TP;
$\quad\quad$ Remove the matched GT box in $\mathcal{A}$ from $\mathcal{B}$, $\mathcal{B} = \mathcal{B} - \mathcal{A}$.
$\quad$ **end**
**end**

**Fig. 10** The algorithm for determining TPs and FPs by greedily matching object detection results to ground truth boxes.

Otherwise, it is considered as a False Positive (FP). The confidence level $p$ is usually compared with some threshold $\beta$ to determine whether the predicted class label $c$ is accepted.

AP is computed separately for each of the object classes, based on *Precision* and *Recall*. For a given object class $c$ and a testing image $\mathbf{I}_i$, let $\{(b_{ij}, p_{ij})\}_{j=1}^M$ denote the detections returned by a detector, ranked by confidence $p_{ij}$ in decreasing order. Each detection $(b_{ij}, p_{ij})$ is either a TP or an FP, which can be determined via the algorithm[4] in Fig. 10. Based on the TP and FP detections, the precision $P(\beta)$ and recall $R(\beta)$ [68] can be computed as a function of the confidence threshold $\beta$, so by varying the confi-

dence threshold different pairs $(P, R)$ can be obtained, in principle allowing precision to be regarded as a function of recall, *i.e.* $P(R)$, from which the Average Precision (AP) [68, 234] can be found.

Since the introduction of MS COCO, more attention has been placed on the accuracy of the bounding box location. Instead of using a fixed IOU threshold, MS COCO introduces a few metrics (summarized in Table 5) for characterizing the performance of an object detector. For instance, in contrast to the traditional mAP computed at a single IoU of 0.5, $AP_{coco}$ is averaged across all object categories and multiple IOU values from 0.5 to 0.95 in steps of 0.05. Because 41% of the objects in MS COCO are small and 24% are large, metrics $AP_{coco}^{small}$, $AP_{coco}^{medium}$ and $AP_{coco}^{large}$ are also introduced. Finally, Table 5 summarizes the main metrics used in the PASCAL, ILSVRC and MS COCO object detection challenges, with metric modifications for the Open Images challenges proposed in [143].

## 5 Detection Frameworks

There has been steady progress in object feature representations and classifiers for recognition, as evidenced by the dramatic change from handcrafted features [276, 52, 72, 98, 275] to learned DCNN features [85, 203, 84, 229, 50]. In contrast, in terms of localization, the basic "sliding window" strategy [52, 74, 72] remains mainstream, although with some efforts to avoid exhaustive search [145, 271]. However, the number of windows is large and grows quadratically with the number of image pixels, and the need to search over multiple scales and aspect ratios further increases the search space. Therefore, the design of efficient and effective detection frameworks plays a key role in reducing this computational cost. Commonly adopted strategies include cascading, sharing feature computation, and reducing per-window computation.

---

[4] It is worth noting that for a given threshold $\beta$, multiple detections of the same object in an image are not considered as all correct detections, and only the detection with the highest confidence level is considered as a TP and the rest as FPs.

**Table 5** Summary of commonly used metrics for evaluating object detectors.

| Metric | Meaning | Definition and Description | | |
|---|---|---|---|---|
| TP | True Positive | A true positive detection, per Fig. 10. | | |
| FP | False Positive | A false positive detection, per Fig. 10. | | |
| $\beta$ | Confidence Threshold | A confidence threshold for computing $P(\beta)$ and $R(\beta)$. | | |
| $\varepsilon$ | IOU Threshold | VOC | Typically around 0.5 | |
| | | ILSVRC | $\min(0.5, \frac{wh}{(w+10)(h+10)})$; $w \times h$ is the size of a GT box. | |
| | | MS COCO | Ten IOU thresholds $\varepsilon \in \{0.5 : 0.05 : 0.95\}$ | |
| $P(\beta)$ | Precision | The fraction of correct detections out of the total detections returned by the detector with confidence of at least $\beta$. | | |
| $R(\beta)$ | Recall | The fraction of all $N_c$ objects detected by the detector having a confidence of at least $\beta$. | | |
| AP | Average Precision | Computed over the different levels of recall achieved by varying the confidence $\beta$. | | |
| mAP | mean Average Precision | VOC | AP at a single IOU and averaged over all classes. | |
| | | ILSVRC | AP at a modified IOU and averaged over all classes. | |
| | | MS COCO | • $AP_{coco}$: mAP averaged over ten IOUs: $\{0.5 : 0.05 : 0.95\}$; <br>• $AP_{coco}^{\text{IOU}=0.5}$: mAP at IOU=0.50 (PASCAL VOC metric); <br>• $AP_{coco}^{\text{IOU}=0.75}$: mAP at IOU=0.75 (strict metric); <br>• $AP_{coco}^{\text{small}}$: mAP for small objects of area smaller than $32^2$; <br>• $AP_{coco}^{\text{medium}}$: mAP for objects of area between $32^2$ and $96^2$; <br>• $AP_{coco}^{\text{large}}$: mAP for large objects of area bigger than $96^2$. | |
| AR | Average Recall | The maximum recall given a fixed number of detections per image, averaged over all categories and IOU thresholds. | | |
| AR | Average Recall | MS COCO | • $AR_{coco}^{\max=1}$: AR given 1 detection per image; <br>• $AR_{coco}^{\max=10}$: AR given 10 detection per image; <br>• $AR_{coco}^{\max=100}$: AR given 100 detection per image; <br>• $AR_{coco}^{\text{small}}$: AR for small objects of area smaller than $32^2$; <br>• $AR_{coco}^{\text{medium}}$: AR for objects of area between $32^2$ and $96^2$; <br>• $AR_{coco}^{\text{large}}$: AR for large objects of area bigger than $96^2$; | |

This section reviews detection frameworks, listed in Fig. 11 and Table 11, the milestone approaches appearing since deep learning entered the field, organized into two main categories:

a. Two stage detection frameworks, which include a preprocessing step for generating object proposals;
b. One stage detection frameworks, or region proposal free frameworks, having a single proposed method which does not separate the process of the detection proposal.

Sections 6 through 9 will discuss fundamental sub-problems involved in detection frameworks in greater detail, including DCNN features, detection proposals, and context modeling.

### 5.1 Region Based (Two Stage) Frameworks

In a region-based framework, category-independent region proposals[5] are generated from an image, CNN [140] features are extracted from these regions, and then category-specific classifiers are used to determine the category labels of the proposals. As can be observed from Fig. 11, DetectorNet [261], OverFeat [239], Multi-Box [67] and RCNN [85] independently and almost simultaneously proposed using CNNs for generic object detection.

**RCNN** [85]: Inspired by the breakthrough image classification results obtained by CNNs and the success of the selective search in region proposal for handcrafted features [271], Girshick *et al.* were among the first to explore CNNs for generic object detection and developed RCNN [85, 87], which integrates AlexNet [140] with a region proposal selective search [271]. As illustrated in detail

in Fig. 12, training an RCNN framework consists of multistage pipelines:

1. *Region proposal computation:* Class agnostic region proposals, which are candidate regions that might contain objects, are obtained via a selective search [271].
2. *CNN model finetuning:* Region proposals, which are cropped from the image and warped into the same size, are used as the input for fine-tuning a CNN model pre-trained using a large-scale dataset such as ImageNet. At this stage, all region proposals with $\geqslant 0.5$ IOU [6] overlap with a ground truth box are defined as positives for that ground truth box's class and the rest as negatives.
3. *Class specific SVM classifiers training:* A set of class-specific linear SVM classifiers are trained using fixed length features extracted with CNN, replacing the softmax classifier learned by fine-tuning. For training SVM classifiers, positive examples are defined to be the ground truth boxes for each class. A region proposal with less than 0.3 IOU overlap with all ground truth instances of a class is negative for that class. Note that the positive and negative examples defined for training the SVM classifiers are different from those for fine-tuning the CNN.
4. *Class specific bounding box regressor training:* Bounding box regression is learned for each object class with CNN features.

In spite of achieving high object detection quality, RCNN has notable drawbacks [84]:

1. Training is a multistage pipeline, slow and hard to optimize because each individual stage must be trained separately.
2. For SVM classifier and bounding box regressor training, it is expensive in both disk space and time, because CNN features need to be extracted from each object proposal in each image, posing great challenges for large scale detection, particularly with very deep networks, such as VGG16 [248].
3. Testing is slow, since CNN features are extracted per object proposal in each test image, without shared computation.

All of these drawbacks have motivated successive innovations, leading to a number of improved detection frameworks such as SPP-Net, Fast RCNN, Faster RCNN *etc.*, as follows.

**SPPNet** [99]: During testing, CNN feature extraction is the main bottleneck of the RCNN detection pipeline, which requires the extraction of CNN features from thousands of warped region proposals per image. As a result, He *et al.* [99] introduced traditional spatial pyramid pooling (SPP) [90, 147] into CNN architectures. Since convolutional layers accept inputs of arbitrary sizes, the requirement of fixed-sized images in CNNs is due only to the Fully Connected (FC) layers, therefore He *et al.* added an SPP layer on top of the last convolutional (CONV) layer to obtain features of fixed length for the FC layers. With this SPPNet, RCNN obtains a significant speedup without sacrificing any detection quality, because it only needs to run the convolutional layers *o*nce on the entire test image to generate fixed-length features for region proposals of arbitrary size. While SPPNet accelerates RCNN evaluation by orders of magnitude, it does not result in a comparable speedup of the detector training. Moreover, fine-tuning in SPPNet [99] is unable to update the convolutional layers before the SPP layer, which limits the accuracy of very deep networks.

---

[5] Object proposals, also called region proposals or detection proposals, are a set of candidate regions or bounding boxes in an image that may potentially contain an object. [27, 110]

[6] Please refer to Section 4.2 for the definition of IOU.
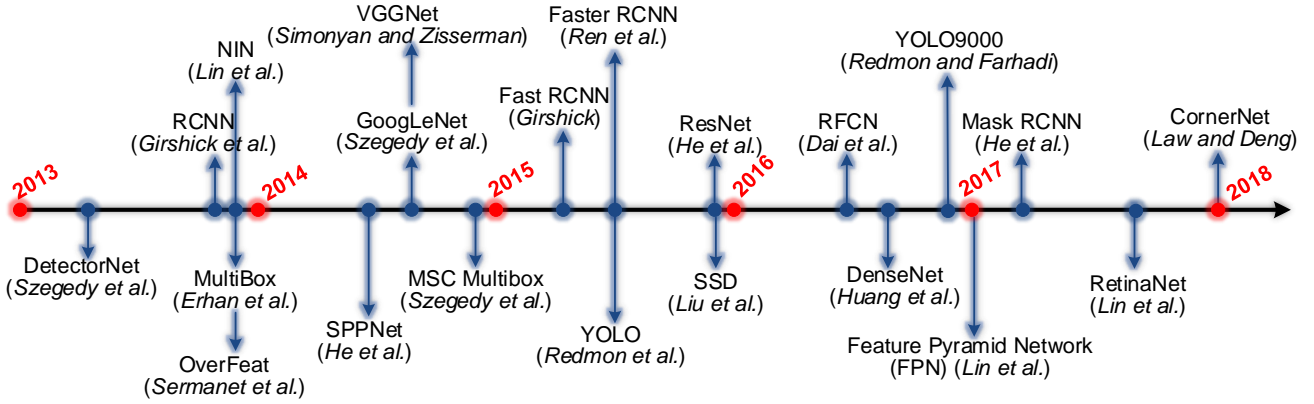
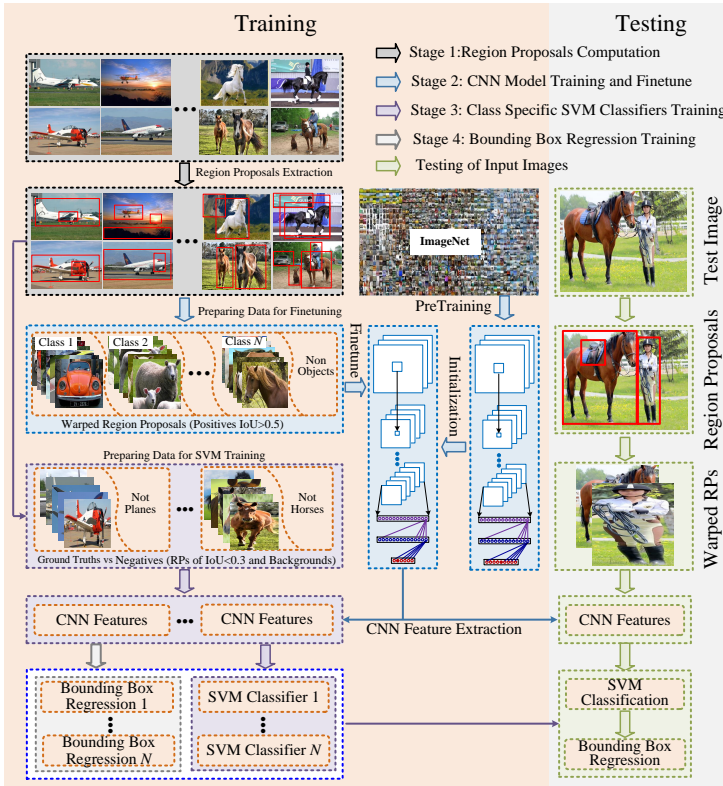**Fig. 11** Milestones in generic object detection.



**Fig. 12** Illustration of the RCNN detection framework [85, 87].

**Fast RCNN** [84]: Girshick proposed Fast RCNN [84] that addresses some of the disadvantages of RCNN and SPPNet, while improving on their detection speed and quality. As illustrated in Fig. 13, Fast RCNN enables end-to-end detector training by developing a streamlined training process that simultaneously learns a softmax classifier and class-specific bounding box regression, rather than separately training a softmax classifier, SVMs, and Bounding Box Regressors (BBRs) as in RCNN/SPPNet. Fast RCNN employs the idea of sharing the computation of convolution across region proposals, and adds a Region of Interest (RoI) pooling layer between the last CONV layer and the first FC layer to extract a fixed-length feature for each region proposal. Essentially, RoI pooling uses warping at the feature level to approximate warping at the image level. The features after the RoI pooling layer are fed into a sequence of FC layers that finally branch into two sibling output layers: softmax probabilities for object category prediction,

and class-specific bounding box regression offsets for proposal refinement. Compared to RCNN/SPPNet, Fast RCNN improves the efficiency considerably – typically 3 times faster in training and 10 times faster in testing. Thus there is higher detection quality, a single training process that updates all network layers, and no storage required for feature caching.

**Faster RCNN** [229, 230]: Although Fast RCNN significantly sped up the detection process, it still relies on external region proposals, whose computation is exposed as the new speed bottleneck in Fast RCNN. Recent work has shown that CNNs have a remarkable ability to localize objects in CONV layers [317, 318, 46, 200, 97], an ability which is weakened in the FC layers. Therefore, the selective search can be replaced by a CNN in producing region proposals. The Faster RCNN framework proposed by Ren *et al.* [229, 230] offered an efficient and accurate Region Proposal Network (RPN) for generating region proposals. They utilize the same backbone network, using features from the last shared convolutional layer to accomplish the task of RPN for region proposal and Fast RCNN for region classification, as shown in Fig. 13.

RPN first initializes $k$ reference boxes (*i.e.* the so called *anchors*) of different scales and aspect ratios at each CONV feature map location. The anchor *p*ositions are image content independent, but the feature vectors themselves, extracted from anchors, are image content dependent. Each anchor is mapped to a lower dimensional vector, which is fed into two sibling FC layers — an object category classification layer and a box regression layer. In contrast to detection in Fast RCNN, the features used for regression in RPN are of the same shape as the anchor box, thus $k$ anchors lead to $k$ regressors. RPN shares CONV features with Fast RCNN, thus enabling highly efficient region proposal computation. RPN is, in fact, a kind of Fully Convolutional Network (FCN) [177, 241]; Faster RCNN is thus a purely CNN based framework without using handcrafted features.

For the VGG16 model [248], Faster RCNN can test at 5 FPS (including all stages) on a GPU, while achieving state-of-the-art object detection accuracy on PASCAL VOC 2007 using 300 proposals per image. The initial Faster RCNN in [229] contains several alternating training stages, later simplified in [230].

Concurrent with the development of Faster RCNN, Lenc and Vedaldi [151] challenged the role of region proposal generation methods such as selective search, studied the role of region proposal generation in CNN based detectors, and found that CNNs contain sufficient geometric information for accurate object detec-
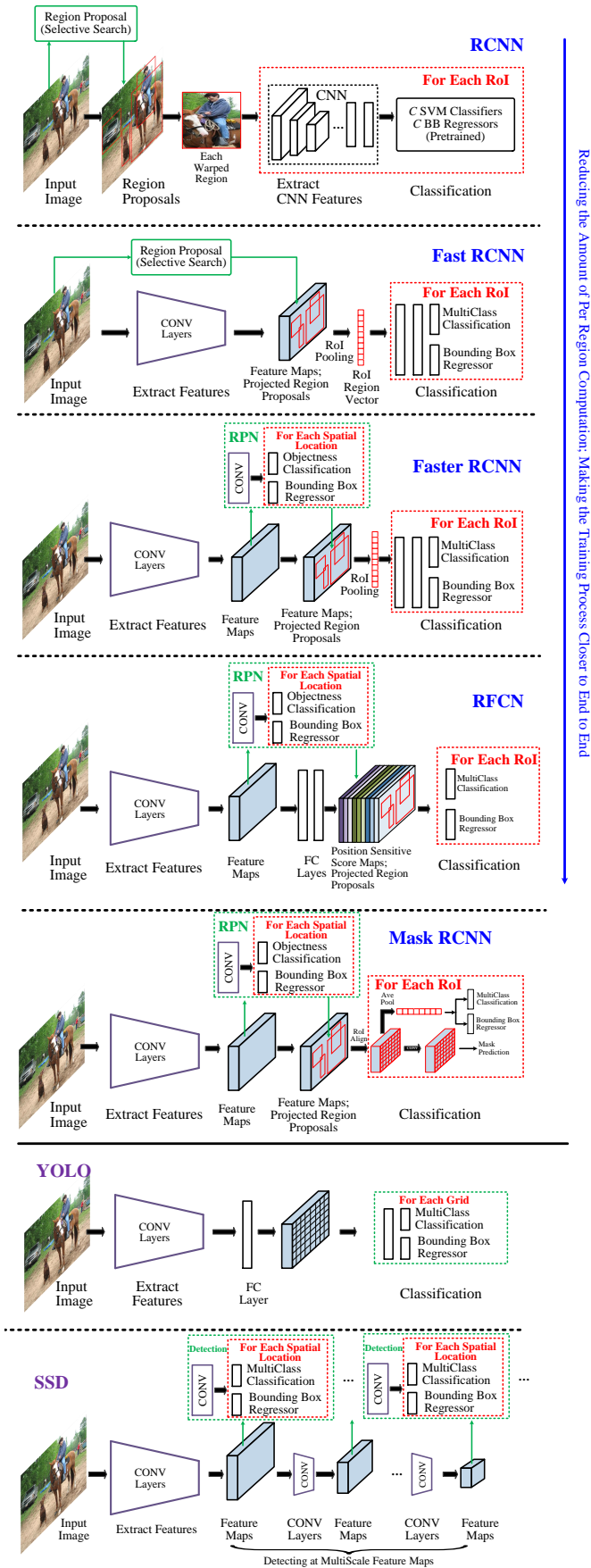
**Fig. 13** High level diagrams of the leading frameworks for generic object detection. The properties of these methods are summarized in Table 11.

tion in the CONV rather than FC layers. They showed the possibility of building integrated, simpler, and faster object detectors that rely exclusively on CNNs, removing region proposal generation methods such as selective search.

**RFCN (Region based Fully Convolutional Network)**: While Faster RCNN is an order of magnitude faster than Fast RCNN, the fact that the region-wise sub-network still needs to be applied per RoI (several hundred RoIs per image) led Dai *et al.* [50] to propose the RFCN detector which is *fully convolutional* (no hidden FC layers) with almost all computations shared over the entire image. As shown in Fig. 13, RFCN differs from Faster RCNN only in the RoI sub-network. In Faster RCNN, the computation after the RoI pooling layer cannot be shared, so Dai *et al.* [50] proposed using all CONV layers to construct a shared RoI sub-network, and RoI crops are taken from the last layer of CONV features prior to prediction. However, Dai *et al.* [50] found that this naive design turns out to have considerably inferior detection accuracy, conjectured to be that deeper CONV layers are more sensitive to category semantics, and less sensitive to translation, whereas object detection needs localization representations that respect translation invariance. Based on this observation, Dai *et al.* [50] constructed a set of position-sensitive score maps by using a bank of specialized CONV layers as the FCN output, on top of which a position-sensitive RoI pooling layer is added. They showed that RFCN with ResNet101 [101] could achieve comparable accuracy to Faster RCNN, often at faster running times.

**Mask RCNN**: He *et al.* [102] proposed Mask RCNN to tackle pixelwise object instance segmentation by extending Faster RCNN. Mask RCNN adopts the same two stage pipeline, with an identical first stage (RPN), but in the second stage, in parallel to predicting the class and box offset, Mask RCNN adds a branch which outputs a binary mask for each RoI. The new branch is a Fully Convolutional Network (FCN) [177, 241] on top of a CNN feature map. In order to avoid the misalignments caused by the original RoI pooling (RoIPool) layer, a RoIAlign layer was proposed to preserve the pixel level spatial correspondence. With a backbone network ResNeXt101-FPN [291, 167], Mask RCNN achieved top results for the COCO object instance segmentation and bounding box object detection. It is simple to train, generalizes well, and adds only a small overhead to Faster RCNN, running at 5 FPS [102].

**Chained Cascade Network and Cascade RCNN**: The essence of cascade [73, 20, 159] is to learn more discriminative classifiers by using multistage classifiers, such that early stages discard a large number of easy negative samples so that later stages can focus on handling more difficult examples. Two-stage object detection can be considered as a cascade, the first detector removing large amounts of background, and the second stage classifying the remaining regions. Recently, end-to-end learning of more than two cascaded classifiers and DCNNs for generic object detection were proposed in the Chained Cascade Network [205], extended in Cascade RCNN [23], and more recently applied for simultaneous object detection and instance segmentation [31], winning the COCO 2018 Detection Challenge.

**Light Head RCNN**: In order to further increase the detection speed of RFCN [50], Li *et al.* [165] proposed Light Head RCNN, making the head of the detection network as light as possible to reduce the RoI computation. In particular, Li *et al.* [165] applied a convolution to produce thin feature maps with small channel
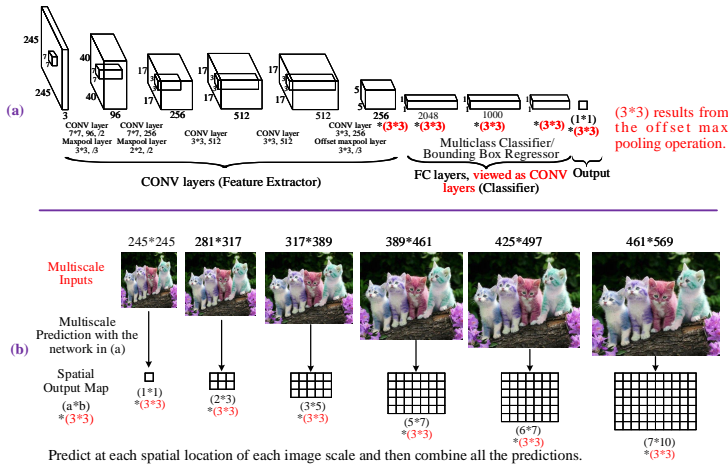
**Fig. 14** Illustration of the OverFeat [239] detection framework.

numbers (*e.g.,* 490 channels for COCO) and a cheap RCNN subnetwork, leading to an excellent trade-off of speed and accuracy.

## 5.2 Unified (One Stage) Frameworks

The region-based pipeline strategies of Section 5.1 have dominated since RCNN [85], such that the leading results on popular benchmark datasets are all based on Faster RCNN [229]. Nevertheless, region-based approaches are computationally expensive for current mobile/wearable devices, which have limited storage and computational capability, therefore instead of trying to optimize the individual components of a complex region-based pipeline, researchers have begun to develop *unified* detection strategies.

Unified pipelines refer to architectures that directly predict class probabilities and bounding box offsets from full images with a single feed-forward CNN in a monolithic setting that does not involve region proposal generation or post classification / feature resampling, encapsulating all computation in a single network. Since the whole pipeline is a single network, it can be optimized end-to-end directly on detection performance.

**DetectorNet**: Szegedy *et al.* [261] were among the first to explore CNNs for object detection. DetectorNet formulated object detection a regression problem to object bounding box masks. They use AlexNet [140] and replace the final softmax classifier layer with a regression layer. Given an image window, they use one network to predict foreground pixels over a coarse grid, as well as four additional networks to predict the object's top, bottom, left and right halves. A grouping process then converts the predicted masks into detected bounding boxes. The network needs to be trained per object type and mask type, and does not scale to multiple classes. DetectorNet must take many crops of the image, and run multiple networks for each part on every crop, thus making it slow.

**OverFeat**, proposed by Sermanet *et al.* [239] and illustrated in Fig. 14, can be considered as one of the first single-stage object detectors based on fully convolutional deep networks. It is one of the most influential object detection frameworks, winning the ILSVRC2013 localization and detection competition. OverFeat performs object detection via a single forward pass through the fully convolutional layers in the network (*i.e.* the "Feature Extrac-

tor", shown in Fig. 14 (a)). The key steps of object detection at test time can be summarized as follows:

1. *Generate object candidates by performing object classification via a sliding window fashion on multiscale images.* OverFeat uses a CNN like AlexNet [140], which would require input images of a fixed size due to its fully connected layers, in order to make the sliding window approach computationally efficient, OverFeat casts the network (as shown in Fig. 14 (a)) into a fully convolutional network, taking inputs of any size, by viewing fully connected layers as convolutions with kernels of size $1 \times 1$. OverFeat leverages multiscale features to improve the overall performance by passing up to six enlarged scales of the original image through the network (as shown in Fig. 14 (b)), resulting in a significantly increased number of evaluated context views. For each of the multiscale inputs, the classifier outputs a grid of predictions (class and confidence).

2. *Increase the number of predictions by offset max pooling.* In order to increase resolution, OverFeat applies offset max pooling after the last CONV layer, *i.e.* performing a subsampling operation at every offset, yielding many more views for voting, increasing robustness while remaining efficient.

3. *Bounding box regression.* Once an object is identified, a single bounding box regressor is applied. The classifier and the regressor share the same feature extraction (CONV) layers, only the FC layers need to be recomputed after computing the classification network.

4. *Combine predictions.* OverFeat uses a greedy merge strategy to combine the individual bounding box predictions across all locations and scales.

OverFeat has a significant speed advantage, but is less accurate than RCNN [85], because it was difficult to train fully convolutional networks at the time. The speed advantage derives from sharing the computation of convolution between overlapping windows in the fully convolutional network. OverFeat is similar to later frameworks such as YOLO [227] and SSD [175], except that the classifier and the regressors in OverFeat are trained sequentially.

**YOLO**: Redmon *et al.* [227] proposed YOLO (You Only Look Once), a unified detector casting object detection as a regression problem from image pixels to spatially separated bounding boxes and associated class probabilities, illustrated in Fig. 13. Since the region proposal generation stage is completely dropped, YOLO directly predicts detections using a small set of candidate regions[7]. Unlike region based approaches (*e.g.* Faster RCNN) that predict detections based on features from a local region, YOLO uses features from an entire image globally. In particular, YOLO divides an image into an $S \times S$ grid, each predicting $C$ class probabilities, $B$ bounding box locations, and confidence scores. By throwing out the region proposal generation step entirely, YOLO is fast by design, running in real time at 45 FPS and Fast YOLO [227] at 155 FPS. Since YOLO sees the entire image when making predictions, it implicitly encodes contextual information about object classes, and is less likely to predict false positives in the background. YOLO makes more localization errors than Fast RCNN, resulting from the coarse division of bounding box location, scale

---

[7] YOLO uses far fewer bounding boxes, only 98 per image, compared to about 2000 from Selective Search.

and aspect ratio. As discussed in [227], YOLO may fail to localize some objects, especially small ones, possibly because of the coarse grid division, and because each grid cell can only contain one object. It is unclear to what extent YOLO can translate to good performance on datasets with many objects per image, such as MS COCO.

**YOLOv2 and YOLO9000**: Redmon and Farhadi [226] proposed YOLOv2, an improved version of YOLO, in which the custom GoogLeNet [263] network is replaced with the simpler DarkNet19, plus batch normalization [100], removing the fully connected layers, and using good anchor boxes[8] learned via *k*means and multiscale training. YOLOv2 achieved state-of-the-art on standard detection tasks. Redmon and Farhadi [226] also introduced YOLO9000, which can detect over 9000 object categories in real time by proposing a joint optimization method to train simultaneously on an ImageNet classification dataset and a COCO detection dataset with WordTree to combine data from multiple sources. Such joint training allows YOLO9000 to perform weakly supervised detection, *i.e.* detecting object classes that do not have bounding box annotations.

**SSD**: In order to preserve real-time speed without sacrificing too much detection accuracy, Liu *et al.* [175] proposed SSD (Single Shot Detector), faster than YOLO [227] and with an accuracy competitive with region-based detectors such as Faster RCNN [229]. SSD effectively combines ideas from RPN in Faster RCNN [229], YOLO [227] and multiscale CONV features [97] to achieve fast detection speed, while still retaining high detection quality. Like YOLO, SSD predicts a fixed number of bounding boxes and scores, followed by an NMS step to produce the final detection. The CNN network in SSD is fully convolutional, whose early layers are based on a standard architecture, such as VGG [248], followed by several auxiliary CONV layers, progressively decreasing in size. The information in the last layer may be too coarse spatially to allow precise localization, so SSD performs detection over multiple scales by operating on multiple CONV feature maps, each of which predicts category scores and box offsets for bounding boxes of appropriate sizes. For a $300 \times 300$ input, SSD achieves 74.3% mAP on the VOC2007 test at 59 FPS versus Faster RCNN 7 FPS / mAP 73.2% or YOLO 45 FPS / mAP 63.4%.

**CornerNet:** Recently, Law *et al.* [146] questioned the dominant role that anchor boxes have come to play in SoA object detection frameworks [84, 102, 227, 175]. Law *et al.* [146] argue that the use of anchor boxes, especially in one stage detectors [77, 168, 175, 227], has drawbacks [146, 168] such as causing a huge imbalance between positive and negative examples, slowing down training and introducing extra hyperparameters. Borrowing ideas from the work on Associative Embedding in multiperson pose estimation [195], Law *et al.* [146] proposed CornerNet by formulating bounding box object detection as detecting paired top-left and bottom-right keypoints[9]. In CornerNet, the backbone network consists of two stacked Hourglass networks [194], with a simple corner pooling approach to better localize corners. CornerNet achieved a 42.1% AP on MS COCO, outperforming all previous one stage detectors; however, the average inference time is

about 4FPS on a Titan X GPU, significantly slower than SSD [175] and YOLO [227]. CornerNet generates incorrect bounding boxes because it is challenging to decide which pairs of keypoints should be grouped into the same objects. To further improve on CornerNet, Duan *et al.* [62] proposed CenterNet to detect each object as a triplet of keypoints, by introducing one extra keypoint at the centre of a proposal, raising the MS COCO AP to 47.0%, but with an inference speed slower than CornerNet.

## 6 Object Representation

As one of the main components in any detector, good feature representations are of primary importance in object detection [56, 85, 82, 324]. In the past, a great deal of effort was devoted to designing local descriptors (*e.g.,* SIFT [178] and HOG [52]) and to explore approaches (*e.g.,* Bag of Words [252] and Fisher Vector [212]) to group and abstract descriptors into higher level representations in order to allow the discriminative parts to emerge; however, these feature representation methods required careful engineering and considerable domain expertise.

In contrast, deep learning methods (especially *deep* CNNs) can learn powerful feature representations with multiple levels of abstraction directly from raw images [13, 149]. As the learning procedure reduces the dependency of specific domain knowledge and complex procedures needed in traditional feature engineering [13, 149], the burden for feature representation has been transferred to the design of better network architectures and training procedures.

The leading frameworks reviewed in Section 5 (RCNN [85], Fast RCNN [84], Faster RCNN [229], YOLO [227], SSD [175]) have persistently promoted detection accuracy and speed, in which it is generally accepted that the CNN architecture (Section 6.1 and Table 15) plays a crucial role. As a result, most of the recent improvements in detection accuracy have been via research into the development of novel networks. Therefore we begin by reviewing popular CNN architectures used in Generic Object Detection, followed by a review of the effort devoted to improving object feature representations, such as developing invariant features to accommodate geometric variations in object scale, pose, viewpoint, part deformation and performing multiscale analysis to improve object detection over a wide range of scales.

### 6.1 Popular CNN Architectures

CNN architectures (Section 3) serve as network backbones used in the detection frameworks of Section 5. Representative frameworks include AlexNet [141], ZFNet [303] VGGNet [248], GoogLeNet [263], Inception series [125, 264, 265], ResNet [101], DenseNet [118] and SENet [115], summarized in Table 6, and where the improvement over time is seen in Fig. 15. A further review of recent CNN advances can be found in [92].

The trend in architecture evolution is for greater depth: AlexNet has 8 layers, VGGNet 16 layers, more recently ResNet and DenseNet both surpassed the 100 layer mark, and it was VGGNet [248] and GoogLeNet [263] which showed that increasing depth can improve the representational power. As can be observed from Table 6, networks such as AlexNet, OverFeat, ZFNet and VGGNet have an

---

[8] Boxes of various sizes and aspect ratios that serve as object candidates.

[9] The idea of using keypoints for object detection appeared previously in DeNet [269].

**Table 6** DCNN architectures that were commonly used for generic object detection. Regarding the statistics for "#Paras" and "#Layers", the final FC prediction layer is not taken into consideration. "Test Error" column indicates the Top 5 classification test error on ImageNet1000. When ambiguous, the "#Paras", "#Layers", and "Test Error" refer to: OverFeat (accurate model), VGGNet16, ResNet101 DenseNet201 (Growth Rate 32, DenseNet-BC), ResNeXt50 (32*4d), and SE ResNet50.

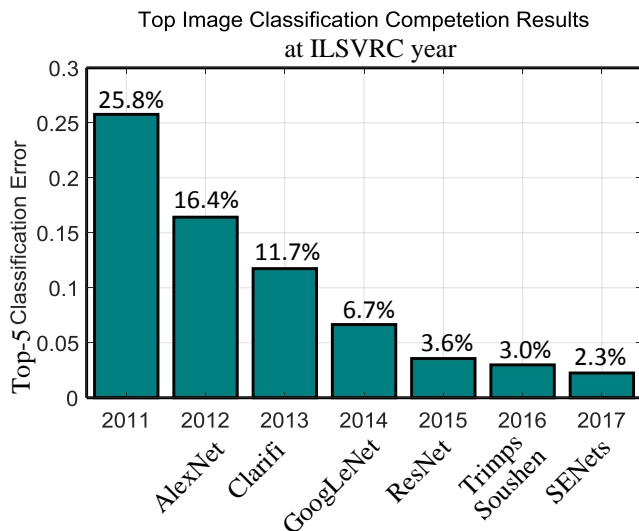| No. | DCNN Architecture | #Paras ($\times 10^6$) | #Layers (CONV+FC) | Test Error (Top 5) | First Used In | Highlights |
|---|---|---|---|---|---|---|
| 1 | AlexNet [141] | 57 | 5 + 2 | 15.3% | [85] | The first DCNN found effective for ImageNet classification; the historical turning point from hand-crafted features to CNN; Winning the ILSVRC2012 Image classification competition. |
| 2 | ZFNet (fast) [303] | 58 | 5 + 2 | 14.8% | [99] | Similar to AlexNet, different in stride for convolution, filter size, and number of filters for some layers. |
| 3 | OverFeat [239] | 140 | 6 + 2 | 13.6% | [239] | Similar to AlexNet, different in stride for convolution, filter size, and number of filters for some layers. |
| 4 | VGGNet [248] | 134 | 13 + 2 | 6.8% | [84] | Increasing network depth significantly by stacking $3 \times 3$ convolution filters and increasing the network depth step by step. |
| 5 | GoogLeNet [263] | 6 | 22 | 6.7% | [263] | Use Inception module, which uses multiple branches of convolutional layers with different filter sizes and then concatenates feature maps produced by these branches. The first inclusion of bottleneck structure and global average pooling. |
| 6 | Inception v2 [125] | 12 | 31 | 4.8% | [112] | Faster training with the introduce of Batch Normalization. |
| 7 | Inception v3 [264] | 22 | 47 | 3.6% | | Inclusion of separable convolution and spatial resolution reduction. |
| 8 | YOLONet [227] | 64 | 24 + 1 | − | [227] | A network inspired by GoogLeNet used in YOLO detector. |
| 9 | ResNet50 [101] | 23.4 | 49 | 3.6% | [101] | With identity mapping, substantially deeper networks can be learned. |
| 10 | ResNet101 [101] | 42 | 100 | (ResNets) | [101] | Requires fewer parameters than VGG by using the global average pooling and bottleneck introduced in GoogLeNet. |
| 11 | InceptionResNet v1 [265] | 21 | 87 | 3.1% | | Combination of identity mapping and Inception module, with similar computational cost of Inception v3, but faster training process. |
| 12 | InceptionResNet v2 [265] | 30 | 95 | (Ensemble) | [120] | A costlier residual version of Inception, with significantly improved recognition performance. |
| 13 | Inception v4 [265] | 41 | 75 | | | An Inception variant without residual connections, with roughly the same recognition performance as InceptionResNet v2, but significantly slower. |
| 14 | ResNeXt [291] | 23 | 49 | 3.0% | [291] | Repeating a building block that aggregates a set of transformations with the same topology. |
| 15 | DenseNet201 [118] | 18 | 200 | − | [321] | Concatenate each layer with every other layer in a feed forward fashion. Alleviate the vanishing gradient problem, encourage feature reuse, reduction in number of parameters. |
| 16 | DarkNet [226] | 20 | 19 | − | [226] | Similar to VGGNet, but with significantly fewer parameters. |
| 17 | MobileNet [112] | 3.2 | 27 + 1 | − | [112] | Light weight deep CNNs using depth-wise separable convolutions. |
| 18 | SE ResNet [115] | 26 | 50 | 2.3% (SENets) | [115] | Channel-wise attention by a novel block called *Squeeze and Excitation*. Complementary to existing backbone CNNs. |



**Fig. 15** Performance of winning entries in the ILSVRC competitions from 2011 to 2017 in the image classification task.

enormous number of parameters, despite being only a few layers deep, since a large fraction of the parameters come from the FC layers. Newer networks like Inception, ResNet, and DenseNet, although having a great depth, actually have far fewer parameters by avoiding the use of FC layers.

With the use of Inception modules [263] in carefully designed topologies, the number of parameters of GoogLeNet is dramatically reduced, compared to AlexNet, ZFNet or VGGNet. Similarly, ResNet demonstrated the effectiveness of skip connections for learning extremely deep networks with hundreds of layers, winning the ILSVRC 2015 classification task. Inspired by ResNet [101], InceptionResNets [265] combined the Inception networks with shortcut connections, on the basis that shortcut connections can significantly accelerate network training. Extending ResNets, Huang *et al.* [118] proposed DenseNets, which are built from dense blocksconnecting each layer to every other layer in a feedforward fashion, leading to compelling advantages such as parameter efficiency, implicit deep supervision[10], and feature reuse. Recently, Hu *et al.* [101] proposed Squeeze and Excitation (SE) blocks, which can

---

[10] DenseNets perform deep supervision in an implicit way, *i.e.* individual layers receive additional supervision from other layers through the shorter connections. The benefits of deep supervision have previously been demonstrated in Deeply Supervised Nets (DSN) [150].

be combined with existing deep architectures to boost their performance at minimal additional computational cost, adaptively recalibrating channel-wise feature responses by explicitly modeling the interdependencies between convolutional feature channels, and which led to winning the ILSVRC 2017 classification task. Research on CNN architectures remains active, with emerging networks such as Hourglass [146], Dilated Residual Networks [299], Xception [45], DetNet [164], Dual Path Networks (DPN) [37], FishNet [257], and GLoRe [38].

The training of a CNN requires a large-scale labeled dataset with intraclass diversity. Unlike image classification, detection requires localizing (possibly many) objects from an image. It has been shown [206] that pretraining a deep model with a large scale dataset having object level annotations (such as ImageNet), instead of only the image level annotations, improves the detection performance. However, collecting bounding box labels is expensive, especially for hundreds of thousands of categories. A common scenario is for a CNN to be pretrained on a large dataset (usually with a large number of visual categories) with image-level labels; the pretrained CNN can then be applied to a small dataset, directly, as a generic feature extractor [223, 8, 60, 296], which can support a wider range of visual recognition tasks. For detection, the pre-trained network is typically fine-tuned[11] on a given detection dataset [60, 85, 87]. Several large scale image classification datasets are used for CNN pre-training, among them ImageNet1000 [54, 234] with 1.2 million images of 1000 object categories, Places [319], which is much larger than ImageNet1000 but with fewer classes, a recent Places-Imagenet hybrid [319], or JFT300M [106, 254].

Pretrained CNNs without fine-tuning were explored for object classification and detection in [60, 87, 1], where it was shown that detection accuracies are different for features extracted from different layers; for example, for AlexNet pre-trained on ImageNet, FC6 / FC7 / Pool5 are in descending order of detection accuracy [60, 87]. Fine-tuning a pre-trained network can increase detection performance significantly [85, 87], although in the case of AlexNet, the fine-tuning performance boost was shown to be much larger for FC6 / FC7 than for Pool5, suggesting that Pool5 features are more general. Furthermore, the relationship between the source and target datasets plays a critical role, for example that ImageNet based CNN features show better performance for object detection than for human action [317, 8].

## 6.2 Methods For Improving Object Representation

Deep CNN based detectors such as RCNN [85], Fast RCNN [84], Faster RCNN [229] and YOLO [227], typically use the deep CNN architectures listed in Table 6 as the backbone network and use features from the top layer of the CNN as object representations; however, detecting objects across a large *range* of scales is a fundamental challenge. A classical strategy to address this issue is to run the detector over a number of scaled input images (*e.g.,* an image pyramid) [74, 85, 99], which typically produces more accurate
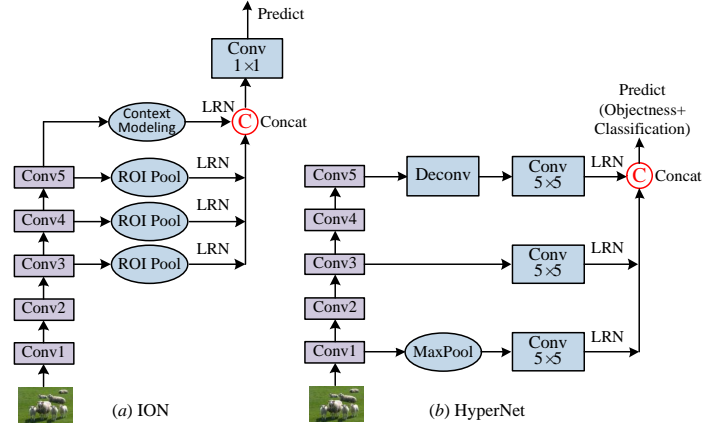
---

**Fig. 16** Comparison of HyperNet and ION. LRN is Local Response Normalization, which performs a kind of "lateral inhibition" by normalizing over local input regions [127].

detection, with, however, obvious limitations of inference time and memory.

### 6.2.1 Handling of Object Scale Variations

Since a CNN computes its feature hierarchy layer by layer, the subsampling layers in the feature hierarchy already lead to an inherent multiscale pyramid, producing feature maps at different spatial resolutions, but subject to challenges [97, 177, 247]. In particular, the higher layers have a large receptive field and strong semantics, and are the most robust to variations such as object pose, illumination and part deformation, but the resolution is low and the geometric details are lost. In contrast, lower layers have a small receptive field and rich geometric details, but the resolution is high and much less sensitive to semantics. Intuitively, semantic concepts of objects can emerge in different layers, depending on the size of the objects. So if a target object is small it requires fine detail information in earlier layers and may very well disappear at later layers, in principle making small object detection very challenging, for which tricks such as dilated or "atrous" convolution [298, 50, 33] have been proposed, increasing feature resolution, but increasing computational complexity. On the other hand, if the target object is large, then the semantic concept will emerge in much later layers. A number of methods [247, 314, 167, 136] have been proposed to improve detection accuracy by exploiting multiple CNN layers, broadly falling into three types of **multiscale object detection**:

1. Detecting with combined features of multiple layers;
2. Detecting at multiple layers;
3. Combinations of the above two methods.

**(1) Detecting with combined features of multiple CNN layers:** Many approaches, including Hypercolumns [97], HyperNet [135], and ION [11], combine features from multiple layers before making a prediction. Such feature combination is commonly accomplished via concatenation, a classic neural network idea that concatenates features from different layers, architectures which have recently become popular for semantic segmentation [177, 241, 97]. As shown in Fig. 16 (a), ION [11] uses RoI pooling to extract RoI features from multiple layers, and then the object proposals generated by selective search and edgeboxes are classified by using the concatenated features. HyperNet [135], shown in Fig. 16

**Table 7** Summary of properties of representative methods in improving DCNN feature representations for generic object detection. Details for Groups (1), (2), and (3) are provided in Section 6.2. Abbreviations: Selective Search (SS), EdgeBoxes (EB), InceptionResNet (IRN). *Conv-Deconv* denotes the use of upsampling and convolutional layers with lateral connections to supplement the standard backbone network. Detection results on VOC07, VOC12 and COCO were reported with mAP@IoU=0.5, and the additional COCO results are computed as the average of mAP for IoU thresholds from 0.5 to 0.95. Training data: "07"←VOC2007 trainval; "07T"←VOC2007 trainval and test; "12"←VOC2012 trainval; CO← COCO trainval. The COCO detection results were reported with COCO2015 Test-Dev, except for MPN [302] which reported with COCO2015 Test-Standard.

| Group | Detector Name | Region Proposal | Backbone DCNN | Pipelined Used | mAP@IoU=0.5 VOC07 | VOC12 | COCO | mAP COCO | Published In | Highlights |
|---|---|---|---|---|---|---|---|---|---|---|
| **(1) Single detection with multilayer features** | ION [11] | SS+EB MCG+RPN | VGG16 | Fast RCNN | 79.4 (07+12) | 76.4 (07+12) | 55.7 | 33.1 | CVPR16 | Use features from multiple layers; use spatial recurrent neural networks for modeling contextual information; the Best Student Entry and the 3$^{rd}$ overall in the COCO detection challenge 2015. |
| | HyperNet [135] | RPN | VGG16 | Faster RCNN | 76.3 (07+12) | 71.4 (07T+12) | – | – | CVPR16 | Use features from multiple layers for both region proposal and region classification. |
| | PVANet [132] | RPN | PVANet | Faster RCNN | **84.9** (07+12+CO) | **84.2** (07T+12+CO) | – | – | NIPSW16 | Deep but lightweight; Combine ideas from concatenated ReLU [240], Inception [263], and HyperNet [135]. |
| **(2) Detection at multiple layers** | SDP+CRC [293] | EB | VGG16 | Fast RCNN | 69.4 (07) | – | – | – | CVPR16 | Use features in multiple layers to reject easy negatives via CRC, and then classify remaining proposals using SDP. |
| | MSCNN [24] | RPN | VGG | Faster RCNN | Only Tested on KITTI | | | | ECCV16 | Region proposal and classification are performed at multiple layers; includes feature upsampling; end to end learning. |
| | MPN [302] | SharpMask [214] | VGG16 | Fast RCNN | – | – | 51.9 | 33.2 | BMVC16 | Concatenate features from different convolutional layers and features of different contextual regions; loss function for multiple overlap thresholds; ranked 2$^{nd}$ in both the COCO15 detection and segmentation challenges. |
| | DSOD [242] | Free | DenseNet | SSD | 77.7 (07+12) | 72.2 (07T+12) | 47.3 | 29.3 | ICCV17 | Concatenate feature sequentially, like DenseNet. Train from scratch on the target dataset without pre-training. |
| | RFBNet [173] | Free | VGG16 | SSD | 82.2 (07+12) | 81.2 (07T+12) | 55.7 | 34.4 | ECCV18 | Propose a multi-branch convolutional block similar to Inception [263], but using dilated convolution. |
| **(3) Combination of (1) and (2)** | DSSD [77] | Free | ResNet101 | SSD | 81.5 (07+12) | 80.0 (07T+12) | 53.3 | 33.2 | 2017 | Use Conv-Deconv, as shown in Fig. 17 (c1, c2). |
| | FPN [167] | RPN | ResNet101 | Faster RCNN | – | – | 59.1 | 36.2 | CVPR17 | Use Conv-Deconv, as shown in Fig. 17 (a1, a2); Widely used in detectors. |
| | TDM [247] | RPN | ResNet101 VGG16 | Faster RCNN | – | – | 57.7 | 36.8 | CVPR17 | Use Conv-Deconv, as shown in Fig. 17 (b2). |
| | RON [136] | RPN | VGG16 | Faster RCNN | 81.3 (07+12+CO) | 80.7 (07T+12+CO) | 49.5 | 27.4 | CVPR17 | Use Conv-deconv, as shown in Fig. 17 (d2); Add the objectness prior to significantly reduce object search space. |
| | ZIP [156] | RPN | Inceptionv2 | Faster RCNN | 79.8 (07+12) | – | – | – | IJCV18 | Use Conv-Deconv, as shown in Fig. 17 (f1). Propose a map attention decision (MAD) unit for features from different layers. |
| | STDN [321] | Free | DenseNet169 | SSD | 80.9 (07+12) | – | 51.0 | 31.8 | CVPR18 | A new scale transfer module, which resizes features of different scales to the same scale in parallel. |
| | RefineDet [308] | RPN | VGG16 ResNet101 | Faster RCNN | 83.8 (07+12) | 83.5 (07T+12) | 62.9 | 41.8 | CVPR18 | Use cascade to obtain better and less anchors. Use Conv-deconv, as shown in Fig. 17 (e2) to improve features. |
| | PANet [174] | RPN | ResNeXt101 +FPN | Mask RCNN | – | – | **67.2** | **47.4** | CVPR18 | Shown in Fig. 17 (g). Based on FPN, add another bottom-up path to pass information between lower and topmost layers; adaptive feature pooling. Ranked 1$^{st}$ and 2$^{nd}$ in COCO 2017 tasks. |
| | DetNet [164] | RPN | DetNet59+FPN | Faster RCNN | – | – | 61.7 | 40.2 | ECCV18 | Introduces dilated convolution into the ResNet backbone to maintain high resolution in deeper layers; Shown in Fig. 17 (i). |
| | FPR [137] | – | VGG16 ResNet101 | SSD | 82.4 (07+12) | 81.1 (07T+12) | 54.3 | 34.6 | ECCV18 | Fuse task oriented features across different spatial locations and scales, globally and locally; Shown in Fig. 17 (h). |
| | M2Det [315] | – | SSD | VGG16 ResNet101 | – | – | 64.6 | 44.2 | AAAI19 | Shown in Fig. 17 (j), newly designed top down path to learn a set of multilevel features, recombined to construct a feature pyramid for object detection. |
| **(4) Model Geometric Transforms** | DeepIDNet [203] | SS+ EB | AlexNet ZFNet OverFeat GoogLeNet | RCNN | 69.0 (07) | – | – | 25.6 | CVPR15 | Introduce a deformation constrained pooling layer, jointly learned with convolutional layers in existing DCNNs. Utilize the following modules that are not trained end to end: cascade, context modeling, model averaging, and bounding box location refinement in the multistage detection pipeline. |
| | DCN [51] | RPN | ResNet101 IRN | RFCN | 82.6 (07+12) | – | 58.0 | 37.5 | CVPR17 | Design deformable convolution and deformable RoI pooling modules that can replace plain convolution in existing DCNNs. |
| | DPFCN [188] | AttractioNet [83] | ResNet | RFCN | 83.3 (07+12) | 81.2 (07T+12) | 59.1 | 39.1 | IJCV18 | Design a deformable part based RoI pooling layer to explicitly select discriminative regions around object proposals. |

(b), follows a similar idea, and integrates deep, intermediate and shallow features to generate object proposals and to predict objects via an end to end joint training strategy. The combined feature is more descriptive, and is more beneficial for localization and classification, but at increased computational complexity.

**(2) Detecting at multiple CNN layers:** A number of recent approaches improve detection by predicting objects of different resolutions at different layers and then combining these predictions: SSD [175] and MSCNN [24], RBFNet [173], and DSOD [242]. SSD [175] spreads out default boxes of different scales to multiple layers within a CNN, and forces each layer to focus on predicting objects of a certain scale. RFBNet [173] replaces the later convolution layers of SSD with a Receptive Field Block (RFB) to enhance the discriminability and robustness of features. The RFB is a multibranch convolutional block, similar to the Inception block [263], but combining multiple branches with different kernels and convolution layers [33]. MSCNN [24] applies deconvolution on multiple layers of a CNN to increase feature map resolution before using the layers to learn region proposals and pool features. Similar to RFBNet [173], TridentNet [163] constructs a parallel multi-branch architecture where each branch shares the same transformation parameters but with different receptive fields; dilated convolution with different dilation rates are used to adapt the receptive fields for objects of different scales.

**(3) Combinations of the above two methods:** Features from different layers are complementary to each other and can improve
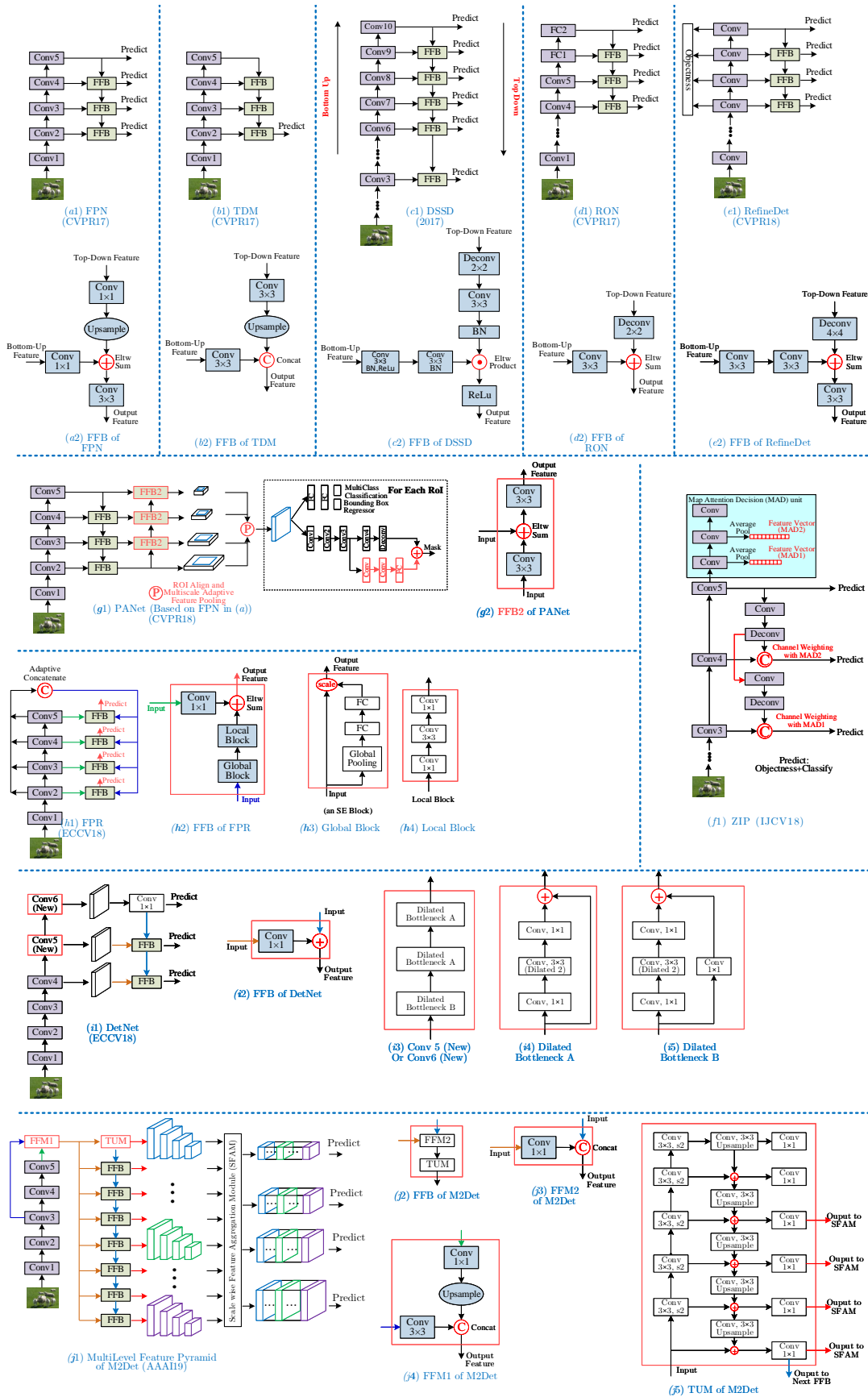
**Fig. 17** Hourglass architectures: Conv1 to Conv5 are the main Conv blocks in backbone networks such as VGG or ResNet. The figure compares a number of Feature Fusion Blocks (FFB) commonly used in recent approaches: FPN [167], TDM [247], DSSD [77], RON [136], RefineDet [308], ZIP [156], PANet [174], FPR [137], DetNet [164] and M2Det [315]. FFM: Feature Fusion Module, TUM: Thinned U-shaped Module

detection accuracy, as shown by Hypercolumns [97], HyperNet [135] and ION [11]. On the other hand, however, it is natural to detect objects of different scales using features of approximately the same size, which can be achieved by detecting large objects from downscaled feature maps while detecting small objects from upscaled feature maps. Therefore, in order to combine the best of both worlds, some recent works propose to detect objects at multiple layers, and the resulting features obtained by combining features from different layers. This approach has been found to be effective for segmentation [177, 241] and human pose estimation [194], has been widely exploited by both one-stage and two-stage detectors to alleviate problems of scale variation across object instances. Representative methods include SharpMask [214], Deconvolutional Single Shot Detector (DSSD) [77], Feature Pyramid Network (FPN) [167], Top Down Modulation (TDM)[247], Reverse connection with Objectness prior Network (RON) [136], ZIP [156], Scale Transfer Detection Network (STDN) [321], RefineDet [308], StairNet [283], Path Aggregation Network (PANet) [174], Feature Pyramid Reconfiguration (FPR) [137], DetNet [164], Scale Aware Network (SAN) [133], Multiscale Location aware Kernel Representation (MLKP) [278] and M2Det [315], as shown in Table 7 and contrasted in Fig. 17.

Early works like FPN [167], DSSD [77], TDM [247], ZIP [156], RON [136] and RefineDet [308] construct the feature pyramid according to the inherent multiscale, pyramidal architecture of the backbone, and achieved encouraging results. As can be observed from Fig. 17 (a1) to (f1), these methods have very similar detection architectures which incorporate a top-down network with lateral connections to supplement the standard bottom-up, feed-forward network. Specifically, after a bottom-up pass the final high level semantic features are transmitted back by the top-down network to combine with the bottom-up features from intermediate layers after lateral processing, and the combined features are then used for detection. As can be seen from Fig. 17 (a2) to (e2), the main differences lie in the design of the simple Feature Fusion Block (FFB), which handles the selection of features from different layers and the combination of multilayer features.

FPN [167] shows significant improvement as a generic feature extractor in several applications including object detection [167, 168] and instance segmentation [102]. Using FPN in a basic Faster RCNN system achieved state-of-the-art results on the COCO detection dataset. STDN [321] used DenseNet [118] to combine features of different layers and designed a scale transfer module to obtain feature maps with different resolutions. The scale transfer module can be directly embedded into DenseNet with little additional cost.

More recent work, such as PANet [174], FPR [137], DetNet [164], and M2Det [315], as shown in Fig. 17 (g-j), propose to further improve on the pyramid architectures like FPN in different ways. Based on FPN, Liu *et al.* designed PANet [174] (Fig. 17 (g1)) by adding another bottom-up path with clean lateral connections from low to top levels, in order to shorten the information path and to enhance the feature pyramid. Then, an adaptive feature pooling was proposed to aggregate features from all feature levels for each proposal. In addition, in the proposal sub-network, a complementary branch capturing different views for each proposal is created to further improve mask prediction. These additional steps bring only slightly extra computational overhead, but are effective and allowed PANet to reach 1st place in the COCO 2017 Challenge Instance Segmentation task and 2nd place in the Object Detection task. Kong *et al.* proposed FPR [137] by explicitly reformulating the feature pyramid construction process (*e.g.* FPN [167]) as feature reconfiguration functions in a highly nonlinear but efficient way. As shown in Fig. 17 (h1), instead of using a top-down path to propagate strong semantic features from the topmost layer down as in FPN, FPR first extracts features from multiple layers in the backbone network by adaptive concatenation, and then designs a more complex FFB module (Fig. 17 (h2)) to spread strong semantics to all scales. Li *et al.* proposed DetNet [164] (Fig. 17 (i1)) by introducing dilated convolutions to the later layers of the backbone network in order to maintain high spatial resolution in deeper layers. Zhao *et al.* [315] proposed a MultiLevel Feature Pyramid Network (MLFPN) to build more effective feature pyramids for detecting objects of different scales. As can be seen from Fig. 17 (j1), features from two different layers of the backbone are first fused as the base feature, after which a top-down path with lateral connections from the base feature is created to build the feature pyramid. As shown in Fig. 17 (j2) and (j5), the FFB module is much more complex than those like FPN, in that FFB involves a Thinned U-shaped Module (TUM) to generate a second pyramid structure, after which the feature maps with equivalent sizes from multiple TUMs are combined for object detection. The authors proposed M2Det by integrating MLFPN into SSD, and achieved better detection performance than other one-stage detectors.

### 6.3 Handling of Other Intraclass Variations

Powerful object representations should combine distinctiveness and robustness. A large amount of recent work has been devoted to handling changes in object scale, as reviewed in Section 6.2.1. As discussed in Section 2.2 and summarized in Fig. 6, object detection still requires robustness to real-world variations other than just scale, which we group into three categories:

- Geometric transformations,
- Occlusions, and
- Image degradations.

To handle these intra-class variations, the most straightforward approach is to augment the training datasets with a sufficient amount of variations; for example, robustness to rotation could be achieved by adding rotated objects at many orientations to the training data. Robustness can frequently be learned this way, but usually at the cost of expensive training and complex model parameters. Therefore, researchers have proposed alternative solutions to these problems.

**Handling of geometric transformations:** DCNNs are inherently limited by the lack of ability to be spatially invariant to geometric transformations of the input data [152, 172, 28]. The introduction of local max pooling layers has allowed DCNNs to enjoy some translation invariance, however the intermediate feature maps are not actually invariant to large geometric transformations of the input data [152]. Therefore, many approaches have been presented to enhance robustness, aiming at learning invariant CNN representations with respect to different types of transformations such as scale [131, 21], rotation [21, 42, 284, 323], or both [126].

One representative work is Spatial Transformer Network (STN) [126], which introduces a new learnable module to handle scaling, cropping, rotations, as well as nonrigid deformations via a global parametric transformation. STN has now been used in rotated text detection [126], rotated face detection and generic object detection [280].

Although rotation invariance may be attractive in certain applications, such as scene text detection [103, 184], face detection [243], and aerial imagery [57, 288], there is limited generic object detection work focusing on rotation invariance because popular benchmark detection datasets (*e.g.* PASCAL VOC, ImageNet, COCO) do not actually present rotated images.

Before deep learning, Deformable Part based Models (DPMs) [74] were successful for generic object detection, representing objects by component parts arranged in a deformable configuration. Although DPMs have been significantly outperformed by more recent object detectors, their spirit still deeply influences many recent detectors. DPM modeling is less sensitive to transformations in object pose, viewpoint and nonrigid deformations, motivating researchers [51, 86, 188, 203, 277] to explicitly model object composition to improve CNN based detection. The first attempts [86, 277] combined DPMs with CNNs by using deep features learned by AlexNet in DPM based detection, but without region proposals. To enable a CNN to benefit from the built-in capability of modeling the deformations of object parts, a number of approaches were proposed, including DeepIDNet [203], DCN [51] and DPFCN [188] (shown in Table 7). Although similar in spirit, deformations are computed in different ways: DeepIDNet [206] designed a deformation constrained pooling layer to replace regular max pooling, to learn the shared visual patterns and their deformation properties across different object classes; DCN [51] designed a deformable convolution layer and a deformable RoI pooling layer, both of which are based on the idea of augmenting regular grid sampling locations in feature maps; and DPFCN [188] proposed a deformable part-based RoI pooling layer which selects discriminative parts of objects around object proposals by simultaneously optimizing latent displacements of all parts.

**Handling of occlusions:** In real-world images, occlusions are common, resulting in information loss from object instances. A deformable parts idea can be useful for occlusion handling, so deformable RoI Pooling [51, 188, 202] and deformable convolution [51] have been proposed to alleviate occlusion by giving more flexibility to the typically fixed geometric structures. Wang *et al.* [280] propose to learn an adversarial network that generates examples with occlusions and deformations, and context may be helpful in dealing with occlusions [309]. Despite these efforts, the occlusion problem is far from being solved; applying GANs to this problem may be a promising research direction.

**Handling of image degradations:** Image noise is a common problem in many real-world applications. It is frequently caused by insufficient lighting, low quality cameras, image compression, or the intentional low-cost sensors on edge devices and wearable devices. While low image quality may be expected to degrade the performance of visual recognition, most current methods are evaluated in a degradation free and clean environment, evidenced by the fact that PASCAL VOC, ImageNet, MS COCO and Open Images all focus on relatively high quality images. To the best of our

knowledge, there is so far very limited work to address this problem.

## 7 Context Modeling

In the physical world, visual objects occur in particular environments and usually coexist with other related objects. There is strong psychological evidence [14, 10] that context plays an essential role in human object recognition, and it is recognized that a proper modeling of context helps object detection and recognition [266, 197, 33, 32, 58, 78], especially when object appearance features are insufficient because of small object size, object occlusion, or poor image quality. Many different types of context have been discussed [58, 78], and can broadly be grouped into one of three categories:

1. Semantic context: The likelihood of an object to be found in some scenes, but not in others;
2. Spatial context: The likelihood of finding an object in some position and not others with respect to other objects in the scene;
3. Scale context: Objects have a limited set of sizes relative to other objects in the scene.

A great deal of work [34, 58, 78, 185, 193, 220, 207] preceded the prevalence of deep learning, and much of this work has yet to be explored in DCNN-based object detectors [35, 114].

The current state of the art in object detection [229, 175, 102] detects objects without explicitly exploiting any contextual information. It is broadly agreed that DCNNs make use of contextual information implicitly [303, 316] since they learn hierarchical representations with multiple levels of abstraction. Nevertheless, there is value in exploring contextual information explicitly in DCNN based detectors [114, 35, 305], so the following reviews recent work in exploiting contextual cues in DCNN- based object detectors, organized into categories of *global* and *local* contexts, motivated by earlier work in [310, 78]. Representative approaches are summarized in Table 8.

### 7.1 Global Context

Global context [310, 78] refers to image or scene level contexts, which can serve as cues for object detection (*e.g.,* a bedroom will predict the presence of a bed). In DeepIDNet [203], the image classification scores were used as contextual features, and concatenated with the object detection scores to improve detection results. In ION [11], Bell *et al.* proposed to use spatial Recurrent Neural Networks (RNNs) to explore contextual information across the entire image. In SegDeepM [326], Zhu *et al.* proposed a Markov random field model that scores appearance as well as context for each detection, and allows each candidate box to select a segment out of a large pool of object segmentation proposals and score the agreement between them. In [245], semantic segmentation was used as a form of contextual priming.

### 7.2 Local Context

Local context [310, 78, 220] considers the relationship among locally nearby objects, as well as the interactions between an ob-

**Table 8** Summary of detectors that exploit context information, with labelling details as in Table 7.

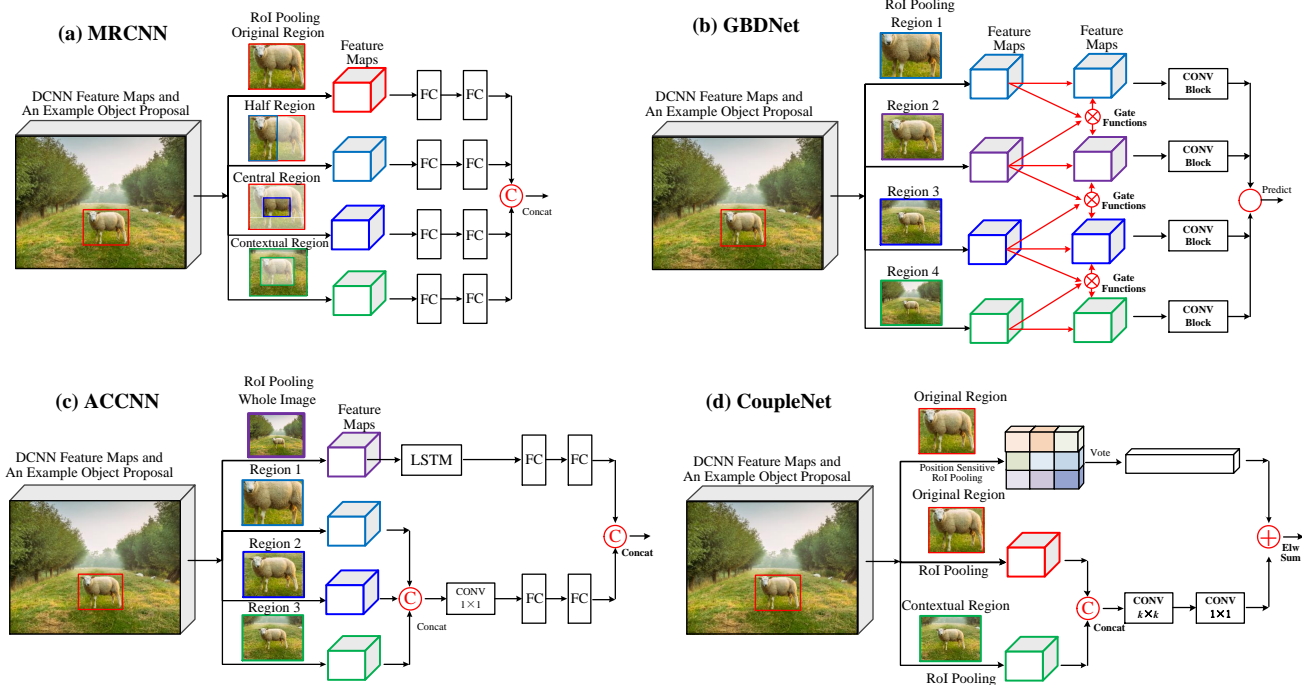| Group | Detector Name | Region Proposal | Backbone DCNN | Pipelined Used | mAP@IoU=0.5 VOC07 | mAP@IoU=0.5 VOC12 | mAP COCO | Published In | Highlights |
|---|---|---|---|---|---|---|---|---|---|
| **Global Context** | SegDeepM [326] | SS+CMPC | VGG16 | RCNN | VOC10 | VOC12 | – | CVPR15 | Additional features extracted from an enlarged object proposal as context information. |
| | DeepIDNet [203] | SS+EB | AlexNet ZFNet | RCNN | 69.0 (07) | – | – | CVPR15 | Use image classification scores as global contextual information to refine the detection scores of each object proposal. |
| | ION [11] | SS+EB | VGG16 | Fast RCNN | 80.1 | 77.9 | 33.1 | CVPR16 | The contextual information outside the region of interest is integrated using spatial recurrent neural networks. |
| | CPF [245] | RPN | VGG16 | Faster RCNN | 76.4 (07+12) | 72.6 (07T+12) | – | ECCV16 | Use semantic segmentation to provide top-down feedback. |
| **Local Context** | MRCNN [82] | SS | VGG16 | SPPNet | 78.2 (07+12) | 73.9 (07+12) | – | ICCV15 | Extract features from multiple regions surrounding or inside the object proposals. Integrate the semantic segmentation-aware features. |
| | GBDNet [304, 305] | CRAFT [292] | Inception v2 ResNet269 PolyNet [311] | Fast RCNN | 77.2 (07+12) | – | 27.0 | ECCV16 TPAMI18 | A GBDNet module to learn the relations of multiscale contextualized regions surrounding an object proposal; GBDNet passes messages among features from different context regions through convolution between neighboring support regions in two directions. |
| | ACCNN[157] | SS | VGG16 | Fast RCNN | 72.0 (07+12) | 70.6 (07T+12) | – | TMM17 | Use LSTM to capture global context. Concatenate features from multi-scale contextual regions surrounding an object proposal. The global and local context features are concatenated for recognition. |
| | CoupleNet[327] | RPN | ResNet101 | RFCN | **82.7** (07+12) | **80.4** (07T+12) | 34.4 | ICCV17 | Concatenate features from multiscale contextual regions surrounding an object proposal. Features of different contextual regions are then combined by convolution and element-wise sum. |
| | SMN [35] | RPN | VGG16 | Faster RCNN | 70.0 (07) | – | – | ICCV17 | Model object-object relationships efficiently through a spatial memory network. Learn the functionality of NMS automatically. |
| | ORN [114] | RPN | ResNet101 +DCN | Faster RCNN | – | – | **39.0** | CVPR18 | Model the relations of a set of object proposals through the interactions between their appearance features and geometry. Learn the functionality of NMS automatically. |
| | SIN [176] | RPN | VGG16 | Faster RCNN | 76.0 (07+12) | 73.1 (07T+12) | 23.2 | CVPR18 | Formulate object detection as graph-structured inference, where objects are graph nodes and relationships the edges. |



**Fig. 18** Representative approaches that explore local surrounding contextual features: MRCNN [82], GBDNet [304, 305], ACCNN [157] and CoupleNet [327]; also see Table 8.

ject and its surrounding area. In general, modeling object relations is challenging, requiring reasoning about bounding boxes of different classes, locations, scales *etc*. Deep learning research that explicitly models object relations is quite limited, with representative ones being Spatial Memory Network (SMN) [35], Object Relation Network [114], and Structure Inference Network (SIN) [176]. In SMN, spatial memory essentially assembles object in-

stances back into a pseudo image representation that is easy to be fed into another CNN for object relations reasoning, leading to a new sequential reasoning architecture where image and memory are processed in parallel to obtain detections which further update memory. Inspired by the recent success of attention modules in natural language processing [274], ORN processes a set of objects simultaneously through the interaction between their appearance

feature and geometry. It does not require additional supervision, and it is easy to embed into existing networks, effective in improving object recognition and duplicate removal steps in modern object detection pipelines, giving rise to the first fully end-to-end object detector. SIN [176] considered two kinds of context: scene contextual information and object relationships within a single image. It formulates object detection as a problem of graph inference, where the objects are treated as nodes in a graph and relationships between objects are modeled as edges.

A wider range of methods has approached the context challenge with a simpler idea: enlarging the detection window size to extract some form of local context. Representative approaches include MRCNN [82], Gated BiDirectional CNN (GBDNet) [304, 305], Attention to Context CNN (ACCNN) [157], CoupleNet [327], and Sermanet *et al.* [238]. In MRCNN [82] (Fig. 18 (a)), in addition to the features extracted from the original object proposal at the last CONV layer of the backbone, Gidaris and Komodakis proposed to extract features from a number of different regions of an object proposal (half regions, border regions, central regions, contextual region and semantically segmented regions), in order to obtain a richer and more robust object representation. All of these features are combined by concatenation.

Quite a number of methods, all closely related to MRCNN, have been proposed since then. The method in [302] used only four contextual regions, organized in a foveal structure, where the classifiers along multiple paths are trained jointly end-to-end. Zeng *et al.* proposed GBDNet [304, 305] (Fig. 18 (b)) to extract features from multiscale contextualized regions surrounding an object proposal to improve detection performance. In contrast to the somewhat naive approach of learning CNN features for each region separately and then concatenating them, GBDNet passes messages among features from different contextual regions. Noting that message passing is not always helpful, but dependent on individual samples, Zeng *et al.* [304] used gated functions to control message transmission. Li *et al.* [157] presented ACCNN (Fig. 18 (c)) to utilize both global and local contextual information: the global context was captured using a Multiscale Local Contextualized (MLC) subnetwork, which recurrently generates an attention map for an input image to highlight promising contextual locations; local context adopted a method similar to that of MRCNN [82]. As shown in Fig. 18 (d), CoupleNet [327] is conceptually similar to ACCNN [157], but built upon RFCN [50], which captures object information with position sensitive RoI pooling, CoupleNet added a branch to encode the global context with RoI pooling.

## 8 Detection Proposal Methods

An object can be located at any position and scale in an image. During the heyday of handcrafted feature descriptors (SIFT [179], HOG [52] and LBP [196]), the most successful methods for object detection (*e.g.* DPM [72]) used *sliding window* techniques [276, 52, 72, 98, 275]. However, the number of windows is huge, growing with the number of pixels in an image, and the need to search at multiple scales and aspect ratios further increases the

search space[12]. Therefore, it is computationally too expensive to apply sophisticated classifiers.

Around 2011, researchers proposed to relieve the tension between computational tractability and high detection quality by using *detection proposals*[13] [273, 271]. Originating in the idea of *objectness* proposed by [2], object proposals are a set of candidate regions in an image that are likely to contain objects, and if high object recall can be achieved with a modest number of object proposals (like one hundred), significant speed-ups over the sliding window approach can be gained, allowing the use of more sophisticated classifiers. Detection proposals are usually used as a pre-processing step, limiting the number of regions that need to be evaluated by the detector, and should have the following characteristics:

1. High recall, which can be achieved with only a few proposals;
2. Accurate localization, such that the proposals match the object bounding boxes as accurately as possible; and
3. Low computational cost.

The success of object detection based on detection proposals [273, 271] has attracted broad interest [25, 7, 3, 43, 330, 65, 138, 186]. A comprehensive review of object proposal algorithms is beyond the scope of this paper, because object proposals have applications beyond object detection [6, 93, 328]. We refer interested readers to the recent surveys [110, 27] which provide in-depth analysis of many classical object proposal algorithms and their impact on detection performance. Our interest here is to review object proposal methods that are based on DCNNs, output class agnostic proposals, and are related to generic object detection.

In 2014, the integration of object proposals [273, 271] and DCNN features [140] led to the milestone RCNN [85] in generic object detection. Since then, detection proposal has quickly become a standard preprocessing step, based on the fact that all winning entries in the PASCAL VOC [68], ILSVRC [234] and MS COCO [166] object detection challenges since 2014 used detection proposals [85, 203, 84, 229, 305, 102].

Among object proposal approaches based on traditional low-level cues (*e.g.,* color, texture, edge and gradients), Selective Search [271], MCG [7] and EdgeBoxes [330] are among the more popular. As the domain rapidly progressed, traditional object proposal approaches [271, 110, 330], which were adopted as external modules independent of the detectors, became the speed bottleneck of the detection pipeline [229]. An emerging class of object proposal algorithms [67, 229, 142, 81, 213, 292] using DCNNs has attracted broad attention.

Recent DCNN based object proposal methods generally fall into two categories: *bounding box* based and *object segment* based, with representative methods summarized in Table 9.
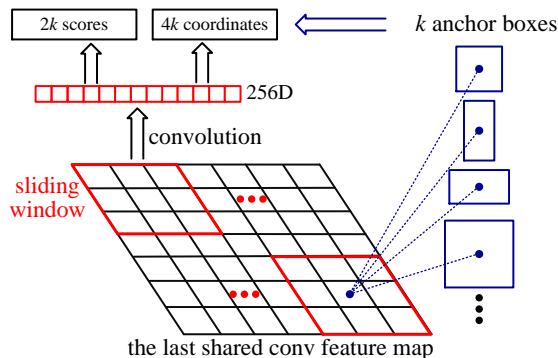
**Bounding Box Proposal Methods** are best exemplified by the RPC method [229] of Ren *et al.*, illustrated in Fig. 19. RPN predicts object proposals by sliding a small network over the feature map of the last shared CONV layer. At each sliding window location, $k$ proposals are predicted by using $k$ anchor boxes, where

---

[12] Sliding window based detection requires classifying around $10^4$-$10^5$ windows per image. The number of windows grows significantly to $10^6$-$10^7$ windows per image when considering multiple scales and aspect ratios.

[13] We use the terminology *detection proposals*, *object proposals* and *region proposals* interchangeably.

**Table 9** Summary of object proposal methods using DCNN. Blue indicates the number of object proposals. The detection results on COCO are based on mAP@IoU[0.5, 0.95], unless stated otherwise.

| Proposer Name | Backbone Network | Detector Tested | Recall@IoU (VOC07) | | | Detection Results (mAP) | | | Published In | Highlights |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 0.5 | 0.7 | 0.9 | VOC07 | VOC12 | COCO | | |
| MultiBox1[67] | AlexNet | RCNN | – | – | – | 29.0 (10) (12) | | | CVPR14 | Learns a class agnostic regressor on a small set of 800 predefined anchor boxes. Do not share features for detection. |
| DeepBox [142] | VGG16 | Fast RCNN | 0.96 (1000) | 0.84 (1000) | 0.15 (1000) | – | – | 37.8 (500) (IoU@0.5) | ICCV15 | Use a lightweight CNN to learn to rerank proposals generated by EdgeBox. Can run at 0.26s per image. Do not share features for detection. |
| RPN[229, 230] | VGG16 | Faster RCNN | 0.97 (300) 0.98 (1000) | 0.79 (300) 0.84 (1000) | 0.04 (300) 0.04 (1000) | 73.2 (300) (07+12) | 70.4 (300) (07++12) | 21.9 (300) | NIPS15 | The first to generate object proposals by sharing full image convolutional features with detection. Most widely used object proposal method. Significant improvements in detection speed. |
| DeepProposal[81] | VGG16 | Fast RCNN | 0.74 (100) 0.92 (1000) | 0.58 (100) 0.80 (1000) | 0.12 (100) 0.16 (1000) | 53.2 (100) (07) | – | – | ICCV15 | Generate proposals inside a DCNN in a multiscale manner. Share features with the detection network. |
| CRAFT [292] | VGG16 | Faster RCNN | 0.98 (300) | 0.90 (300) | 0.13 (300) | 75.7 (07+12) | 71.3 (12) | – | CVPR16 | Introduced a classification network (*i.e.* two class Fast RCNN) cascade that comes after the RPN. Not sharing features extracted for detection. |
| AZNet [181] | VGG16 | Fast RCNN | 0.91 (300) | 0.71 (300) | 0.11 (300) | 70.4 (07) | – | 22.3 | CVPR16 | Use coarse-to-fine search: start from large regions, then recursively search for subregions that may contain objects. Adaptively guide computational resources to focus on likely subregions. |
| ZIP [156] | Inception v2 | Faster RCNN | 0.85 (300) COCO | 0.74 (300) COCO | 0.35 (300) COCO | 79.8 (07+12) | – | – | IJCV18 | Generate proposals using conv-deconv network with multilayers; Proposed a map attention decision (MAD) unit to assign the weights for features from different layers. |
| DeNet[269] | ResNet101 | Fast RCNN | 0.82 (300) | 0.74 (300) | 0.48 (300) | 77.1 (07+12) | 73.9 (07++12) | 33.8 | ICCV17 | A lot faster than Faster RCNN; Introduces a bounding box corner estimation for predicting object proposals efficiently to replace RPN; Does not require predefined anchors. |
| Proposer Name | Backbone Network | Detector Tested | Box Proposals (AR, COCO) | | | Segment Proposals (AR, COCO) | | | Published In | Highlights |
| DeepMask [213] | VGG16 | Fast RCNN | 0.33 (100), 0.48 (1000) | | | 0.26 (100), 0.37 (1000) | | | NIPS15 | First to generate object mask proposals with DCNN; Slow inference time; Need segmentation annotations for training; Not sharing features with detection network; Achieved mAP of 69.9% (500) with Fast RCNN. |
| InstanceFCN [48] | VGG16 | – | – | | | 0.32 (100), 0.39 (1000) | | | ECCV16 | Combines ideas of FCN [177] and DeepMask [213]. Introduces instance sensitive score maps. Needs segmentation annotations to train the network. |
| SharpMask [214] | MPN [302] | Fast RCNN | 0.39 (100), 0.53 (1000) | | | 0.30 (100), 0.39 (1000) | | | ECCV16 | Leverages features at multiple convolutional layers by introducing a top-down refinement module. Does not share features with detection network. Needs segmentation annotations for training. |
| FastMask[113] | ResNet39 | – | 0.43 (100), 0.57 (1000) | | | 0.32 (100), 0.41 (1000) | | | CVPR17 | Generates instance segment proposals efficiently in one-shot manner similar to SSD [175]. Uses multiscale convolutional features. Uses segmentation annotations for training. |

*The table rows are grouped under two side labels:* **Bounding Box Object Proposal Methods** *(upper section)* and **Segment Proposal Methods** *(lower section).*



**Fig. 19** Illustration of the Region Proposal Network (RPN) introduced in [229].

each anchor box[14] is centered at some location in the image, and is associated with a particular scale and aspect ratio. Ren *et al.* [229] proposed integrating RPN and Fast RCNN into a single network by sharing their convolutional layers, leading to Faster RCNN, the first end-to-end detection pipeline. RPN has been broadly selected as the proposal method by many state-of-the-art object detectors, as can be observed from Tables 7 and 8.

Instead of fixing *a priori* a set of anchors as MultiBox [67, 262] and RPN [229], Lu *et al.* [181] proposed generating anchor locations by using a recursive search strategy which can adaptively guide computational resources to focus on sub-regions likely to

contain objects. Starting with the whole image, all regions visited during the search process serve as anchors. For any anchor region encountered during the search procedure, a scalar zoom indicator is used to decide whether to further partition the region, and a set of bounding boxes with objectness scores are computed by an Adjacency and Zoom Network (AZNet), which extends RPN by adding a branch to compute the scalar zoom indicator in parallel with the existing branch.

Further work attempts to generate object proposals by exploiting multilayer convolutional features. Concurrent with RPN [229], Ghodrati *et al.* [81] proposed DeepProposal, which generates object proposals by using a cascade of multiple convolutional features, building an inverse cascade to select the most promising object locations and to refine their boxes in a coarse-to-fine manner. An improved variant of RPN, HyperNet [135] designs Hyper Features which aggregate multilayer convolutional features and shares them both in generating proposals and detecting objects via an end-to-end joint training strategy. Yang *et al.* proposed CRAFT [292] which also used a cascade strategy, first training an RPN network to generate object proposals and then using them to train another binary Fast RCNN network to further distinguish objects from background. Li *et al.* [156] proposed ZIP to improve RPN by predicting object proposals with multiple convolutional feature maps at different network depths to integrate both low level details and high level semantics. The backbone used in ZIP is a "zoom out and in" network inspired by the conv and deconv structure [177].

---

[14] The concept of "anchor" first appeared in [229].

Finally, recent work which deserves mention includes Deep-box [142], which proposed a lightweight CNN to learn to rerank proposals generated by EdgeBox, and DeNet [269] which introduces bounding box corner estimation to predict object proposals efficiently to replace RPN in a Faster RCNN style detector.

**Object Segment Proposal Methods** [213, 214] aim to generate segment proposals that are likely to correspond to objects. Segment proposals are more informative than bounding box proposals, and take a step further towards object instance segmentation [96, 49, 162]. In addition, using instance segmentation supervision can improve the performance of bounding box object detection. The pioneering work of DeepMask, proposed by Pinheiro *et al.* [213], segments proposals learnt directly from raw image data with a deep network. Similarly to RPN, after a number of shared convolutional layers DeepMask splits the network into two branches in order to predict a class agnostic mask and an associated objectness score. Also similar to the efficient sliding window strategy in Over-Feat [239], the trained DeepMask network is applied in a sliding window manner to an image (and its rescaled versions) during inference. More recently, Pinheiro *et al.* [214] proposed SharpMask by augmenting the DeepMask architecture with a refinement module, similar to the architectures shown in Fig. 17 (b1) and (b2), augmenting the feed-forward network with a top-down refinement process. SharpMask can efficiently integrate spatially rich information from early features with strong semantic information encoded in later layers to generate high fidelity object masks.

Motivated by Fully Convolutional Networks (FCN) for semantic segmentation [177] and DeepMask [213], Dai *et al.* proposed InstanceFCN [48] to generate instance segment proposals. Similar to DeepMask, the InstanceFCN network is split into two fully convolutional branches, one to generate instance sensitive score maps, the other to predict the objectness score. Hu *et al.* proposed Fast-Mask [113] to efficiently generate instance segment proposals in a one-shot manner, similar to SSD [175], in order to make use of multiscale convolutional features. Sliding windows extracted densely from multiscale convolutional feature maps were input to a scale-tolerant attentional head module in order to predict segmentation masks and objectness scores. FastMask is claimed to run at 13 FPS on $800 \times 600$ images.

# 9 Other Issues

**Data Augmentation.** Performing data augmentation for learning DCNNs [26, 84, 85] is generally recognized to be important for visual recognition. Trivial data augmentation refers to perturbing an image by transformations that leave the underlying category unchanged, such as cropping, flipping, rotating, scaling, translating, color perturbations, and adding noise. By artificially enlarging the number of samples, data augmentation helps in reducing overfitting and improving generalization. It can be used at training time, at test time, or both. Nevertheless, it has the obvious limitation that the time required for training increases significantly. Data augmentation may synthesize completely new training images [210, 280], however it is hard to guarantee that the synthetic images generalize well to real ones. Some researchers [64, 94] proposed augmenting datasets by pasting real segmented objects into natural images; indeed, Dvornik *et al.* [63] showed that appropriately modeling the
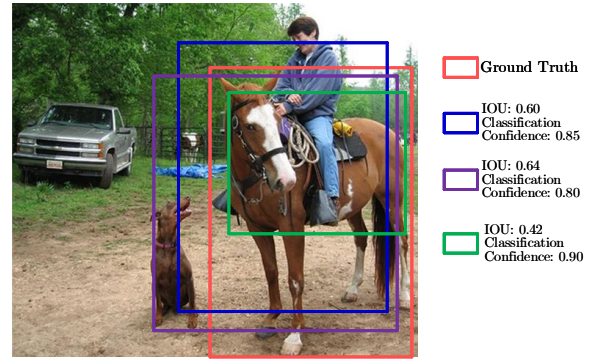


**Fig. 20** Localization error could stem from insufficient overlap or duplicate detections. Localization error is a frequent cause of false positives.

visual context surrounding objects is crucial to place them in the right environment, and proposed a context model to automatically find appropriate locations on images to place new objects for data augmentation.

**Novel Training Strategies.** Detecting objects under a wide range of scale variations, especially the detection of very small objects, stands out as a key challenge. It has been shown [120, 175] that image resolution has a considerable impact on detection accuracy, therefore scaling is particularly commonly used in data augmentation, since higher resolutions increase the possibility of detecting small objects [120]. Recently, Singh *et al.* proposed advanced and efficient data argumentation methods SNIP [249] and SNIPER [251] to illustrate the scale invariance problem, as summarized in Table 10. Motivated by the intuitive understanding that small and large objects are difficult to detect at smaller and larger scales, respectively, SNIP introduces a novel training scheme that can reduce scale variations during training, but without reducing training samples; SNIPER allows for efficient multiscale training, only processing context regions around ground truth objects at the appropriate scale, instead of processing a whole image pyramid. Peng *et al.* [209] studied a key factor in training, the minibatch size, and proposed MegDet, a Large MiniBatch Object Detector, to enable the training with a much larger minibatch size than before (from 16 to 256). To avoid the failure of convergence and significantly speed up the training process, Peng *et al.* [209] proposed a learning rate policy and Cross GPU Batch Normalization, and effectively utilized 128 GPUs, allowing MegDet to finish COCO training in 4 hours on 128 GPUs, and winning the COCO 2017 Detection Challenge.

**Reducing Localization Error.** In object detection, the Intersection Over Union[15] (IOU) between a detected bounding box and its ground truth box is the most popular evaluation metric, and an IOU threshold (*e.g.* typical value of 0.5) is required to define positives and negatives. From Fig. 13, in most state of the art detectors [84, 175, 102, 229, 227] object detection is formulated as a multi-task learning problem, *i.e.,* jointly optimizing a softmax classifier which assigns object proposals with class labels and bounding box regressors, localizing objects by maximizing IOU or other metrics between detection results and ground truth. Bounding boxes are only a crude approximation for articulated objects, consequently background pixels are almost invariably included in a bounding box, which affects the accuracy of classification and localization.

---

[15] Please refer to Section 4.2 for more details on the definition of IOU.

**Table 10** Representative methods for training strategies and class imbalance handling. Results on COCO are reported with Test Dev. The detection results on COCO are based on mAP@IoU[0.5, 0.95].

| Detector Name | Region Proposal | Backbone DCNN | Pipelined Used | VOC07 Results | VOC12 Results | COCO Results | Published In | Highlights |
|---|---|---|---|---|---|---|---|---|
| MegDet [209] | RPN | ResNet50 +FPN | Faster RCNN | – | – | 52.5 | CVPR18 | Allow training with much larger minibatch size than before by introducing cross GPU batch normalization; Can finish the COCO training in 4 hours on 128 GPUs and achieved improved accuracy; Won COCO2017 detection challenge. |
| SNIP [251] | RPN | DPN [37] +DCN [51] | RFCN | – | – | 48.3 | CVPR18 | A new multiscale training scheme. Empirically examined the effect of up-sampling for small object detection. During training, only select objects that fit the scale of features as positive samples. |
| SNIPER [251] | RPN | ResNet101 +DCN | Faster RCNN | – | – | 47.6 | 2018 | An efficient multiscale training strategy. Process context regions around ground-truth instances at the appropriate scale. |
| OHEM [246] | SS | VGG16 | Fast RCNN | 78.9 (07+12) | 76.3 (07++12) | 22.4 | CVPR16 | A simple and effective Online Hard Example Mining algorithm to improve training of region based detectors. |
| FactorNet [204] | SS | GooglNet | RCNN | – | – | – | CVPR16 | Identify the imbalance in the number of samples for different object categories; propose a divide-and-conquer feature learning scheme. |
| Chained Cascade [23] | SS CRAFT | VGG Inceptionv2 | Fast RCNN, Faster RCNN | 80.4 (07+12) (SS+VGG) | – | – | ICCV17 | Jointly learn DCNN and multiple stages of cascaded classifiers. Boost detection accuracy on PASCAL VOC 2007 and ImageNet for both fast RCNN and Faster RCNN using different region proposal methods. |
| Cascade RCNN [23] | RPN | VGG ResNet101 +FPN | Faster RCNN | – | – | 42.8 | CVPR18 | Jointly learn DCNN and multiple stages of cascaded classifiers, which are learned using different localization accuracy for selecting positive samples. Stack bounding box regression at multiple stages. |
| RetinaNet [168] | – | ResNet101 +FPN | RetinaNet | – | – | 39.1 | ICCV17 | Propose a novel Focal Loss which focuses training on hard examples. Handles well the problem of imbalance of positive and negative samples when training a one-stage detector. |

The study in [108] shows that object localization error is one of the most influential forms of error, in addition to confusion between similar objects. Localization error could stem from insufficient overlap (smaller than the required IOU threshold, such as the green box in Fig. 20) or duplicate detections (*i.e.,* multiple overlapping detections for an object instance). Usually, some postprocessing step like NonMaximum Suppression (NMS) [18, 111] is used for eliminating duplicate detections. However, due to misalignments the bounding box with better localization could be suppressed during NMS, leading to poorer localization quality (such as the purple box shown in Fig. 20). Therefore, there are quite a few methods aiming at improving detection performance by reducing localization error.

MRCNN [82] introduces iterative bounding box regression, where an RCNN is applied several times. CRAFT [292] and AttractioNet [83] use a multi-stage detection sub-network to generate accurate proposals, to forward to Fast RCNN. Cai and Vasconcelos proposed Cascade RCNN [23], a multistage extension of RCNN, in which a sequence of detectors is trained sequentially with increasing IOU thresholds, based on the observation that the output of a detector trained with a certain IOU is a good distribution to train the detector of the next higher IOU threshold, in order to be sequentially more selective against close false positives. This approach can be built with any RCNN-based detector, and is demonstrated to achieve consistent gains (about 2 to 4 points) independent of the baseline detector strength, at a marginal increase in computation. There is also recent work [128, 232, 121] formulating IOU directly as the optimization objective, and in proposing improved NMS results [18, 104, 111, 270], such as Soft NMS [18] and learning NMS [111].

**Class Imbalance Handling.** Unlike image classification, object detection has another unique problem: the serious imbalance between the number of labeled object instances and the number of background examples (image regions not belonging to any object class of interest). Most background examples are easy negatives, however this imbalance can make the training very inefficient, and the large number of easy negatives tends to overwhelm the training. In the past, this issue has typically been addressed via techniques such as bootstrapping [259]. More recently, this problem has also seen some attention [153, 168, 246]. Because the region proposal stage rapidly filters out most background regions and proposes a small number of object candidates, this class imbalance issue is mitigated to some extent in two-stage detectors [85, 84, 229, 102], although example mining approaches, such as Online Hard Example Mining (OHEM) [246], may be used to maintain a reasonable balance between foreground and background. In the case of one-stage object detectors [227, 175], this imbalance is extremely serious (*e.g.* 100,000 background examples to every object). Lin *et al.* [168] proposed Focal Loss to address this by rectifying the Cross Entropy loss, such that it down-weights the loss assigned to correctly classified examples. Li *et al.* [153] studied this issue from the perspective of gradient norm distribution, and proposed a Gradient Harmonizing Mechanism (GHM) to handle it.

## 10 Discussion and Conclusion

Generic object detection is an important and challenging problem in computer vision and has received considerable attention. Thanks to remarkable developments in deep learning techniques, the field of object detection has dramatically evolved. As a comprehensive survey on deep learning for generic object detection, this paper has highlighted the recent achievements, provided a structural taxonomy for methods according to their roles in detection, summarized existing popular datasets and evaluation criteria, and discussed performance for the most representative methods. We conclude this

review with a discussion of the state of the art in Section 10.1, an overall discussion of key issues in Section 10.2, and finally suggested future research directions in Section 10.3.

## 10.1 State of the Art Performance

A large variety of detectors has appeared in the last few years, and the introduction of standard benchmarks, such as PASCAL VOC [68, 69], ImageNet [234] and COCO [166], has made it easier to compare detectors. As can be seen from our earlier discussion in Sections 5 through 9, it may be misleading to compare detectors in terms of their originally reported performance (*e.g.* accuracy, speed), as they can differ in fundamental / contextual respects, including the following choices:

- Meta detection frameworks, such as RCNN [85], Fast RCNN [84], Faster RCNN [229], RFCN [50], Mask RCNN [102], YOLO [227] and SSD [175];
- Backbone networks such as VGG [248], Inception [263, 125, 264], ResNet [101], ResNeXt [291], and Xception [45] *etc.* listed in Table 6;
- Innovations such as multilayer feature combination [167, 247, 77], deformable convolutional networks [51], deformable RoI pooling [203, 51], heavier heads [231, 209], and lighter heads [165];
- Pretraining with datasets such as ImageNet [234], COCO [166], Places [319], JFT [106] and Open Images [139];
- Different detection proposal methods and different numbers of object proposals;
- Train/test data augmentation, novel multiscale training strategies [249, 251] *etc*, and model ensembling.

Although it may be impractical to compare every recently proposed detector, it is nevertheless valuable to integrate representative and publicly available detectors into a common platform and to compare them in a unified manner. There has been very limited work in this regard, except for Huang's study [120] of the three main families of detectors (Faster RCNN [229], RFCN [50] and SSD [175]) by varying the backbone network, image resolution, and the number of box proposals.

As can be seen from Tables 7, 8, 9, 10, 11, we have summarized the best reported performance of many methods on three widely used standard benchmarks. The results of these methods were reported on the same test benchmark, despite their differing in one or more of the aspects listed above.

Figs. 3 and 21 present a very brief overview of the state of the art, summarizing the best detection results of the PASCAL VOC, ILSVRC and MSCOCO challenges; more results can be found at detection challenge websites [124, 189, 208]. The competition winner of the open image challenge object detection task achieved 61.71% mAP in the public leader board and 58.66% mAP on the private leader board, obtained by combining the detection results of several two-stage detectors including Fast RCNN [84], Faster RCNN [229], FPN [167], Deformable RCNN [51], and Cascade RCNN [23]. In summary, the backbone network, the detection framework, and the availability of large scale datasets are the three most important factors in detection accuracy. Ensembles of multiple models, the incorporation of context features, and data augmentation all help to achieve better accuracy.
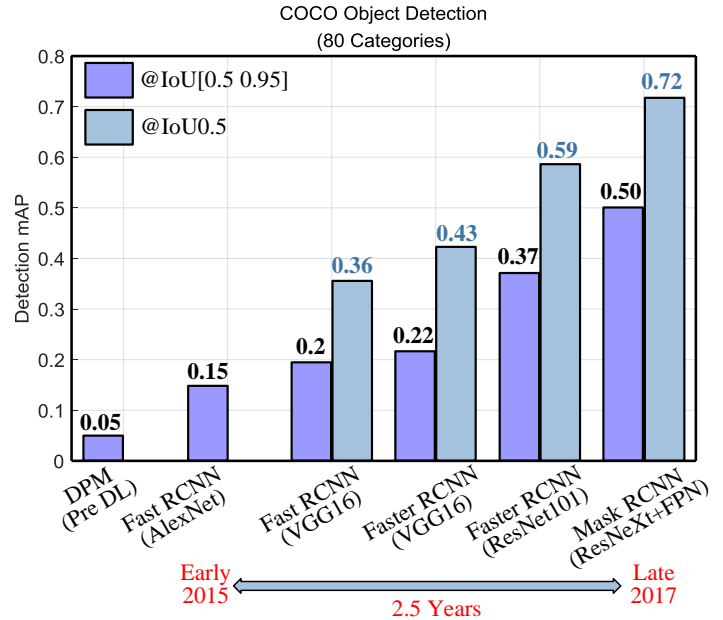


**Fig. 21** Evolution of object detection performance on COCO (Test-Dev results). Results are quoted from [84, 102, 230]. The backbone network, the design of detection framework and the availability of good and large scale datasets are the three most important factors in detection accuracy.

In less than five years, since AlexNet [140] was proposed, the Top5 error on ImageNet classification [234] with 1000 classes has dropped from 16% to 2%, as shown in Fig. 15. However, the mAP of the best performing detector [209] on COCO [166], trained to detect only 80 classes, is only at 73%, even at 0.5 IoU, illustrating how object detection is much harder than image classification. The accuracy and robustness achieved by the state-of-the-art detectors far from satisfies the requirements of real world applications, so there remains significant room for future improvement.

## 10.2 Summary and Discussion

With hundreds of references and many dozens of methods discussed throughout this paper, we would now like to focus on the key factors which have emerged in generic object detection based on deep learning.

**(1) Detection Frameworks: Two Stage vs. One Stage**
In Section 5 we identified two major categories of detection frameworks: region based (two stage) and unified (one stage):

- When large computational cost is allowed, two-stage detectors generally produce higher detection accuracies than one-stage, evidenced by the fact that most winning approaches used in famous detection challenges like are predominantly based on two-stage frameworks, because their structure is more flexible and better suited for region based classification. The most widely used frameworks are Faster RCNN [229], RFCN [50] and Mask RCNN [102].
- It has been shown in [120] that the detection accuracy of one-stage SSD [175] is less sensitive to the quality of the backbone network than representative two-stage frameworks.
- One-stage detectors like YOLO [227] and SSD [175] are generally faster than two-stage ones, because of avoiding preprocessing algorithms, using lightweight backbone networks, per-

forming prediction with fewer candidate regions, and making the classification subnetwork fully convolutional. However, two-stage detectors can run in real time with the introduction of similar techniques. In any event, whether one stage or two, the most time consuming step is the feature extractor (backbone network) [146, 229].

- It has been shown [120, 227, 175] that one-stage frameworks like YOLO and SSD typically have much poorer performance when detecting small objects than two-stage architectures like Faster RCNN and RFCN, but are competitive in detecting large objects.

There have been many attempts to build better (faster, more accurate, or more robust) detectors by attacking each stage of the detection framework. No matter whether one, two or multiple stages, the design of the detection framework has converged towards a number of crucial design choices:

- A fully convolutional pipeline
- Exploring complementary information from other correlated tasks, *e.g.*, Mask RCNN [102]
- Sliding windows [229]
- Fusing information from different layers of the backbone.

The evidence from recent success of cascade for object detection [23, 40, 41] and instance segmentation on COCO [31] and other challenges has shown that multistage object detection could be a future framework for a speed-accuracy trade-off. A teaser investigation is being done in the 2019 WIDER Challenge [180].

### (2) Backbone Networks

As discussed in Section 6.1, backbone networks are one of the main driving forces behind the rapid improvement of detection performance, because of the key role played by discriminative object feature representation. Generally, deeper backbones such as ResNet [101], ResNeXt [291], InceptionResNet [265] perform better; however, they are computationally more expensive and require much more data and massive computing for training. Some backbones [112, 123, 312] were proposed for focusing on speed instead, such as MobileNet [112] which has been shown to achieve VGGNet16 accuracy on ImageNet with only $\frac{1}{30}$ the computational cost and model size. Backbone training from scratch may become possible as more training data and better training strategies are available [285, 183, 182].

### (3) Improving the Robustness of Object Representation

The variation of real world images is a key challenge in object recognition. The variations include lighting, pose, deformations, background clutter, occlusions, blur, resolution, noise, and camera distortions.

### (3.1) Object Scale and Small Object Size

Large variations of object scale, particularly those of small objects, pose a great challenge. Here a summary and discussion on the main strategies identified in Section 6.2:

- Using image pyramids: They are simple and effective, helping to enlarge small objects and to shrink large ones. They are computationally expensive, but are nevertheless commonly used during inference for better accuracy.
- Using features from convolutional layers of different resolutions: In early work like SSD [175], predictions are performed

independently, and no information from other layers is combined or merged. Now it is quite standard to combine features from different layers, e.g. in FPN [167].

- Using dilated convolutions [164, 163]: A simple and effective method to incorporate broader context and maintain high resolution feature maps.
- Using anchor boxes of different scales and aspect ratios: Drawbacks of having many parameters, and scales and aspect ratios of anchor boxes are usually heuristically determined.
- Up-scaling: Particularly for the detection of small objects, high-resolution networks [255, 256] can be developed. It remains unclear whether super-resolution techniques improve detection accuracy or not.

Despite recent advances, the detection accuracy for small objects is still much lower than that of larger ones. Therefore, the detection of small objects remains one of the key challenges in object detection. Perhaps localization requirements need to be generalized as a function of scale, since certain applications, e.g. autonomous driving, only require the identification of the existence of small objects within a larger region, and exact localization is not necessary.

### (3.2) Deformation, Occlusion, and other factors

As discussed in Section 2.2, there are approaches to handling geometric transformation, occlusions, and deformation mainly based on two paradigms. The first is a spatial transformer network, which uses regression to obtain a deformation field and then warp features according to the deformation field [51]. The second is based on a deformable part-based model [74], which finds the maximum response to a part filter with spatial constraints taken into consideration [203, 86, 277].

Rotation invariance may be attractive in certain applications, but there are limited generic object detection work focusing on rotation invariance, because popular benchmark detection datasets (PASCAL VOC, ImageNet, COCO) do not have large variations in rotation. Occlusion handling is intensively studied in face detection and pedestrian detection, but very little work has been devoted to occlusion handling for generic object detection. In general, despite recent advances, deep networks are still limited by the lack of robustness to a number of variations, which significantly constrains their real-world applications.

### (4) Context Reasoning

As introduced in Section 7, objects in the wild typically coexist with other objects and environments. It has been recognized that contextual information (object relations, global scene statistics) helps object detection and recognition [197], especially for small objects, occluded objects, and with poor image quality. There was extensive work preceding deep learning [185, 193, 220, 58, 78], and also quite a few works in the era of deep learning [82, 304, 305, 35, 114]. How to efficiently and effectively incorporate contextual information remains to be explored, possibly guided by how human vision uses context, based on scene graphs [161], or via the full segmentation of objects and scenes using panoptic segmentation [134].

### (5) Detection Proposals

Detection proposals significantly reduce search spaces. As recommended in [110], future detection proposals will surely have to improve in repeatability, recall, localization accuracy, and speed. Since the success of RPN [229], which integrated proposal gen-

eration and detection into a common framework, CNN based detection proposal generation methods have dominated region proposal. It is recommended that new detection proposals should be assessed for object detection, instead of evaluating detection proposals alone.

**(6) Other Factors**
As discussed in Section 9, there are many other factors affecting object detection quality: data augmentation, novel training strategies, combinations of backbone models, multiple detection frameworks, incorporating information from other related tasks, methods for reducing localization error, handling the huge imbalance between positive and negative samples, mining of hard negative samples, and improving loss functions.

## 10.3 Research Directions

Despite the recent tremendous progress in the field of object detection, the technology remains significantly more primitive than human vision and cannot yet satisfactorily address real-world challenges like those of Section 2.2. We see a number of long-standing challenges:

- Working in an open world: being robust to any number of environmental changes, being able to evolve or adapt.
- Object detection under constrained conditions: learning from weakly labeled data or few bounding box annotations, wearable devices, unseen object categories etc.
- Object detection in other modalities: video, RGBD images, 3D point clouds, lidar, remotely sensed imagery *etc*.

Based on these challenges, we see the following directions of future research:

**(1) Open World Learning:** The ultimate goal is to develop object detection capable of accurately and efficiently recognizing and localizing instances in thousands or more object categories in open-world scenes, at a level competitive with the human visual system. Object detection algorithms are unable, in general, to recognize object categories outside of their training dataset, although ideally there should be the ability to recognize novel object categories [144, 95]. Current detection datasets [68, 234, 166] contain only a few dozen to hundreds of categories, significantly fewer than those which can be recognized by humans. New larger-scale datasets [107, 250, 226] with significantly more categories will need to be developed.

**(2) Better and More Efficient Detection Frameworks:** One of the reasons for the success in generic object detection has been the development of superior detection frameworks, both region-based (RCNN [85], Fast RCNN [84], Faster RCNN [229], Mask RCNN [102]) and one-stage detectors (YOLO [227], SSD [175]). Region-based detectors have higher accuracy, one-stage detectors are generally faster and simpler. Object detectors depend heavily on the underlying backbone networks, which have been optimized for image classification, possibly causing a learning bias; learning object detectors from scratch could be helpful for new detection frameworks.

**(3) Compact and Efficient CNN Features:** CNNs have increased remarkably in depth, from several layers (AlexNet [141])

to hundreds of layers (ResNet [101], DenseNet [118]). These networks have millions to hundreds of millions of parameters, requiring massive data and GPUs for training. In order reduce or remove network redundancy, there has been growing research interest in designing compact and lightweight networks [29, 4, 119, 112, 169, 300] and network acceleration [44, 122, 253, 155, 158, 282].

**(4) Automatic Neural Architecture Search:** Deep learning bypasses manual feature engineering which requires human experts with strong domain knowledge, however DCNNs require similarly significant expertise. It is natural to consider automated design of detection backbone architectures, such as the recent Automated Machine Learning (AutoML) [219], which has been applied to image classification and object detection [22, 39, 80, 171, 331, 332].

**(5) Object Instance Segmentation:** For a richer and more detailed understanding of image content, there is a need to tackle pixel-level object instance segmentation [166, 102, 117], which can play an important role in potential applications that require the precise boundaries of individual objects.

**(6) Weakly Supervised Detection:** Current state-of-the-art detectors employ fully supervised models learned from labeled data with object bounding boxes or segmentation masks [69, 166, 234, 166]. However, fully supervised learning has serious limitations, particularly where the collection of bounding box annotations is labor intensive and where the number of images is large. Fully supervised learning is not scalable in the absence of fully labeled training data, so it is essential to understand how the power of CNNs can be leveraged where only weakly / partially annotated data are provided [17, 55, 244].

**(7) Few / Zero Shot Object Detection:** The success of deep detectors relies heavily on gargantuan amounts of annotated training data. When the labeled data are scarce, the performance of deep detectors frequently deteriorates and fails to generalize well. In contrast, humans (even children) can learn a visual concept quickly from very few given examples and can often generalize well [16, 144, 71]. Therefore, the ability to learn from only few examples, *few* shot detection, is very appealing [30, 61, 75, 129, 144, 228, 237]. Even more constrained, *zero* shot object detection localizes and recognizes object classes that have never been seen[16] before [9, 53, 222, 221], essential for life-long learning machines that need to intelligently and incrementally discover new object categories.

**(8) Object Detection in Other Modalities:** Most detectors are based on still 2D images; object detection in other modalities can be highly relevant in domains such as autonomous vehicles, unmanned aerial vehicles, and robotics. These modalities raise new challenges in effectively using depth [36, 211, 289, 286], video [70, 130], and point clouds [217, 218].

**(9) Universal Object Detection:** Recently, there has been increasing effort in learning *universal representations*, those which are effective in multiple image domains, such as natural images, videos, aerial images, and medical CT images [224, 225]. Most such research focuses on image classification, rarely targeting object detection [281], and developed detectors are usually domain

---

[16] Although side information may be provided, such as a wikipedia page or an attributes vector.

specific. Object detection independent of image domain and cross-domain object detection represent important future directions.

The research field of generic object detection is still far from complete. However given the breakthroughs over the past five years we are optimistic of future developments and opportunities.

## 11 Acknowledgments

## References

1. Agrawal P., Girshick R., Malik J. (2014) Analyzing the performance of multilayer neural networks for object recognition. In: ECCV, pp. 329–344 16

2. Alexe B., Deselaers T., Ferrari V. (2010) What is an object? In: CVPR, pp. 73–80 22

3. Alexe B., Deselaers T., Ferrari V. (2012) Measuring the objectness of image windows. IEEE TPAMI 34(11):2189–2202 22

4. Alvarez J., Salzmann M. (2016) Learning the number of neurons in deep networks. In: NIPS, pp. 2270–2278 28

5. Andreopoulos A., Tsotsos J. (2013) 50 years of object recognition: Directions forward. Computer Vision and Image Understanding 117(8):827–891 2, 3, 4

6. Arbeláez P., Hariharan B., Gu C., Gupta S., Bourdev L., Malik J. (2012) Semantic segmentation using regions and parts. In: CVPR, pp. 3378–3385 22

7. Arbeláez P., Pont-Tuset J., Barron J., Marques F., Malik J. (2014) Multi-scale combinatorial grouping. In: CVPR, pp. 328–335 22

8. Azizpour H., Razavian A., Sullivan J., Maki A., Carlsson S. (2016) Factors of transferability for a generic convnet representation. IEEE TPAMI 38(9):1790–1802 16

9. Bansal A., Sikka K., Sharma G., Chellappa R., Divakaran A. (2018) Zero shot object detection. In: ECCV 28

10. Bar M. (2004) Visual objects in context. Nature Reviews Neuroscience 5(8):617–629 20

11. Bell S., Lawrence Z., Bala K., Girshick R. (2016) Inside Outside Net: Detecting objects in context with skip pooling and recurrent neural networks. In: CVPR, pp. 2874–2883 16, 17, 19, 20, 21

12. Belongie S., Malik J., Puzicha J. (2002) Shape matching and object recognition using shape contexts. IEEE TPAMI 24(4):509–522 5

13. Bengio Y., Courville A., Vincent P. (2013) Representation learning: A review and new perspectives. IEEE TPAMI 35(8):1798–1828 2, 3, 6, 14

14. Biederman I. (1972) Perceiving real world scenes. IJCV 177(7):77–80 20

15. Biederman I. (1987) Recognition by components: a theory of human image understanding. Psychological review 94(2):115 6

16. Biederman I. (1987) Recognition by components: a theory of human image understanding. Psychological review 94(2):115 28

17. Bilen H., Vedaldi A. (2016) Weakly supervised deep detection networks. In: CVPR, pp. 2846–2854 28

18. Bodla N., Singh B., Chellappa R., Davis L. S. (2017) SoftNMS improving object detection with one line of code. In: ICCV, pp. 5562–5570 25

19. Borji A., Cheng M., Jiang H., Li J. (2014) Salient object detection: A survey. arXiv: 14115878v1 1:1–26 3

20. Bourdev L., Brandt J. (2005) Robust object detection via soft cascade. In: CVPR, vol 2, pp. 236–243 12

21. Bruna J., Mallat S. (2013) Invariant scattering convolution networks. IEEE TPAMI 35(8):1872–1886 19

22. Cai H., Yang J., Zhang W., Han S., Yu Y. (2018) Path level network transformation for efficient architecture search 28

23. Cai Z., Vasconcelos N. (2018) Cascade RCNN: Delving into high quality object detection. In: CVPR 12, 25, 26, 27

24. Cai Z., Fan Q., Feris R., Vasconcelos N. (2016) A unified multiscale deep convolutional neural network for fast object detection. In: ECCV, pp. 354–370 17

25. Carreira J., Sminchisescu C. (2012) CMPC: Automatic object segmentation using constrained parametric mincuts. IEEE TPAMI 34(7):1312–1328 22

26. Chatfield K., Simonyan K., Vedaldi A., Zisserman A. (2014) Return of the devil in the details: Delving deep into convolutional nets. In: BMVC 24

27. Chavali N., Agrawal H., Mahendru A., Batra D. (2016) Object proposal evaluation protocol is gameable. In: CVPR, pp. 835–844 10, 22

28. Chellappa R. (2016) The changing fortunes of pattern recognition and computer vision. Image and Vision Computing 55:3–5 19

29. Chen G., Choi W., Yu X., Han T., Chandraker M. (2017) Learning efficient object detection models with knowledge distillation. In: NIPS 28

30. Chen H., Wang Y., Wang G., Qiao Y. (2018) LSTD: A low shot transfer detector for object detection. In: AAAI 28

31. Chen K., Pang J., Wang J., Xiong Y., Li X., Sun S., Feng W., Liu Z., Shi J., Ouyang W., et al. (2019) Hybrid task cascade for instance segmentation. In: CVPR 12, 27

32. Chen L., Papandreou G., Kokkinos I., Murphy K., Yuille A. (2015) Semantic image segmentation with deep convolutional nets and fully connected CRFs. In: ICLR 20

33. Chen L., Papandreou G., Kokkinos I., Murphy K., Yuille A. (2018) DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. IEEE TPAMI 40(4):834–848 16, 17, 20

34. Chen Q., Song Z., Dong J., Huang Z., Hua Y., Yan S. (2015) Contextualizing object detection and classification. IEEE TPAMI 37(1):13–27 20

35. Chen X., Gupta A. (2017) Spatial memory for context reasoning in object detection. In: ICCV 20, 21, 27

36. Chen X., Kundu K., Zhu Y., Berneshawi A. G., Ma H., Fidler S., Urtasun R. (2015) 3d object proposals for accurate object class detection. In: NIPS, pp. 424–432 28

37. Chen Y., Li J., Xiao H., Jin X., Yan S., Feng J. (2017) Dual path networks. In: NIPS, pp. 4467–4475 16, 25

38. Chen Y., Rohrbach M., Yan Z., Yan S., Feng J., Kalantidis Y. (2019) Graph based global reasoning networks. In: CVPR 16

39. Chen Y., Yang T., Zhang X., Meng G., Pan C., Sun J. (2019) DetNAS: Neural architecture search on object detection. arXiv:190310979 28

40. Cheng B., Wei Y., Shi H., Feris R., Xiong J., Huang T. (2018) Decoupled classification refinement: Hard false positive suppression for object detection. arXiv:181004002 27

41. Cheng B., Wei Y., Shi H., Feris R., Xiong J., Huang T. (2018) Revisiting RCNN: on awakening the classification power of faster RCNN. In: ECCV 27

42. Cheng G., Zhou P., Han J. (2016) RIFDCNN: Rotation invariant and fisher discriminative convolutional neural networks for object detection. In: CVPR, pp. 2884–2893 19

43. Cheng M., Zhang Z., Lin W., Torr P. (2014) BING: Binarized normed gradients for objectness estimation at 300fps. In: CVPR, pp. 3286–3293 22

44. Cheng Y., Wang D., Zhou P., Zhang T. (2018) Model compression and acceleration for deep neural networks: The principles, progress, and challenges. IEEE Signal Processing Magazine 35(1):126–136 28

45. Chollet F. (2017) Xception: Deep learning with depthwise separable convolutions. In: CVPR, pp. 1800–1807 16, 26

46. Cinbis R., Verbeek J., Schmid C. (2017) Weakly supervised object localization with multi-fold multiple instance learning. IEEE TPAMI 39(1):189–203 11

47. Csurka G., Dance C., Fan L., Willamowski J., Bray C. (2004) Visual categorization with bags of keypoints. In: ECCV Workshop on statistical learning in computer vision 3, 5

48. Dai J., He K., Li Y., Ren S., Sun J. (2016) Instance sensitive fully convolutional networks. In: ECCV, pp. 534–549 23, 24

49. Dai J., He K., Sun J. (2016) Instance aware semantic segmentation via multitask network cascades. In: CVPR, pp. 3150–3158 24

50. Dai J., Li Y., He K., Sun J. (2016) RFCN: object detection via region based fully convolutional networks. In: NIPS, pp. 379–387 9, 12, 16, 22, 26, 35

51. Dai J., Qi H., Xiong Y., Li Y., Zhang G., Hu H., Wei Y. (2017) Deformable convolutional networks. In: ICCV 17, 20, 25, 26, 27

52. Dalal N., Triggs B. (2005) Histograms of oriented gradients for human detection. In: CVPR, vol 1, pp. 886–893 3, 5, 9, 14, 22

53. Demirel B., Cinbis R. G., Ikizler-Cinbis N. (2018) Zero shot object detection by hybrid region embedding. In: BMVC 28

54. Deng J., Dong W., Socher R., Li L., Li K., Li F. (2009) ImageNet: A large scale hierarchical image database. In: CVPR, pp. 248–255 5, 7, 16

55. Diba A., Sharma V., Pazandeh A. M., Pirsiavash H., Van Gool L. (2017) Weakly supervised cascaded convolutional networks. In: CVPR, vol 3, p. 9 28

56. Dickinson S., Leonardis A., Schiele B., Tarr M. (2009) The Evolution of Object Categorization and the Challenge of Image Abstraction in *Object Categorization: Computer and Human Vision Perspectives*. Cambridge University Press 3, 14

57. Ding J., Xue N., Long Y., Xia G., Lu Q. (2018) Learning RoI transformer for detecting oriented objects in aerial images. In: CVPR 20

58. Divvala S., Hoiem D., Hays J., Efros A., Hebert M. (2009) An empirical study of context in object detection. In: CVPR, pp. 1271–1278 20, 27

59. Dollar P., Wojek C., Schiele B., Perona P. (2012) Pedestrian detection: An evaluation of the state of the art. IEEE TPAMI 34(4):743–761 2, 3

60. Donahue J., Jia Y., Vinyals O., Hoffman J., Zhang N., Tzeng E., Darrell T. (2014) DeCAF: A deep convolutional activation feature for generic visual recognition. In: ICML, vol 32, pp. 647–655 16

61. Dong X., Zheng L., Ma F., Yang Y., Meng D. (2018) Few example object detection with model communication. IEEE TPAMI 28

62. Duan K., Bai S., Xie L., Qi H., Huang Q., Tian Q. (2019) CenterNet: Keypoint triplets for object detection. arXiv preprint arXiv:190408189 14

63. Dvornik N., Mairal J., Schmid C. (2018) Modeling visual context is key to augmenting object detection datasets. In: ECCV, pp. 364–380 24

64. Dwibedi D., Misra I., Hebert M. (2017) Cut, paste and learn: Surprisingly easy synthesis for instance detection. In: ICCV, pp. 1301–1310 24

65. Endres I., Hoiem D. (2010) Category independent object proposals 22

66. Enzweiler M., Gavrila D. M. (2009) Monocular pedestrian detection: Survey and experiments. IEEE TPAMI 31(12):2179–2195 2, 3

67. Erhan D., Szegedy C., Toshev A., Anguelov D. (2014) Scalable object detection using deep neural networks. In: CVPR, pp. 2147–2154 10, 22, 23

68. Everingham M., Gool L. V., Williams C., Winn J., Zisserman A. (2010) The pascal visual object classes (voc) challenge. IJCV 88(2):303–338 1, 3, 4, 5, 7, 8, 9, 22, 26, 28

69. Everingham M., Eslami S., Gool L. V., Williams C., Winn J., Zisserman A. (2015) The pascal visual object classes challenge: A retrospective. IJCV 111(1):98–136 7, 8, 26, 28

70. Feichtenhofer C., Pinz A., Zisserman A. (2017) Detect to track and track to detect. In: ICCV, pp. 918–927 28

71. FeiFei L., Fergus R., Perona P. (2006) One shot learning of object categories. IEEE TPAMI 28(4):594–611 28

72. Felzenszwalb P., McAllester D., Ramanan D. (2008) A discriminatively trained, multiscale, deformable part model. In: CVPR, pp. 1–8 9, 22

73. Felzenszwalb P., Girshick R., McAllester D. (2010) Cascade object detection with deformable part models. In: CVPR, pp. 2241–2248 12

74. Felzenszwalb P., Girshick R., McAllester D., Ramanan D. (2010) Object detection with discriminatively trained part based models. IEEE TPAMI 32(9):1627–1645 3, 9, 16, 20, 27

75. Finn C., Abbeel P., Levine S. (2017) Model agnostic meta learning for fast adaptation of deep networks. In: ICML, pp. 1126–1135 28

76. Fischler M., Elschlager R. (1973) The representation and matching of pictorial structures. IEEE Transactions on computers 100(1):67–92 1, 5

77. Fu C.-Y., Liu W., Ranga A., Tyagi A., Berg A. C. (2017) DSSD: Deconvolutional single shot detector. In: arXiv preprint arXiv:1701.06659 14, 17, 18, 19, 26

78. Galleguillos C., Belongie S. (2010) Context based object categorization: A critical survey. Computer Vision and Image Understanding 114:712–722 3, 20, 27

79. Geronimo D., Lopez A. M., Sappa A. D., Graf T. (2010) Survey of pedestrian detection for advanced driver assistance systems. IEEE TPAMI 32(7):1239–1258 2, 3

80. Ghiasi G., Lin T., Pang R., Le Q. (2019) NASFPN: learning scalable feature pyramid architecture for object detection. arXiv:190407392 28

81. Ghodrati A., Diba A., Pedersoli M., Tuytelaars T., Van Gool L. (2015) DeepProposal: Hunting objects by cascading deep convolutional layers. In: ICCV, pp. 2578–2586 22, 23

82. Gidaris S., Komodakis N. (2015) Object detection via a multiregion and semantic segmentation aware CNN model. In: ICCV, pp. 1134–1142 14, 21, 22, 25, 27

83. Gidaris S., Komodakis N. (2016) Attend refine repeat: Active box proposal generation via in out localization. In: BMVC 17, 25

84. Girshick R. (2015) Fast R-CNN. In: ICCV, pp. 1440–1448 2, 9, 10, 11, 14, 15, 16, 22, 24, 25, 26, 28, 35

85. Girshick R., Donahue J., Darrell T., Malik J. (2014) Rich feature hierarchies for accurate object detection and semantic segmentation. In: CVPR, pp. 580–587 2, 3, 5, 9, 10, 11, 13, 14, 15, 16, 22, 24, 25, 26, 28, 35

86. Girshick R., Iandola F., Darrell T., Malik J. (2015) Deformable part models are convolutional neural networks. In: CVPR, pp. 437–446 20, 27

87. Girshick R., Donahue J., Darrell T., Malik J. (2016) Region-based convolutional networks for accurate object detection and segmentation. IEEE TPAMI 38(1):142–158 10, 11, 16

88. Goodfellow I., Shlens J., Szegedy C. (2015) Explaining and harnessing adversarial examples. In: ICLR 7

89. Goodfellow I., Bengio Y., Courville A. (2016) Deep Learning. MIT press 6

90. Grauman K., Darrell T. (2005) The pyramid match kernel: Discriminative classification with sets of image features. In: ICCV, vol 2, pp. 1458–1465 10

91. Grauman K., Leibe B. (2011) Visual object recognition. Synthesis lectures on artificial intelligence and machine learning 5(2):1–181 1, 2, 3

92. Gu J., Wang Z., Kuen J., Ma L., Shahroudy A., Shuai B., Liu T., Wang X., Wang G., Cai J., Chen T. (2017) Recent advances in convolutional neural networks. Pattern Recognition pp. 1–24 2, 3, 6, 14

93. Guillaumin M., Küttel D., Ferrari V. (2014) Imagenet autoannotation with segmentation propagation. International Journal of Computer Vision 110(3):328–348 22

94. Gupta A., Vedaldi A., Zisserman A. (2016) Synthetic data for text localisation in natural images. In: CVPR, pp. 2315–2324 24

95. Hariharan B., Girshick R. B. (2017) Low shot visual recognition by shrinking and hallucinating features. In: ICCV, pp. 3037–3046 28

96. Hariharan B., Arbeláez P., Girshick R., Malik J. (2014) Simultaneous detection and segmentation. In: ECCV, pp. 297–312 24

97. Hariharan B., Arbeláez P., Girshick R., Malik J. (2016) Object instance segmentation and fine grained localization using hypercolumns. IEEE TPAMI 11, 14, 16, 19

98. Harzallah H., Jurie F., Schmid C. (2009) Combining efficient object localization and image classification. In: ICCV, pp. 237–244 9, 22

99. He K., Zhang X., Ren S., Sun J. (2014) Spatial pyramid pooling in deep convolutional networks for visual recognition. In: ECCV, pp. 346–361 2, 10, 15, 16, 35

100. He K., Zhang X., Ren S., Sun J. (2015) Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. In: ICCV, pp. 1026–1034 14

101. He K., Zhang X., Ren S., Sun J. (2016) Deep residual learning for image recognition. In: CVPR, pp. 770–778 3, 12, 14, 15, 26, 27, 28

102. He K., Gkioxari G., Dollár P., Girshick R. (2017) Mask RCNN. In: ICCV 12, 14, 19, 20, 22, 24, 25, 26, 27, 28, 35

103. He T., Tian Z., Huang W., Shen C., Qiao Y., Sun C. (2018) An end to end textspotter with explicit alignment and attention. In: CVPR, pp. 5020–5029 20

104. He Y., Zhu C., Wang J., Savvides M., Zhang X. (2019) Bounding box regression with uncertainty for accurate object detection. In: CVPR 25

105. Hinton G., Salakhutdinov R. (2006) Reducing the dimensionality of data with neural networks. science 313(5786):504–507 1

106. Hinton G., Vinyals O., Dean J. (2015) Distilling the knowledge in a neural network. arXiv:150302531 16, 26

107. Hoffman J., Guadarrama S., Tzeng E. S., Hu R., Donahue J., Girshick R., Darrell T., Saenko K. (2014) LSDA: large scale detection through adaptation. In: NIPS, pp. 3536–3544 28

108. Hoiem D., Chodpathumwan Y., Dai Q. (2012) Diagnosing error in object detectors. In: ECCV, pp. 340–353 8, 25

109. Hosang J., Omran M., Benenson R., Schiele B. (2015) Taking a deeper look at pedestrians. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4073–4082 2

110. Hosang J., Benenson R., Dollr P., Schiele B. (2016) What makes for effective detection proposals? IEEE TPAMI 38(4):814–829 10, 22, 27

111. Hosang J., Benenson R., Schiele B. (2017) Learning nonmaximum suppression. In: ICCV 25

112. Howard A., Zhu M., Chen B., Kalenichenko D., Wang W., Weyand T., Andreetto M., Adam H. (2017) Mobilenets: Efficient convolutional neural networks for mobile vision applications. In: CVPR 15, 27, 28

113. Hu H., Lan S., Jiang Y., Cao Z., Sha F. (2017) FastMask: Segment multiscale object candidates in one shot. In: CVPR, pp. 991–999 23, 24

114. Hu H., Gu J., Zhang Z., Dai J., Wei Y. (2018) Relation networks for object detection. In: CVPR 20, 21, 27

115. Hu J., Shen L., Sun G. (2018) Squeeze and excitation networks. In: CVPR 14, 15

116. Hu P., Ramanan D. (2017) Finding tiny faces. In: CVPR, pp. 1522–1530 2

117. Hu R., Dollár P., He K., Darrell T., Girshick R. (2018) Learning to segment every thing. In: CVPR 28

118. Huang G., Liu Z., Weinberger K. Q., van der Maaten L. (2017) Densely connected convolutional networks. In: CVPR 14, 15, 19, 28

119. Huang G., Liu S., van der Maaten L., Weinberger K. (2018) CondenseNet: An efficient densenet using learned group convolutions. In: CVPR 28

120. Huang J., Rathod V., Sun C., Zhu M., Korattikara A., Fathi A., Fischer I., Wojna Z., Song Y., Guadarrama S., Murphy K. (2017) Speed/accuracy trade offs for modern convolutional object detectors. In: CVPR 15, 24, 26, 27

121. Huang Z., Huang L., Gong Y., Huang C., Wang X. (2019) Mask scoring rcnn. In: CVPR 25

122. Hubara I., Courbariaux M., Soudry D., ElYaniv R., Bengio Y. (2016) Binarized neural networks. In: NIPS, pp. 4107–4115 28

123. Iandola F., Han S., Moskewicz M., Ashraf K., Dally W., Keutzer K. (2016) SqueezeNet: Alexnet level accuracy with 50x fewer parameters and 0.5 mb model size. In: arXiv preprint arXiv:1602.07360 27

124. ILSVRC detection challenge results (2018) http://www.image-net.org/challenges/LSVRC/ 26

125. Ioffe S., Szegedy C. (2015) Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: International Conference on Machine Learning, pp. 448–456 14, 15, 26

126. Jaderberg M., Simonyan K., Zisserman A., et al. (2015) Spatial transformer networks. In: NIPS, pp. 2017–2025 19, 20

127. Jia Y., Shelhamer E., Donahue J., Karayev S., Long J., Girshick R., Guadarrama S., Darrell T. (2014) Caffe: Convolutional architecture for fast feature embedding. In: ACM MM, pp. 675–678 16

128. Jiang B., Luo R., Mao J., Xiao T., Jiang Y. (2018) Acquisition of localization confidence for accurate object detection. In: ECCV, pp. 784–799 25

129. Kang B., Liu Z., Wang X., Yu F., Feng J., Darrell T. (2018) Few shot object detection via feature reweighting. arXiv preprint arXiv:181201866 28

130. Kang K., Ouyang W., Li H., Wang X. (2016) Object detection from video tubelets with convolutional neural networks. In: CVPR, pp. 817–825 28

131. Kim A., Sharma A., Jacobs D. (2014) Locally scale invariant convolutional neural networks. In: NIPS 19

132. Kim K., Hong S., Roh B., Cheon Y., Park M. (2016) PVANet: Deep but lightweight neural networks for real time object detection. In: NIPSW 17

133. Kim Y., Kang B.-N., Kim D. (2018) SAN: learning relationship between convolutional features for multiscale object detection. In: ECCV, pp. 316–331 19

134. Kirillov A., He K., Girshick R., Rother C., Dollár P. (2018) Panoptic segmentation. arXiv preprint arXiv:180100868 27

135. Kong T., Yao A., Chen Y., Sun F. (2016) HyperNet: towards accurate region proposal generation and joint object detection. In: CVPR, pp. 845–853 16, 17, 19, 23

136. Kong T., Sun F., Yao A., Liu H., Lu M., Chen Y. (2017) RON: Reverse connection with objectness prior networks for object detection. In: CVPR 16, 17, 18, 19

137. Kong T., Sun F., Tan C., Liu H., Huang W. (2018) Deep feature pyramid reconfiguration for object detection. In: ECCV, pp. 169–185 17, 18, 19

138. Krähenbühl1 P., Koltun V. (2014) Geodesic object proposals. In: ECCV 22

139. Krasin I., Duerig T., Alldrin N., Ferrari V., AbuElHaija S., Kuznetsova A., Rom H., Uijlings J., Popov S., Kamali S., Malloci M., PontTuset J., Veit A., Belongie S., Gomes V., Gupta A., Sun C., Chechik G., Cai D., Feng Z., Narayanan D., Murphy K. (2017) OpenImages: A public dataset for large scale multilabel and multiclass image classification. Dataset available from https://storagegoogleapiscom/openimages/web/indexhtml 26

140. Krizhevsky A., Sutskever I., Hinton G. (2012) ImageNet classification with deep convolutional neural networks. In: NIPS, pp. 1097–1105 2, 3, 5, 10, 13, 22, 26

141. Krizhevsky A., Sutskever I., Hinton G. (2012) ImageNet classification with deep convolutional neural networks. In: NIPS, pp. 1097–1105 14, 15, 28

142. Kuo W., Hariharan B., Malik J. (2015) DeepBox: Learning objectness with convolutional networks. In: ICCV, pp. 2479–2487 22, 23, 24

143. Kuznetsova A., Rom H., Alldrin N., Uijlings J., Krasin I., PontTuset J., Kamali S., Popov S., Malloci M., Duerig T., et al. (2018) The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. arXiv preprint arXiv:181100982 7, 8, 9

144. Lake B., Salakhutdinov R., Tenenbaum J. (2015) Human level concept learning through probabilistic program induction. Science 350(6266):1332–1338 28

145. Lampert C. H., Blaschko M. B., Hofmann T. (2008) Beyond sliding windows: Object localization by efficient subwindow search. In: CVPR, pp. 1–8 9

146. Law H., Deng J. (2018) CornerNet: Detecting objects as paired keypoints. In: ECCV 14, 16, 27

147. Lazebnik S., Schmid C., Ponce J. (2006) Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In: CVPR, vol 2, pp. 2169–2178 3, 5, 10

148. LeCun Y., Bottou L., Bengio Y., Haffner P. (1998) Gradient based learning applied to document recognition. Proceedings of the IEEE 86(11):2278–2324 2

149. LeCun Y., Bengio Y., Hinton G. (2015) Deep learning. Nature 521:436–444 1, 2, 3, 6, 14

150. Lee C., Xie S., Gallagher P., Zhang Z., Tu Z. (2015) Deeply supervised nets. In: Artificial Intelligence and Statistics, pp. 562–570 15

151. Lenc K., Vedaldi A. (2015) R-CNN minus R. In: BMVC15 11, 35

152. Lenc K., Vedaldi A. (2018) Understanding image representations by measuring their equivariance and equivalence. IJCV 19

153. Li B., Liu Y., Wang X. (2019) Gradient harmonized single stage detector. In: AAAI 25

154. Li H., Lin Z., Shen X., Brandt J., Hua G. (2015) A convolutional neural network cascade for face detection. In: CVPR, pp. 5325–5334 2

155. Li H., Kadav A., Durdanovic I., Samet H., Graf H. P. (2017) Pruning filters for efficient convnets. In: ICLR 28

156. Li H., Liu Y., Ouyang W., XiaogangWang (2018) Zoom out and in network with map attention decision for region proposal and object detection. IJCV 17, 18, 19, 23

157. Li J., Wei Y., Liang X., Dong J., Xu T., Feng J., Yan S. (2017) Attentive contexts for object detection. IEEE Transactions on Multimedia 19(5):944–954 21, 22

158. Li Q., Jin S., Yan J. (2017) Mimicking very efficient network for object detection. In: CVPR, pp. 7341–7349 28

159. Li S. Z., Zhang Z. (2004) Floatboost learning and statistical face detection. IEEE TPAMI 26(9):1112–1123 12

160. Li Y., Wang S., Tian Q., Ding X. (2015) Feature representation for statistical learning based object detection: A review. Pattern Recognition 48(11):3542–3559 3

161. Li Y., Ouyang W., Zhou B., Wang K., Wang X. (2017) Scene graph generation from objects, phrases and region captions. In: ICCV, pp. 1261–1270 27

162. Li Y., Qi H., Dai J., Ji X., Wei Y. (2017) Fully convolutional instance aware semantic segmentation. In: CVPR, pp. 4438–4446 24

163. Li Y., Chen Y., Wang N., Zhang Z. (2019) Scale aware trident networks for object detection. arXiv preprint arXiv:190101892 17, 27

164. Li Z., Peng C., Yu G., Zhang X., Deng Y., Sun J. (2018) DetNet: A backbone network for object detection. In: ECCV 16, 17, 18, 19, 27

165. Li Z., Peng C., Yu G., Zhang X., Deng Y., Sun J. (2018) Light head RCNN: In defense of two stage object detector. In: CVPR 12, 26

166. Lin T., Maire M., Belongie S., Hays J., Perona P., Ramanan D., Dollár P., Zitnick L. (2014) Microsoft COCO: Common objects in context. In: ECCV, pp. 740–755 3, 4, 5, 7, 8, 22, 26, 28

167. Lin T., Dollár P., Girshick R., He K., Hariharan B., Belongie S. (2017) Feature pyramid networks for object detection. In: CVPR 12, 16, 17, 18, 19, 26, 27

168. Lin T., Goyal P., Girshick R., He K., Dollár P. (2017) Focal loss for dense object detection. In: ICCV 14, 19, 25

169. Lin X., Zhao C., Pan W. (2017) Towards accurate binary convolutional neural network. In: NIPS, pp. 344–352 28

170. Litjens G., Kooi T., Bejnordi B., Setio A., Ciompi F., Ghafoorian M., J. van der Laak B. v., Sánchez C. (2017) A survey on deep learning in medical image analysis. Medical Image Analysis 42:60–88 2, 3, 6

171. Liu C., Zoph B., Neumann M., Shlens J., Hua W., Li L., FeiFei L., Yuille A., Huang J., Murphy K. (2018) Progressive neural architecture search. In: ECCV, pp. 19–34 28

172. Liu L., Fieguth P., Guo Y., Wang X., Pietikäinen M. (2017) Local binary features for texture classification: Taxonomy and experimental study. Pattern Recognition 62:135–160 19

173. Liu S., Huang D., Wang Y. (2018) Receptive field block net for accurate and fast object detection. In: ECCV 17

174. Liu S., Qi L., Qin H., Shi J., Jia J. (2018) Path aggregation network for instance segmentation. In: CVPR, pp. 8759–8768 17, 18, 19

175. Liu W., Anguelov D., Erhan D., Szegedy C., Reed S., Fu C., Berg A. (2016) SSD: single shot multibox detector. In: ECCV, pp. 21–37 13, 14, 17, 20, 23, 24, 25, 26, 27, 28, 35

176. Liu Y., Wang R., Shan S., Chen X. (2018) Structure Inference Net: Object detection using scene level context and instance level relationships. In: CVPR, pp. 6985–6994 21, 22

177. Long J., Shelhamer E., Darrell T. (2015) Fully convolutional networks for semantic segmentation. In: CVPR, pp. 3431–3440 11, 12, 16, 19, 23, 24

178. Lowe D. (1999) Object recognition from local scale invariant features. In: ICCV, vol 2, pp. 1150–1157 3, 5, 14

179. Lowe D. (2004) Distinctive image features from scale-invariant keypoints. IJCV 60(2):91–110 3, 5, 22

180. Loy C., Lin D., Ouyang W., Xiong Y., Yang S., Huang Q., Zhou D., Xia W., Li Q., Luo P., et al. (2019) WIDER face and pedestrian challenge 2018: Methods and results. arXiv:190206854 27

181. Lu Y., Javidi T., Lazebnik S. (2016) Adaptive object detection using adjacency and zoom prediction. In: CVPR, pp. 2351–2359 23

182. Luo P., Wang X., Shao W., Peng Z. (2018) Towards understanding regularization in batch normalization. In: ICLR 27

183. Luo P., Ren J., Peng Z., Zhang R., Li J. (2019) Switchable normalization for learning to normalize deep representation. IEEE TPAMI 27

184. Ma J., Shao W., Ye H., Wang L., Wang H., Zheng Y., Xue X. (2018) Arbitrary oriented scene text detection via rotation proposals. IEEE TMM 20(11):3111–3122 20

185. Malisiewicz T., Efros A. (2009) Beyond categories: The visual memex model for reasoning about object relationships. In: NIPS 20, 27

186. Manen S., Guillaumin M., Van Gool L. (2013) Prime object proposals with randomized prim's algorithm. In: CVPR, pp. 2536–2543 22

187. Mikolajczyk K., Schmid C. (2005) A performance evaluation of local descriptors. IEEE TPAMI 27(10):1615–1630 5

188. Mordan T., Thome N., Henaff G., Cord M. (2018) End to end learning of latent deformable part based representations for object detection. IJCV pp. 1–21 17, 20

189. MS COCO detection leaderboard (2018) http://cocodataset.org/#detection-leaderboard 26

190. Mundy J. (2006) Object recognition in the geometric era: A retrospective. in book Toward Category Level Object Recognition edited by J Ponce, M Hebert, C Schmid and A Zisserman pp. 3–28 5

191. Murase H., Nayar S. (1995) Visual learning and recognition of 3D objects from appearance. IJCV 14(1):5–24 5

192. Murase H., Nayar S. (1995) Visual learning and recognition of 3d objects from appearance. IJCV 14(1):5–24 5

193. Murphy K., Torralba A., Freeman W. (2003) Using the forest to see the trees: a graphical model relating features, objects and scenes. In: NIPS 20, 27

194. Newell A., Yang K., Deng J. (2016) Stacked hourglass networks for human pose estimation. In: ECCV, pp. 483–499 14, 19

195. Newell A., Huang Z., Deng J. (2017) Associative embedding: end to end learning for joint detection and grouping. In: NIPS, pp. 2277–2287 14

196. Ojala T., Pietikäinen M., Maenpää T. (2002) Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. IEEE TPAMI 24(7):971–987 5, 22

197. Oliva A., Torralba A. (2007) The role of context in object recognition. Trends in cognitive sciences 11(12):520–527 20, 27

198. Opelt A., Pinz A., Fussenegger M., Auer P. (2006) Generic object recognition with boosting. IEEE TPAMI 28(3):416–431 4

199. Oquab M., Bottou L., Laptev I., Sivic J. (2014) Learning and transferring midlevel image representations using convolutional neural networks. In: CVPR, pp. 1717–1724 6

200. Oquab M., Bottou L., Laptev I., Sivic J. (2015) Is object localization for free? weakly supervised learning with convolutional neural networks. In: CVPR, pp. 685–694 11

201. Osuna E., Freund R., Girosit F. (1997) Training support vector machines: an application to face detection. In: CVPR, pp. 130–136 5

202. Ouyang W., Wang X. (2013) Joint deep learning for pedestrian detection. In: ICCV, pp. 2056–2063 20

203. Ouyang W., Wang X., Zeng X., Qiu S., Luo P., Tian Y., Li H., Yang S., Wang Z., Loy C.-C., et al. (2015) DeepIDNet: Deformable deep convolutional neural networks for object detection. In: CVPR, pp. 2403–2412 9, 17, 20, 21, 22, 26, 27

204. Ouyang W., Wang X., Zhang C., Yang X. (2016) Factors in finetuning deep model for object detection with long tail distribution. In: CVPR, pp. 864–873 25

205. Ouyang W., Wang K., Zhu X., Wang X. (2017) Chained cascade network for object detection. ICCV 12

206. Ouyang W., Zeng X., Wang X., Qiu S., Luo P., Tian Y., Li H., Yang S., Wang Z., Li H., Wang K., Yan J., Loy C. C., Tang X. (2017) DeepID-Net: Object detection with deformable part based convolutional neural networks. IEEE TPAMI 39(7):1320–1334 16, 20

207. Parikh D., Zitnick C., Chen T. (2012) Exploring tiny images: The roles of appearance and contextual information for machine and human object recognition. IEEE TPAMI 34(10):1978–1991 20

208. PASCAL VOC detection leaderboard (2018) http://host.robots.ox.ac.uk:8080/leaderboard/main_bootstrap.php 26

209. Peng C., Xiao T., Li Z., Jiang Y., Zhang X., Jia K., Yu G., Sun J. (2018) MegDet: A large minibatch object detector. In: CVPR, 24, 25, 26

210. Peng X., Sun B., Ali K., Saenko K. (2015) Learning deep object detectors from 3d models. In: ICCV, pp. 1278–1286 24

211. Pepik B., Benenson R., Ritschel T., Schiele B. (2015) What is holding back convnets for detection? In: German Conference on Pattern Recognition, pp. 517–528 28

212. Perronnin F., Sánchez J., Mensink T. (2010) Improving the fisher kernel for large scale image classification. In: ECCV, pp. 143–156 3, 5, 14

213. Pinheiro P., Collobert R., Dollar P. (2015) Learning to segment object candidates. In: NIPS, pp. 1990–1998 22, 23, 24

214. Pinheiro P., Lin T., Collobert R., Dollár P. (2016) Learning to refine object segments. In: ECCV, pp. 75–91 17, 19, 23, 24

215. Ponce J., Hebert M., Schmid C., Zisserman A. (2007) Toward Category Level Object Recognition. Springer 3, 5

216. Pouyanfar S., Sadiq S., Yan Y., Tian H., Tao Y., Reyes M. P., Shyu M., Chen S., Iyengar S. (2018) A survey on deep learning: Algorithms, techniques, and applications. ACM Computing Surveys 51(5):92:1–92:36 6

217. Qi C. R., Su H., Mo K., Guibas L. J. (2017) PointNet: Deep learning on point sets for 3D classification and segmentation. In: CVPR, pp. 652–660 28

218. Qi C. R., Liu W., Wu C., Su H., Guibas L. J. (2018) Frustum pointnets for 3D object detection from RGBD data. In: CVPR, pp. 918–927 28

219. Quanming Y., Mengshuo W., Hugo J. E., Isabelle G., Yiqi H., Yufeng L., Weiwei T., Qiang Y., Yang Y. (2018) Taking human out of learning applications: A survey on automated machine learning. arXiv:181013306 28

220. Rabinovich A., Vedaldi A., Galleguillos C., Wiewiora E., Belongie S. (2007) Objects in context. In: ICCV 20, 27

221. Rahman S., Khan S., Barnes N. (2018) Polarity loss for zero shot object detection. arXiv preprint arXiv:181108982 28

222. Rahman S., Khan S., Porikli F. (2018) Zero shot object detection: Learning to simultaneously recognize and localize novel concepts. In: ACCV 28

223. Razavian R., Azizpour H., Sullivan J., Carlsson S. (2014) CNN features off the shelf: an astounding baseline for recognition. In: CVPR Workshops, pp. 806–813 16

224. Rebuffi S., Bilen H., Vedaldi A. (2017) Learning multiple visual domains with residual adapters. In: Advances in Neural Information Processing Systems, pp. 506–516 28

225. Rebuffi S., Bilen H., Vedaldi A. (2018) Efficient parametrization of multidomain deep neural networks. In: CVPR, pp. 8119–8127 28

226. Redmon J., Farhadi A. (2017) YOLO9000: Better, faster, stronger. In: CVPR 14, 15, 28, 35

227. Redmon J., Divvala S., Girshick R., Farhadi A. (2016) You only look once: Unified, real time object detection. In: CVPR, pp. 779–788 13, 14, 15, 16, 24, 25, 26, 27, 28, 35

228. Ren M., Triantafillou E., Ravi S., Snell J., Swersky K., Tenenbaum J. B., Larochelle H., Zemel R. S. (2018) Meta learning for semisupervised few shot classification. In: ICLR 28

229. Ren S., He K., Girshick R., Sun J. (2015) Faster R-CNN: Towards real time object detection with region proposal networks. In: NIPS, pp. 91–99 9, 11, 13, 14, 16, 20, 22, 23, 24, 25, 26, 27, 28, 35

230. Ren S., He K., Girshick R., Sun J. (2017) Faster RCNN: Towards real time object detection with region proposal networks. IEEE TPAMI 39(6):1137–1149 2, 11, 23, 26

231. Ren S., He K., Girshick R., Zhang X., Sun J. (2017) Object detection networks on convolutional feature maps. IEEE TPAMI 26

232. Rezatofighi H., Tsoi N., Gwak J., Sadeghian A., Reid I., Savarese S. (2019) Generalized intersection over union: A metric and a loss for bounding box regression. In: CVPR 25

233. Rowley H., Baluja S., Kanade T. (1998) Neural network based face detection. IEEE TPAMI 20(1):23–38 5

234. Russakovsky O., Deng J., Su H., Krause J., Satheesh S., Ma S., Huang Z., Karpathy A., Khosla A., Bernstein M., Berg A., Li F. (2015) ImageNet large scale visual recognition challenge. IJCV 115(3):211–252 1, 2, 3, 4, 5, 7, 8, 9, 16, 22, 26, 28

235. Russell B., Torralba A., Murphy K., Freeman W. (2008) LabelMe: A database and web based tool for image annotation. IJCV 77(1-3):157–173 4

236. Schmid C., Mohr R. (1997) Local grayvalue invariants for image retrieval. IEEE TPAMI 19(5):530–535 5

237. Schwartz E., Karlinsky L., Shtok J., Harary S., Marder M., Pankanti S., Feris R., Kumar A., Giries R., Bronstein A. (2019) RepMet: Representative based metric learning for classification and one shot object detection. In: CVPR 28

238. Sermanet P., Kavukcuoglu K., Chintala S., LeCun Y. (2013) Pedestrian detection with unsupervised multistage feature learning. In: CVPR, pp. 3626–3633 5, 22

239. Sermanet P., Eigen D., Zhang X., Mathieu M., Fergus R., LeCun Y. (2014) OverFeat: Integrated recognition, localization and detection using convolutional networks. In: ICLR 2, 3, 10, 13, 15, 24, 35

240. Shang W., Sohn K., Almeida D., Lee H. (2016) Understanding and improving convolutional neural networks via concatenated rectified linear units. In: ICML, pp. 2217–2225 17

241. Shelhamer E., Long J., Darrell T. (2017) Fully convolutional networks for semantic segmentation. IEEE TPAMI 11, 12, 16, 19

242. Shen Z., Liu Z., Li J., Jiang Y., Chen Y., Xue X. (2017) DSOD: Learning deeply supervised object detectors from scratch. In: ICCV 17

243. Shi X., Shan S., Kan M., Wu S., Chen X. (2018) Real time rotation invariant face detection with progressive calibration networks. In: CVPR 20

244. Shi Z., Yang Y., Hospedales T., Xiang T. (2017) Weakly supervised image annotation and segmentation with objects and attributes. IEEE TPAMI 39(12):2525–2538 28

245. Shrivastava A., Gupta A. (2016) Contextual priming and feedback for Faster RCNN. In: ECCV, pp. 330–348 20, 21

246. Shrivastava A., Gupta A., Girshick R. (2016) Training region based object detectors with online hard example mining. In: CVPR, pp. 761–769 25

247. Shrivastava A., Sukthankar R., Malik J., Gupta A. (2017) Beyond skip connections: Top down modulation for object detection. In: CVPR 16, 17, 18, 19, 26

248. Simonyan K., Zisserman A. (2015) Very deep convolutional networks for large scale image recognition. In: ICLR 3, 6, 10, 11, 14, 15, 26

249. Singh B., Davis L. (2018) An analysis of scale invariance in object detection-SNIP. In: CVPR 8, 24, 26

250. Singh B., Li H., Sharma A., Davis L. S. (2018) RFCN 3000 at 30fps: Decoupling detection and classification. In: CVPR 28

251. Singh B., Najibi M., Davis L. S. (2018) SNIPER: Efficient multiscale training. arXiv:180509300 24, 25, 26

252. Sivic J., Zisserman A. (2003) Video google: A text retrieval approach to object matching in videos. In: International Conference on Computer Vision (ICCV), vol 2, pp. 1470–1477 3, 5, 14

253. Song Han W. J. D. Huizi Mao (2016) Deep Compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. In: ICLR 28

254. Sun C., Shrivastava A., Singh S., Gupta A. (2017) Revisiting unreasonable effectiveness of data in deep learning era. In: ICCV, pp. 843–852 16

255. Sun K., Xiao B., Liu D., Wang J. (2019) Deep high resolution representation learning for human pose estimation. In: CVPR 27

256. Sun K., Zhao Y., Jiang B., Cheng T., Xiao B., Liu D., Mu Y., Wang X., Liu W., Wang J. (2019) High resolution representations for labeling pixels and regions. CoRR abs/1904.04514 27

257. Sun S., Pang J., Shi J., Yi S., Ouyang W. (2018) FishNet: A versatile backbone for image, region, and pixel level prediction. In: NIPS, pp. 754–764 16

258. Sun Z., Bebis G., Miller R. (2006) On road vehicle detection: A review. IEEE TPAMI 28(5):694–711 2, 3

259. Sung K., , Poggio T. (1994) Learning and example selection for object and pattern detection. MIT AI Memo (1521) 25

260. Swain M., Ballard D. (1991) Color indexing. IJCV 7(1):11–32 5

261. Szegedy C., Toshev A., Erhan D. (2013) Deep neural networks for object detection. In: NIPS, pp. 2553–2561 10, 13

262. Szegedy C., Reed S., Erhan D., Anguelov D., Ioffe S. (2014) Scalable, high quality object detection. In: arXiv preprint arXiv:1412.1441 23

263. Szegedy C., Liu W., Jia Y., Sermanet P., Reed S., Anguelov D., Erhan D., Vanhoucke V., Rabinovich A. (2015) Going deeper with convolutions. In: CVPR, pp. 1–9 3, 14, 15, 17, 26

264. Szegedy C., Vanhoucke V., Ioffe S., Shlens J., Wojna Z. (2016) Rethinking the inception architecture for computer vision. In: CVPR, pp. 2818–2826 14, 15, 26

265. Szegedy C., Ioffe S., Vanhoucke V., Alemi A. (2017) Inception v4, inception resnet and the impact of residual connections on learning. AAAI pp. 4278–4284 14, 15, 27

266. Torralba A. (2003) Contextual priming for object detection. IJCV 53(2):169–191 20

267. Turk M. A., Pentland A. (1991) Face recognition using eigenfaces. In: CVPR, pp. 586–591 5

268. Tuzel O., Porikli F., Meer P. (2006) Region covariance: A fast descriptor for detection and classification. In: ECCV, pp. 589–600 5

269. TychsenSmith L., Petersson L. (2017) DeNet: scalable real time object detection with directed sparse sampling. In: ICCV 14, 23, 24

270. TychsenSmith L., Petersson L. (2018) Improving object localization with fitness nms and bounded iou loss. In: CVPR 25

271. Uijlings J., van de Sande K., Gevers T., Smeulders A. (2013) Selective search for object recognition. IJCV 104(2):154–171 3, 9, 10, 22

272. Vaillant R., Monrocq C., LeCun Y. (1994) Original approach for the localisation of objects in images. IEE Proceedings Vision, Image and Signal Processing 141(4):245–250 5

273. Van de Sande K., Uijlings J., Gevers T., Smeulders A. (2011) Segmentation as selective search for object recognition. In: ICCV, pp. 1879–1886 22

274. Vaswani A., Shazeer N., Parmar N., Uszkoreit J., Jones L., Gomez A. N., Kaiser Ł., Polosukhin I. (2017) Attention is all you need. In: NIPS, pp. 6000–6010 21

275. Vedaldi A., Gulshan V., Varma M., Zisserman A. (2009) Multiple kernels for object detection. In: ICCV, pp. 606–613 9, 22

276. Viola P., Jones M. (2001) Rapid object detection using a boosted cascade of simple features. In: CVPR, vol 1, pp. 1–8 3, 5, 9, 22

277. Wan L., Eigen D., Fergus R. (2015) End to end integration of a convolution network, deformable parts model and nonmaximum suppression. In: CVPR, pp. 851–859 20, 27

278. Wang H., Wang Q., Gao M., Li P., Zuo W. (2018) Multiscale location aware kernel representation for object detection. In: CVPR 19

279. Wang X., Han T., Yan S. (2009) An HOG-LBP human detector with partial occlusion handling. In: International Conference on Computer Vision, pp. 32–39 3

280. Wang X., Shrivastava A., Gupta A. (2017) A Fast RCNN: Hard positive generation via adversary for object detection. In: CVPR 20, 24

281. Wang X., Cai Z., Gao D., Vasconcelos N. (2019) Towards universal object detection by domain attention. arXiv:190404402 28

282. Wei Y., Pan X., Qin H., Ouyang W., Yan J. (2018) Quantization mimic: Towards very tiny CNN for object detection. In: ECCV, pp. 267–283 28

283. Woo S., Hwang S., Kweon I. (2018) StairNet: Top down semantic aggregation for accurate one shot detection. In: WACV, pp. 1093–1102 19

284. Worrall D. E., Garbin S. J., Turmukhambetov D., Brostow G. J. (2017) Harmonic networks: Deep translation and rotation equivariance. In: CVPR, vol 2 19

285. Wu Y., He K. (2018) Group normalization. In: ECCV, pp. 3–19 27

286. Wu Z., Song S., Khosla A., Yu F., Zhang L., Tang X., Xiao J. (2015) 3D ShapeNets: A deep representation for volumetric shapes. In: CVPR, pp. 1912–1920 28

287. Wu Z., Pan S., Chen F., Long G., Zhang C., Yu P. S. (2019) A comprehensive survey on graph neural networks. arXiv preprint arXiv:190100596 6

288. Xia G., Bai X., Ding J., Zhu Z., Belongie S., Luo J., Datcu M., Pelillo M., Zhang L. (2018) DOTA: a large-scale dataset for object detection in aerial images. In: CVPR, pp. 3974–3983 20

289. Xiang Y., Mottaghi R., Savarese S. (2014) Beyond PASCAL: A benchmark for 3D object detection in the wild. In: WACV, pp. 75–82 28

290. Xiao R., Zhu L., Zhang H. (2003) Boosting chain learning for object detection. In: ICCV, pp. 709–715 5

291. Xie S., Girshick R., Dollár P., Tu Z., He K. (2017) Aggregated residual transformations for deep neural networks. In: CVPR 12, 15, 26, 27

292. Yang B., Yan J., Lei Z., Li S. (2016) CRAFT objects from images. In: CVPR, pp. 6043–6051 21, 22, 23, 25

293. Yang F., Choi W., Lin Y. (2016) Exploit all the layers: Fast and accurate CNN object detector with scale dependent pooling and cascaded rejection classifiers. In: CVPR, pp. 2129–2137 17

294. Yang M., Kriegman D., Ahuja N. (2002) Detecting faces in images: A survey. IEEE TPAMI 24(1):34–58 2, 3

295. Ye Q., Doermann D. (2015) Text detection and recognition in imagery: A survey. IEEE TPAMI 37(7):1480–1500 2, 3

296. Yosinski J., Clune J., Bengio Y., Lipson H. (2014) How transferable are features in deep neural networks? In: NIPS, pp. 3320–3328 16

297. Young T., Hazarika D., Poria S., Cambria E. (2018) Recent trends in deep learning based natural language processing. IEEE Computational Intelligence Magazine 13(3):55–75 6

298. Yu F., Koltun V. (2016) Multiscale context aggregation by dilated convolutions 16

299. Yu F., Koltun V., Funkhouser T. (2017) Dilated residual networks. In: CVPR, vol 2, p. 3 16

300. Yu R., Li A., Chen C., Lai J., et al. (2018) NISP: Pruning networks using neuron importance score propagation. CVPR 28

301. Zafeiriou S., Zhang C., Zhang Z. (2015) A survey on face detection in the wild: Past, present and future. Computer Vision and Image Understanding 138:1–24 2, 3

302. Zagoruyko S., Lerer A., Lin T., Pinheiro P., Gross S., Chintala S., Dollár P. (2016) A multipath network for object detection. In: BMVC 17, 22, 23

303. Zeiler M., Fergus R. (2014) Visualizing and understanding convolutional networks. In: ECCV, pp. 818–833 6, 14, 15, 20

304. Zeng X., Ouyang W., Yang B., Yan J., Wang X. (2016) Gated bidirectional cnn for object detection. In: ECCV, pp. 354–369 21, 22, 27

305. Zeng X., Ouyang W., Yan J., Li H., Xiao T., Wang K., Liu Y., Zhou Y., Yang B., Wang Z., Zhou H., Wang X. (2017) Crafting gbdnet for object detection. IEEE TPAMI 20, 21, 22, 27

306. Zhang K., Zhang Z., Li Z., Qiao Y. (2016) Joint face detection and alignment using multitask cascaded convolutional networks. IEEE SPL 23(10):1499–1503 2

307. Zhang L., Lin L., Liang X., He K. (2016) Is faster RCNN doing well for pedestrian detection? In: ECCV, pp. 443–457 2

308. Zhang S., Wen L., Bian X., Lei Z., Li S. (2018) Single shot refinement neural network for object detection. In: CVPR 17, 18, 19

309. Zhang S., Yang J., Schiele B. (2018) Occluded pedestrian detection through guided attention in CNNs. In: CVPR, pp. 2056–2063 20

310. Zhang X., Yang Y., Han Z., Wang H., Gao C. (2013) Object class detection: A survey. ACM Computing Surveys 46(1):10:1–10:53 1, 2, 3, 4, 20

311. Zhang X., Li Z., Change Loy C., Lin D. (2017) PolyNet: a pursuit of structural diversity in very deep networks. In: CVPR, pp. 718–726 21

312. Zhang X., Zhou X., Lin M., Sun J. (2018) ShuffleNet: an extremely efficient convolutional neural network for mobile devices. In: CVPR 27

313. Zhang Z., Geiger J., Pohjalainen J., Mousa A. E., Jin W., Schuller B. (2018) Deep learning for environmentally robust speech recognition: An overview of recent developments. ACM Trans Intell Syst Technol 9(5):49:1–49:28 6

314. Zhang Z., Qiao S., Xie C., Shen W., Wang B., Yuille A. (2018) Single shot object detection with enriched semantics. In: CVPR 16

315. Zhao Q., Sheng T., Wang Y., Tang Z., Chen Y., Cai L., Ling H. (2019) M2Det: A single shot object detector based on multilevel feature pyramid network. In: AAAI 17, 18, 19

316. Zheng S., Jayasumana S., Romera-Paredes B., Vineet V., Su Z., Du D., Huang C., Torr P. (2015) Conditional random fields as recurrent neural networks. In: ICCV, pp. 1529–1537 20

317. Zhou B., Khosla A., Lapedriza A., Oliva A., Torralba A. (2015) Object detectors emerge in deep scene CNNs. In: ICLR 11, 16

318. Zhou B., Khosla A., Lapedriza A., Oliva A., Torralba A. (2016) Learning deep features for discriminative localization. In: CVPR, pp. 2921–2929 11

319. Zhou B., Lapedriza A., Khosla A., Oliva A., Torralba A. (2017) Places: A 10 million image database for scene recognition. IEEE Trans Pattern Analysis and Machine Intelligence 8, 16, 26

320. Zhou J., Cui G., Zhang Z., Yang C., Liu Z., Sun M. (2018) Graph neural networks: A review of methods and applications. arXiv preprint arXiv:181208434 6

321. Zhou P., Ni B., Geng C., Hu J., Xu Y. (2018) Scale transferrable object detection. In: CVPR 15, 17, 19

322. Zhou Y., Liu L., Shao L., Mellor M. (2016) DAVE: A unified framework for fast vehicle detection and annotation. In: ECCV, pp. 278–293 2

323. Zhou Y., Ye Q., Qiu Q., Jiao J. (2017) Oriented response networks. In: CVPR, pp. 4961–4970 19

324. Zhu X., Vondrick C., Fowlkes C., Ramanan D. (2016) Do we need more training data? IJCV 119(1):76–92 14

325. Zhu X., Tuia D., Mou L., Xia G., Zhang L., Xu F., Fraundorfer F. (2017) Deep learning in remote sensing: A comprehensive review and list of resources. IEEE Geoscience and Remote Sensing Magazine 5(4):8–36 6

326. Zhu Y., Urtasun R., Salakhutdinov R., Fidler S. (2015) SegDeepM: Exploiting segmentation and context in deep neural networks for object detection. In: CVPR, pp. 4703–4711 20, 21

327. Zhu Y., Zhao C., Wang J., Zhao X., Wu Y., Lu H. (2017) CoupleNet: Coupling global structure with local parts for object detection. In: ICCV 21, 22

328. Zhu Y., Zhou Y., Ye Q., Qiu Q., Jiao J. (2017) Soft proposal networks for weakly supervised object localization. In: ICCV, pp. 1841–1850 22

329. Zhu Z., Liang D., Zhang S., Huang X., Li B., Hu S. (2016) Traffic sign detection and classification in the wild. In: CVPR, pp. 2110–2118 2

330. Zitnick C., Dollár P. (2014) Edge boxes: Locating object proposals from edges. In: ECCV, pp. 391–405 22

331. Zoph B., Le Q. (2017) Neural architecture search with reinforcement learning 28

332. Zoph B., Vasudevan V., Shlens J., Le Q. (2018) Learning transferable architectures for scalable image recognition. In: CVPR, pp. 8697–8710 28

**Table 11** Summary of properties and performance of milestone detection frameworks for generic object detection. See Section 5 for a detailed discussion. Some architectures are illustrated in Fig. 13. The properties of the backbone DCNNs can be found in Table 6.

| | Detector Name | RP | Backbone DCNN | Input ImgSize | VOC07 Results | VOC12 Results | Speed (FPS) | Published In | Source Code | Highlights and Disadvantages |
|---|---|---|---|---|---|---|---|---|---|---|
| **Region based (Section 5.1)** | RCNN [85] | SS | AlexNet | Fixed | 58.5 (07) | 53.3 (12) | < 0.1 | CVPR14 | Caffe Matlab | **Highlights:** First to integrate CNN with RP methods; Dramatic performance improvement over previous state of the art. **Disadvantages:** Multistage pipeline of sequentially-trained (External RP computation, CNN finetuning, each warped RP passing through CNN, SVM and BBR training); Training is expensive in space and time; Testing is slow. |
| | SPPNet [99] | SS | ZFNet | Arbitrary | 60.9 (07) | — | < 1 | ECCV14 | Caffe Matlab | **Highlights:** First to introduce SPP into CNN architecture; Enable convolutional feature sharing; Accelerate RCNN evaluation by orders of magnitude without sacrificing performance; Faster than Overfeat. **Disadvantages:** Inherit disadvantages of RCNN; Does not result in much training speedup; Finetuning not able to update the CONV layers before SPP layer. |
| | Fast RCNN [84] | SS | AlexNet VGGM VGG16 | Arbitrary | 70.0 (VGG) (07+12) | 68.4 (VGG) (07++12) | < 1 | ICCV15 | Caffe Python | **Highlights:** First to enable end-to-end detector training (ignoring RP generation); Design a RoI pooling layer; Much faster and more accurate than SPPNet; No disk storage required for feature caching. **Disadvantages:** External RP computation is exposed as the new bottleneck; Still too slow for real time applications. |
| | Faster RCNN [229] | RPN | ZFnet VGG | Arbitrary | 73.2 (VGG) (07+12) | 70.4 (VGG) (07++12) | < 5 | NIPS15 | Caffe Matlab Python | **Highlights:** Propose RPN for generating nearly cost-free and high quality RPs instead of selective search; Introduce translation invariant and multiscale anchor boxes as references in RPN; Unify RPN and Fast RCNN into a single network by sharing CONV layers; An order of magnitude faster than Fast RCNN without performance loss; Can run testing at 5 FPS with VGG16. **Disadvantages:** Training is complex, not a streamlined process; Still falls short of real time. |
| | RCNN⊖R [151] | New | ZFNet +SPP | Arbitrary | 59.7 (07) | — | < 5 | BMVC15 | — | **Highlights:** Replace selective search with static RPs; Prove the possibility of building integrated, simpler and faster detectors that rely exclusively on CNN. **Disadvantages:** Falls short of real time; Decreased accuracy from poor RPs. |
| | RFCN [50] | RPN | ResNet101 | Arbitrary | 80.5 (07+12) 83.6 (07+12+CO) | 77.6 (07++12) 82.0 (07++12+CO) | < 10 | NIPS16 | Caffe Matlab | **Highlights:** Fully convolutional detection network; Design a set of position sensitive score maps using a bank of specialized CONV layers; Faster than Faster RCNN without sacrificing much accuracy. **Disadvantages:** Training is not a streamlined process; Still falls short of real time. |
| | Mask RCNN [102] | RPN | ResNet101 ResNeXt101 | Arbitrary | 50.3 (ResNeXt101) (COCO Result) | | < 5 | ICCV17 | Caffe Matlab Python | **Highlights:** A simple, flexible, and effective framework for object instance segmentation; Extends Faster RCNN by adding another branch for predicting an object mask in parallel with the existing branch for BB prediction; Feature Pyramid Network (FPN) is utilized; Outstanding performance. **Disadvantages:** Falls short of real time applications. |
| **Unified (Section 5.2)** | OverFeat [239] | — | AlexNet like | Arbitrary | — | — | < 0.1 | ICLR14 | c++ | **Highlights:** Convolutional feature sharing; Multiscale image pyramid CNN feature extraction; Won the ISLVRC2013 localization competition; Significantly faster than RCNN. **Disadvantages:** Multi-stage pipeline sequentially trained; Single bounding box regressor; Cannot handle multiple object instances of the same class; Too slow for real time applications. |
| | YOLO [227] | — | GoogLeNet like | Fixed | 66.4 (07+12) | 57.9 (07++12) | < 25 (VGG) | CVPR16 | DarkNet | **Highlights:** First efficient unified detector; Drop RP process completely; Elegant and efficient detection framework; Significantly faster than previous detectors; YOLO runs at 45 FPS, Fast YOLO at 155 FPS; **Disadvantages:** Accuracy falls far behind state of the art detectors; Struggle to localize small objects. |
| | YOLOv2 [226] | — | DarkNet | Fixed | 78.6 (07+12) | 73.5 (07++12) | < 50 | CVPR17 | DarkNet | **Highlights:** Propose a faster DarkNet19; Use a number of existing strategies to improve both speed and accuracy; Achieve high accuracy and high speed; YOLO9000 can detect over 9000 object categories in real time. **Disadvantages:** Not good at detecting small objects. |
| | SSD [175] | — | VGG16 | Fixed | 76.8 (07+12) 81.5 (07+12+CO) | 74.9 (07++12) 80.0 (07++12+CO) | < 60 | ECCV16 | Caffe Python | **Highlights:** First accurate and efficient unified detector; Effectively combine ideas from RPN and YOLO to perform detection at multi-scale CONV layers; Faster and significantly more accurate than YOLO; Can run at 59 FPS; **Disadvantages:** Not good at detecting small objects. |

Abbreviations in this table: Region Proposal (RP), Selective Search (SS), Region Proposal Network (RPN), RCNN⊖R represents "RCNN minus R" and used a trivial RP method. Training data: "07"←VOC2007 trainval; "07T"←VOC2007 trainval and test; "12"←VOC2012 trainval; "CO"←COCO trainval; The "Speed" column roughly estimates the detection speed with a single Nvidia Titan X GPU.