



# Deep Learning Approaches for Object Detection

Sushma Jaiswal<sup>1</sup>, Tarun Jaiswal<sup>2\*</sup>

<sup>1</sup>Guru Ghasidas Central University, Bilaspur (C.G.), India

<sup>2</sup>Department of Computer Applications, NIT Raipur, Raipur, India

E-mail: tjaiswal\_1207@yahoo.com

---

**Abstract:** In computer vision, object detection is a very important, exciting and mind-blowing study. Object detection work in numerous fields such as observing security, independently/autonomous driving and etc. Deep-learning based object detection techniques have developed at a very fast pace and have attracted the attention of many researchers. The main focus of the 21st century is the development of the object-detection framework, comprehensively and genuinely. In this investigation, we initially investigate and evaluate the various object detection approaches and designate the benchmark datasets. We also delivered the wide-ranging general idea of object detection approaches in an organized way. We covered the first and second stage detectors of object detection methods. And lastly, we consider the construction of these object detection approaches to give dimensions for further research.

**Keywords:** object detection methods, deep learning, convolutional-neural-network (CNN), computer vision, recurrent neural network (RNN)

---

## 1. Introduction

Object-detection is a scientifically and mechanically very encouraging and essentially valuable field of the machine/computer vision. In an image, the object-detection task classifies different types of objects<sup>[1]</sup>. A huge victory has been achieved in the controlled framework for object detection problems but in an uncontrolled environment, the unruly continue mysterious such as occlusion, arbitrary viewpoint and cluttered environments. For example, it is very easy to train the robot for recognizing the tea/coffee machine with nothing else in the image, and here the problem arises when the robot detects other apparatuses of the kitchen, in this situation the recognition development in such consequence is very challenging. Yet there is not precise clarification has been originated. So the object detection and recognition area attract the researcher or investigator from the last 02 eras. Object detection is a multi-disciplinary investigation area and often comprises the pitches of digital image processing, computer/machine vision, machine/deep learning, cognitive science, graph/topology, statistics/probability, computer-oriented-optimization-techniques, etc. Because the object detection approaches development is so varied and time tedious job, so we need all state-of-the-art approaches under the one umbrella. This investigation is an effort to concisely encapsulate the numerous features of object detection and the foremost stages elaborate for the greatest object detection procedure or structure. The object detection provinces consist of multi-grouping detection, edge, point, salient, pose, scene, text face and pedestrian detection, etc. Currently, the scene understanding in fashioning of the recent scenario, i.e. safekeeping, armed, robotics, traffic monitoring, shipping and medicinal ground. In recent times, deep learning (DL) techniques<sup>[2]</sup> have materialized as commanding approaches for learning feature map mechanically from the information. In specific, these systems have delivered noteworthy enhancement for the discovery of the object. The problematic situation has involved huge consideration in the preceding 10 years, whereas it has been considered used for eras thru cognitive behaviour, neuroscientists, and engineers. In prevalent studies<sup>[3]</sup> domain-specific object detection categorizes in 02 ways, named as one-stage (You Only Look Once (YOLO), Single Shot Detector (SSD)) and two-stage-detection (Faster Regions based on Convolutional Neural Networks (R-CNN)). The two-stage detector is superior accuracy as well as localization, however, the one-stage attain a superior extrapolation rate. The two-stage categories obtained thru Region of Interest (ROI) Pooling Layer. Region-Proposal-Network (RPN) is known as the first phase of Faster R-CNN that suggested bounding box, and in the second phase, feature maps are taken out thru ROI Pooling task of every candidate box designated for classification and regression scheme. Detection of object majorly classified in 02 parts, named as generative and discriminative<sup>[4]</sup>. Generative classifier contains a probability scheme designated for the pose irregularity

---

and composed by appearance framework. The discriminative classifier can discriminate the images/sub-images covered by the object. And to reduce the error during the training time, a type of classifier is used that works on regularization bias to avoid overfitting. The authors [5] concentrated on satellite image object detection that plays a vital role in numerous areas, such as surveillance of the airport, vessel traffic-intensive care, armed and defense, transportation fortitude and urban studies. Through satellite sensors, remote sensing images are obtained, and they are very crucial and intricate the reason behind are arbitrary variations, clutter, occlusion, illumination, interference and high altitude.

Object detection has a wide variety of applications in enterprises and industries, with use cases ranging from personal/individual security to efficiency and productivity in the workplace scenario. Object detection is useful in many domains of Computer vision and image processing (see Figure 1). Considerable challenges stay on the pitch of object detection. The opportunities are endless and incalculable when it derives to upcoming use cases for object detection.

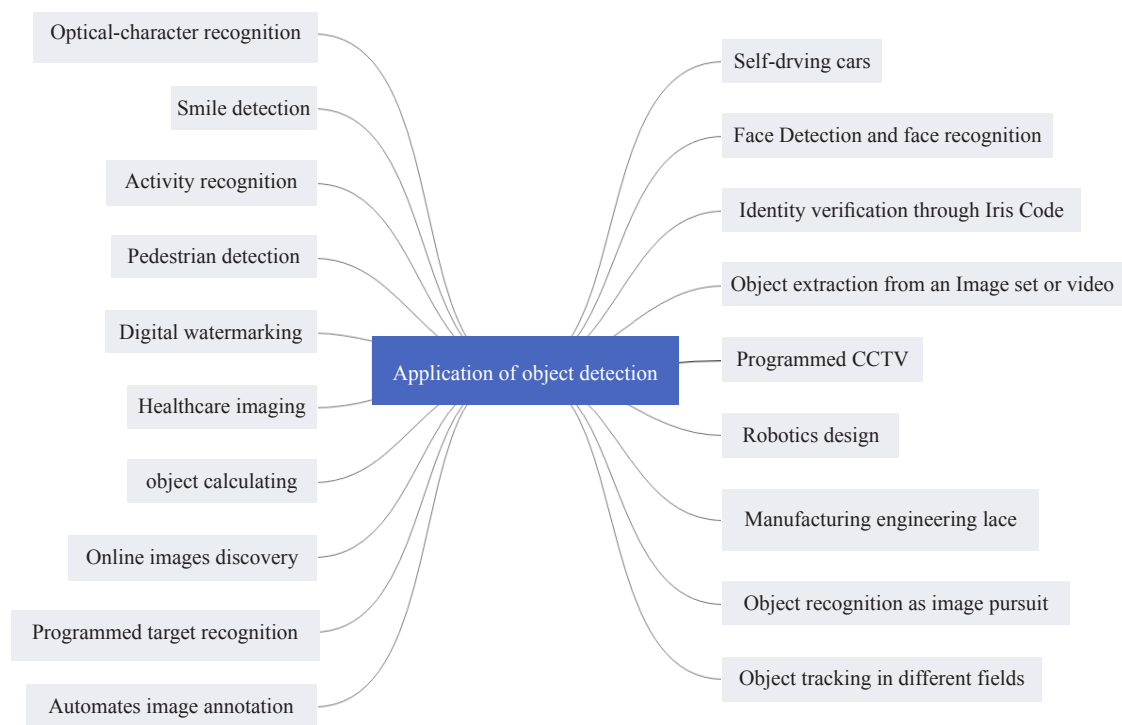


Figure 1. Real-time applications of object detection

The rest of this paper is systematized as follows. In section-II, we review the related approaches of various object detection methods. The experiment and evaluation are presented in section III, whereas section IV concluded the research.

## 2. Object detection approaches

Object-detection is a very popular technique of computer visualization by which we can locate and identify objects which are embedded in a picture or video. With this kind of object identification and localization approach, we can even count the object as well as find the actual position of the object. Object detection allows us to classify the types of things found while also locating instances of them within the image. For this purpose, various object-detection approaches are available named as CNN, YOLO, and fast R-CNN, etc. In this segment, we review the antiquity of object detection in manifold scenarios, with revolutionary detectors, object-detection datasets, metrics and scales, and the development of crucial techniques.

There are numerous kinds of object detection approaches, such as *Conventional approaches* and *Modern approaches*. These frameworks quietly differ from each other in terms of precision, speed and H/W configuration.

### 2.1 Conventional approach

The conventional approaches contain two phases-Feature Extraction and object classification.

### 2.1.1 Feature extraction phase

This phase includes various techniques such as Viola-Jones, Scale Invariant Feature Transform (SIFT) and Histograms of oriented gradients (HOG) etc. to extract the feature for object detection purposes.

*Viola-Jones object detection.* The very most famous method for image-based face detection is known as Viola-Jones object detection (see Algorithm 1), works on full view frontal-upright image faces. In this method, the whole face image requisite point concerning the camera and have not to be slanting to any crosswise. The viola jones follows the four main phases:

1. Haar feature selection follows the demesne of image-based object detection.
2. Integral image creation, in this phase image detector, is calculated first.
3. Adaboost is training or learning algorithm, this phase selects a trifling number of precarious pictorial structures from an enormous set of possible features.
4. Cascading classifier permits background regions of an image to be rapidly detached while more time spent on computation<sup>[6]</sup>.

Algorithm 1: Viola-Jones Face Detection Algorithm

1. Input: original test image
2. Output: image with face indicators as rectangles
3. For  $i \leftarrow 1$  to num of scales in the pyramid of images do
4.     Downsample image to create  $image_i$
5.     Compute integral image,  $image_{ii}$
6. For  $j \leftarrow 1$  to num of shift steps of sub-window do
7.     For  $k \leftarrow 1$  to num of stages in cascade classifier do
8.         For  $l \leftarrow 1$  to num of filters of stage  $k$  do
9.             Filter detection sub-window
10.             Accumulate filter outputs
11.             end for
12.             If accumulation fails per-stage threshold then
13.                 Reject sub-window as face
14.                 Break this  $k$  for loop
15.             end if
16.         end for
17.         If sub-window passed all per-stage checks then
18.             Accept this sub-window as a face
19.         end if
20.     end for
21. end for

*Scale Invariant Feature Transform (SIFT).* The authors developed an image descriptor for matching and recognition for the 3D scene and view-based object recognition. The method is invariant to transformation (translation, rotations and scaling) in the image region and illumination variants. The developed method is very beneficial for matching and recognition of an object in actual circumstances<sup>[7]</sup>.

*Histograms of oriented gradients (HOG).* HOG is stimulated via the great discriminatory function of local position-dependent histograms, and the SIFT descriptor described as customary gradient orientation histograms and calculated over a grid in the image. SIFT descriptor is a local image descriptor while HOG descriptor is a regional image-descriptor. In this sense, the HOG descriptor is precisely correlated to the regional receptive arena histograms, which are simplified over sub-regions in the image domain. The authors described two variants of HOG operator-the first one is local variants calculated over R-HOG while the second one is accumulated over C-HOG<sup>[7]</sup>.

*PVANET.* This investigation<sup>[8]</sup> aimed to enhance the accurateness in multi-grouping object detection jobs whereas it decreases the cost thru familiarizing and integration of modern research mechanisms with prevalent techniques. For this purpose, the authors used 03 main frameworks for innovations, firstly, CNN feature abstraction; secondly, region-proposal; thirdly, the region of interest (ROI)-Classification. The authors also reconstruct the feature extraction module. The intended net is trained via batch normalization and residual connection. The network learning speed is grounded on plateau detection. For this purpose authors configured the system Intel i7-6700KCPU thru a single-core and 46 ms per image on NVIDIATitan-XGPU, captivating merely 750 ms/image. The experiment validates the method on 02 datasets VOC 2007

and 2012, achieved 83.8% mAP and 82.5% mAP.

### 2.1.2 Classification phase

This phase includes the support-vector-machine (SVM) and the Over-Feat method for object classification.

The authors used the support-vector-machine (SVM) to identify the object (see Algorithm 2), and this method uses the reduction of the future vector by kernel principal component analysis [9].

Algorithm 2: The Adaptive One-Class SVM Algorithm (AOSVM)

1. Input:  $C$ , other kernel hyper parameters,  $\lambda$ ,  $M (= 10^9)$ ,  $c_1 = x_1$
2. Algorithm initialization:  $K_R = k(c_1, c_1)$ ,  $P^{(0)} = K_R^{-1}$ ,  $q^{(0)} = 0$ ,  $\beta^{(0)} = 0$
3. For  $n = 1, 2, \dots$  (as new data is available)
  - 3.1  $k_n = [k(c_1, x_n), \dots, k(c_R, x_n)]^T$
  - 3.2  $f^{(n)}(x_n) = k_n^T \beta^{(n-1)} - 1$  and  $e_n^{(n)} = -f^{(n)}(x_n)$

$$3.3 \ a_n^{(n)} = \begin{cases} 0 & e_n^{(n)} < 0 \\ M & 0 \leq e_n^{(n)} < C/M \\ \frac{C}{e_n^{(n)}} & e_n^{(n)} \geq C/M \end{cases}$$

$$3.4 \ q^{(n)} = \lambda q^{(n-1)} + 2k_n a_n^{(n)}$$

$$3.5 \ k_n \leftarrow k_n + g_n, \text{ where } g_n \sim N(0, (1 - \lambda)K_R)$$

$$3.6 \ \text{Update } P^{(n)} \text{ according to (16) to update solution } \beta^{(n)} = P^{(n)} q^{(n)}$$

3.7 If  $x_n$  is support vector ( $e_n^{(n)} > 0$ ) then

\* compute residuals and update  $A_n^{bag}$  and  $A_n^{base}$

\* If  $\{A_n^{bag}\}_r > v^{grow}$  then add  $x_n$  to base

\* if the base has changed, update structures

*OverFeat*. The authors [10] introduced the novel method based on multi-scale, sliding window techniques and its achieved 4th position in the classification, 1st position in localization and detection with ILSVRC 2013 datasets. The authors also showed ConvNets, which were used for localization and detection. There is numerous significant enhancement in the proposed methods: -used  $\ell_2$  loss instead of elevating the intersection-over-union (IOU) standard, Substitute Parameterizations of the bounding-box may help to integrate the results and the entire the net. ConvNets can handle further stimulating chores and a slight adjustment to net intended for classification.

## 2.2 Modern approach

The modern approach can reduce the drawback of conventional approaches. The modern approach contains two phases-region proposals based approach and Regression or Classification based approach.

### 2.2.1 Region proposal based approach

*R-CNN*. The authors [11] considered the very easy and accessible object detection approach and improved the performance over prevalent results. The authors measured the accurateness by two scales. First, applied huge dimension CNN to inner section and second, the prototype for training huge dimension CNN when captioning data is rare. By integrating the inner region and CNN achieved significant outcomes.

*SPP-Net*. The authors [12] introduced the spatial pyramid pooling strategy and pyramid pooling vigorously for object deformations, so net framework known as SPP-net able to produce a fixed-length demonstration nevertheless of image-size or scale. SPP-net enhanced all CNN-based image classification approaches. The experimental validation based on the 02 datasets, namely PascalVO C2007 & Caltech101. SPP-net accomplishes prevalent classification outcomes using an exclusive full-image demonstration and nope fine-tuning. SPP-net extracts the feature-map from the entire image just once. The authors used SPP-net and decrease the computational complexity. The practically SPP-net is better as compared to the R-CNN technique. In ImageNet-Large-Scale-Visual-Recognition-Challenge (ILSVRC) year 2014, SPP-net achieved 2nd position in detection and 3rd in classification from a total 38th entries.

*Deep Conv. Net via Bayesian-Optimization and Structured-Prediction*. The authors [13] represented the high capacity CNN framework and described some drawbacks of the defined method. The authors also reported the localization problem via-used a search algorithm arranged with Bayesian optimization. With experimental validation, the authors proved the represented approach enhanced the detection accuracy over the baseline technique on 02 datasets, i.e. PASCAL VOC

2007/2012 datasets.

*MR-CNN*. The authors <sup>[14]</sup> performed an object detection system that encrypts the semantic segmentation-aware features with the multi-region (MR) deep-convolutional NN network (CNN). The central objective of the applied CNN approach is-a diverse set of discriminative presence elements and demonstrations localization sensitivity. The proposed approach has given superior localization accurateness. The Authors applied the proposed method on 02 standard datasets, i.e. PASCAL VOC2007 & PASCAL VOC2012 and succeed mAP 78.2 and 73.9, respectively.

*DeepBox*. The authors introduced <sup>[15]</sup> the novel methods DeepBox for object detection. In this method, CNN is used by the DeepBox for further improvement of the bottom-up approach. The authors expended 04 layers CNN for huge network arrangement and faster access. DeepBox approach enhanced the bottom-up position. This enhancement oversimplifies to groups of the CNN and has a central expansion of 4.5-point in detection mAP. The proposed execution attains excellent performance.

*Fast R-CNN*. The authors <sup>[16]</sup> considered the Fast-Region-based Convolutional-Network technique (Fast R-CNN) that is able to perceive an object. Fast R-CNN enhanced the rapidity of testing/training and detection accurateness. The proposed method is proficient on the very deep VGG16 network and 213× quicker at testing-time and attains the great maps on the standard dataset, i.e. PASCALVOC 2012. Fast R-CNN trains VGG16 3× faster, tests 10× faster, and is more precise than SPPnet approach.

*DeepProposal*. The Authors <sup>[17]</sup> assessed the superiority of the activation layers of a convolutional-neural-network (CNN). For this purpose, the authors created the assumption in a sliding-window method over diverse activation layers. The authors also showed the last layer of a convolution had high precision value and poor localization while the first layer of convolution has superior localization and low precision value. On grounded assumptions, the authors construct the inverse cascade, which is able to work on first to last convolution layer of the CNN Network. After the extensive experiments, authors concluded-The method used similar features for detection purpose, combined features by integral images and eluded the unnecessary evaluation owed to the inverse coarse-to-fine cascade. The proposed method beat its previous method.

*Faster R-CNN, RPN*. In this investigation, the authors <sup>[18]</sup> focused on Region-Proposal-Network (RPN) for object detection. And this method took the entire-image convolutional topographies by the network. RPN sequence forecast the region and its scores at every scale. RPN used superior region schemes and give to Fast R-CNN for discovery. RPN and Fast R-CNN trained towards convolutional features. The Authors validates the proposed method by applying VGG-16 and found the 5 fps on a Graphics Processing Unit (GPU) whereas PASCALVOC 2007 achieved 73.2% mAP and PASCAL VOC 2012 70.4% mAP by 300 proposals each image.

*AZNet*. The authors <sup>[19]</sup> represented the search policy that adaptively shows the computational resources to segments in the direction of comprehending objects (see Algorithm 3). The authors differentiate the approaches that grounded on immovable anchor direction, and the proposed method familiarizes to situations of sparse and slight objects. The authors performed the experimental evaluation and provided the proposed method outperformed than Faster R-CNN method, whereas the used conditions 02 orders of enormousness smaller quantity anchors on mediocre.

Algorithm 3: Adaptive search with AZNet

Data: Input image  $x$  (the whole image region  $b_x$ ).  $Y_k$  is the region proposed at step  $k$ ,  $Y^K$  are the accumulated region proposals up to step  $x$ ,  $Z_k$  are the regions to further zoom in to at step  $k$ ,  $B_k$  are anchor regions at step  $k$ .

Result: Region proposal at termination  $Y^K$ .

Initialization :  $B_o \leftarrow \{b_x\}$ ,  $Y^0 \leftarrow 0$ ,  $k \leftarrow 0$

While ( $B_k$  is not an emptyset) do

Initialize  $Y_k$  and  $Z_k$  as empty sets.

For each  $b \in B_k$  do

Compute adjacency predictions  $A_b$  and the zoom indicator  $zb$  using AZNet.

Include all  $a \in A_b$  confidence scores into  $Y_k$ .

Include  $b$  into  $Z_k$  if  $zb$  is above the threshold.

End

$Y^k \leftarrow Y^{k-1} \cup Y_k$

$B_{k+1} \leftarrow \text{Divide-Regions}(Z_k)$

$k \leftarrow k + 1$

End

$k \leftarrow k - 1$

*ION.* In the interior and exterior ROI, the authors considered the Inside-Outside Net (ION) object detection approach<sup>[20]</sup>. Exterior ROI cohesively exhausted by spatial recurrent neural networks with background information, whereas interior ROI costumed skip pooling to find out material at manifold balances and levels of abstraction. The extensive experiments proved in terms of proposed methods enhanced the object detection on PASCAL VOC 2012 as of 73.9% to 76.4% mAP and MS COCO dataset as of 19.7% to 33.1% mAP. During the competition of Microsoft Common Objects in COntext (MS COCO) in the year 2015, the proposed method secured 3rd place.

*HyperNet.* The authors discovered HyperNet<sup>[21]</sup> based on the deep hierarchical network for object detection and the treated ROI Scheme. The framework of HyperNet is grounded on hyper features and worked on the computational mechanism (see Algorithm 4). This framework is a collection of hierarchical feature maps and wrapped them into consistent space. The suggested method has various hyper features for object detection such as improved semantic, balancing and superficial features of the image. Hyper-Feature used joined training stratagem. The proposed method accurateness on PASCAL VOC 2007/2012 expending first 100 proposals each image, speed of 5 fps on a GPU and best for actual dealing out.

Algorithm 4: HyperNet training process. After 6 steps, the proposal and detection modules from a unified network.

Step 1: Pre-train a deep CNN model for initializing basic layers in Step 2 and Step 3.

Step 2: Train HyperNet for region proposal generation.

Step 3: Train HyperNet for object detection using region proposals obtained from Step 2.

Step 4: Fine-tune HyperNet for region proposal generation sharing Hyper Feature layers trained in Step 3.

Step 5: Fine-tune HyperNet for object detection using region proposals obtained from Step 4, with shared Hyper Feature layers fixed.

Step 6: Output the unified HyperNet jointly trained in Step 4 & Step 5 as the final model.

*OHEM.* The region centered ConvNet detections, and the authors, developed a novel approach known as the online hard example mining (OHEM) method<sup>[22]</sup>. The OHEM, excludes numerous heuristics and hyper-parameters in mutual usage. The authors tested the proposed method on easy as well as hard examples. It harvests reliable and noteworthy improvements in detection enactment on benchmarks of PASCAL VOC 2007/2012 as of 78.9% and 76.3% mAP. On the MS COCO dataset, whereas datasets are massive, the proposed method achieved superior performance.

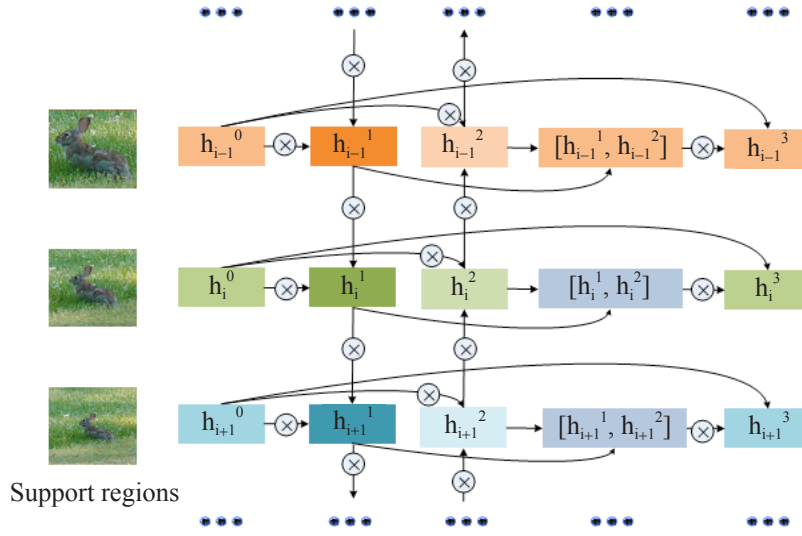
*CRAFT.* The authors<sup>[23]</sup> investigated Cascade-Region-proposal-network and FasT-rcnn (CRAFT), which handles every segment by a wisely deliberate network cascade. The authors used the “divide and conquer” method and further divide every job into two sub-jobs. It decreases the false positive via apprehending composed to inter and intra-category alterations. As compared to the 03 benchmark standard, namely PASCAL VOC 07/12 and ILSVRC, the proposed method attains reliable and substantial enhancement over the state-of-the-art on object-detection.

*MPN.* There are 03 alterations to Fast R-CNN given by the authors, these are, firstly frisk influences to the elasticity of the network entrée to multi-scale features. Secondly, for delivering context, foveal regions are used. Thirdly, essential damage to expand localization. The authors tied the proposed outcome with MultiPath classifier and DeepMask proposals and attained a 66% enhancement over the reference point Fast R-CNN through Selective Search<sup>[24]</sup>.

*GBDNet.* Gated bi-directional CNN (GBD-Net) presented by the authors<sup>[25]</sup>. GBD-Net delivered the feature information from it sustenance regions for the period of learning and extraction. That mechanism has done via the convolution among neighbouring sustenance regions in 02 ways and can be directed in numerous layers. The mechanism not always helpful but reliant on separable samples, so for controlling the mechanism needed extra-Gated function (GBD-v1 and GBD-v2 illustrated below). Extensive experiments applied on 03 datasets, i.e., ImageNet, Pascal VOC 2007 and Microsoft COCO and showed the effectiveness of the method. The method won the ImageNet object detection challenge in the year 2016.

#### (a). Structure of GBD-v1

Bi-direction arrangement: Figure 2 shows the framework of the bi-directional net. It takes features  $f^{-0.2}$ ,  $f^{0.2}$ ,  $f^{0.8}$  and  $f^{1.7}$  as input and produces features  $h_1^3$ ,  $h_2^3$ ,  $h_3^3$  and  $h_4^3$  for a single candidate box.



**Figure 2. Bi-directional arrangement:**  $\otimes$  represents the convolution. The input of this structure is the features  $\{h_j^0\}$  of contextual regions. The bi-directional links between these features are used for sending the information to the contexts. The outcome  $h_i^3$  modify the features for diverse resolutions

The features  $\{h_i^3\}$  construct two-directional links. The first link begins from the features and have the smallest region size and terminate at the feature. While the second one works oppositely. For a single candidate box  $b^0$ ,  $h_i^0 = f^{p_i}$  demonstrate features with context-pad value  $p_i$ . The forward propagation for the proposed bidirectional scheme can be described as follows

$$b^0 = [x^0, y^0, w^0, h^0] \quad (1)$$

(where center location  $(x^0, y^0)$ , width  $w^0$  and height  $h^0$ )

$$h_i^1 = \sigma(h_i^0 \otimes w_i^1 + b_i^{0,1}) + \sigma(h_{i-1}^1 \otimes w_{i-1,i}^1 + b_i^1), \quad (2)$$

(high res. To low pass)

$$h_i^2 = \sigma(h_i^0 \otimes w_i^2 + b_i^{0,2}) + \sigma(h_{i-1}^2 \otimes w_{i-1,i}^2 + b_i^2), \quad (3)$$

(low res. To high pass)

$$h_i^3 = \sigma(\text{cat}(h_i^1, h_i^2) \otimes w_i^3 + b_i^3), \quad (4)$$

(message integration)

There are 04 various resolutions,  $i = 1, 2, 3, 4$ .

$h_i^1$  shows the modified features after obtaining information from  $h_{i-1}^1$  with excellent resolution and the lowest support-region. Suppose that  $h_0^1 = 0$ , while  $h_i^1$  has the lowest support-region and obtains no information.

$h_i^2$  shows the modified features after obtaining information from  $h_{i+1}^2$  with a smaller resolution and a broader support-region. Suppose that  $h_5^2 = 0$ , while  $h_4^2$  has the broadest support region and obtains no information, and  $\text{cat}()$ -concatenates CNN features maps along the channel direction. The features  $h_i^1$  and  $h_i^2$  after information passing are combined into  $h_i^3$  by the convolutional filters  $w_i^3$ .  $\otimes$  Show the convolution operation. While the biases and filters of convolutional layers are represented by  $b^*$  &  $w^*$ . For non-linear function  $\sigma(\cdot)$  RELU is used. From the above-given equations the features in  $h_i^1$  obtain the information from the lower-context features. The features  $h_i^2$  obtain information from the higher-context features. Then  $h_i^3$  gathers information via 02 directions to have a great demonstration of the  $i$ th resolution.

### (b). Structure of GBD-v2

The amended GBD-Net structure has the following formulation.

$$h_i^1 = \sigma(h_i^0 \otimes w_i^1 + b_i^{0,1}) + G_i^1 \cdot \sigma(h_{i-1}^1 \otimes w_{i-1,i}^1 + b_i^1), \quad (5)$$

$$h_i^2 = \sigma(h_i^0 \otimes w_i^2 + b_i^{0,2}) + G_i^2 \cdot \sigma(h_{i+1}^2 \otimes w_{i,i+1}^2 + b_i^2), \quad (6)$$

$$h_i^{3,m} = \max(h_i^1, h_i^2), \quad (7)$$

$$h_i^3 = h_i^0 + \beta h_i^{3,m}, \quad (8)$$

Gate functions work as convolution layers with the sigmoid non-linearity to make the information passing-rate in the range of (0, 1).

$$h_i^1 = \sigma(h_i^0 \otimes w_i^1 + b_i^{0,1}) + G_i^1 \cdot \sigma(h_{i-1}^1 \otimes w_{i-1,i}^1 + b_i^1), \quad (9)$$

$$h_i^2 = \sigma(h_i^0 \otimes w_i^2 + b_i^{0,2}) + G_i^2 \cdot \sigma(h_{i+1}^2 \otimes w_{i,i+1}^2 + b_i^2), \quad (10)$$

$$G_i^1 = \text{sigm}(h_{i-1}^0 \otimes w_{i-1,i}^g + b_{i-1,i}^g) \quad (11)$$

$$G_i^2 = \text{sigm}(h_{i+1}^0 \otimes w_{i+1,i}^g + b_{i+1,i}^g) \quad (12)$$

Where  $\text{sigm}(x) = 1 / [1 + \exp(x)]$  is the element-wise sigmoid function and  $(\cdot)$  indicates element-wise product.  $G$  represents the gate function to control information passing. It has learnable convolutional parameters  $W_*^g, b$  and uses features from the colocated regions to decide the rates of information passing. When  $G(x, w, b)$  is 0, the information is not sent. The formulation for obtaining  $h_i^3$  is unaffected (Figure 3).

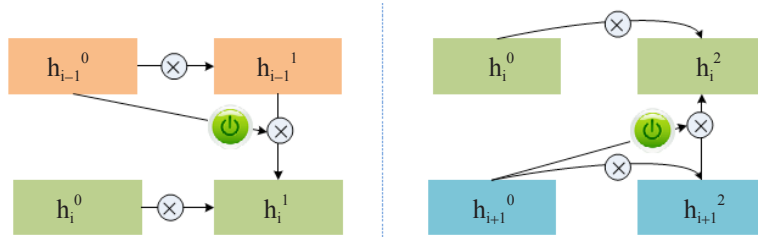


Figure 3. Bidirectional scheme with gate functions

Where  $G_i^1$  and  $G_i^2$  are defined in equation 11 & 12. Figure 4 shows the updated GBD scheme. The operations essential for gaining  $h_i^1$  and  $h_i^2$  are the same as before. The core changes for obtaining  $h_i^3$  are as follows.

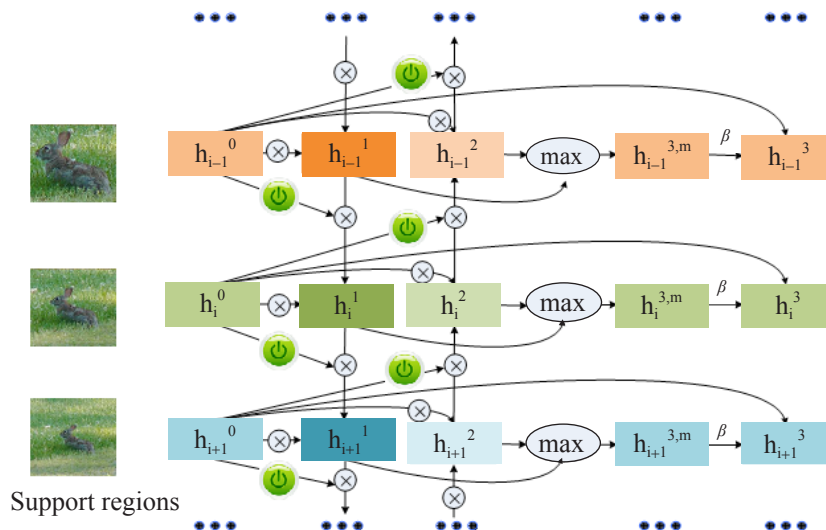


Figure 4. Updated bi-directional scheme



$h_i^1$  and  $h_i^2$  are concatenated and then convoluted by filters to produce the result  $h_i^3$ . In the updated scheme, max-pooling is used for combining the information from  $h_i^1$  and  $h_i^2$ . This saves the memory and computation required by the convolution in the old GBD structure. Then, improve an identity mapping layer in the structure, which similar to the  $h_i^3 = h_i^0 + \dots$ . The main objective of the GBD scheme is to refine the input by using the information from the other contextual features. The pre-trained model provides the output for the layer, and this led to create the learning difficulties for a good model. Careful selection or initialization of the convolution parameters & the gate function is necessary to train the model for efficient learning.

*CPF*. The authors <sup>[26]</sup> focused on 03 tasks, supplement Faster RCNN through a semantic separation network, the separation network designed for at the top contextual grooming and deliver at the top iterative response expending 02 phase training. The authors recognized the integration of top-down and responses in the modern Faster R-CNN structure. The proposed methods achieved excellent performance on object-detection, semantic-segmentation and region-proposal unit.

*MS-CNN*. Fast multi-scale object detection methods, a cohesive deep NN is known as multi-scale CNN (MS-CNN) introduced by the authors <sup>[27]</sup>. MS-CNN contains subnetwork, detection executed at multiple output layers and receptive field matched the object, and it learned the network via optimizing a multi-task loss. MS-CNN outperformed over 02 datasets KITTI and Caltech.

*R-FCN*. There are 02 aspects of the proposed methods translation-invariance in image cataloging and translation-variance in object-detection <sup>[28]</sup>. The proposed technique is region-based and fully-convolutional-networks aimed at precise and well-organized detection of the object. The extensive experimental outcome showed 83.6% mAP on PASCAL VOC 2007 thru 101 layer ResNet, and computational speed 170 ms per image it is 2.5-20× quicker than the Faster R-CNN equivalent.

*DeepID-Net*. The authors focused on the novel deformation constrained pooling (def-pool) layer facsimiles, whereas the deformable-object proportion integrated with geometric restriction and forfeit <sup>[29]</sup>. The preprocessed scheme worked on feature learning via varying the model and trained strategy. The slack and surplus are the main apparatuses in the detection cylinder. The whole structure available for investigators comprehends the deep learning object detection mechanism. The proposed method outperformed on ILSVRC 2014 dataset and achieve 31% to 50.3% accuracy. With 6.1% the proposed method became champion of ILSVRC 2014 and GoogLeNet.

*NoC*. The authors considered deep network “Networks on Convolutional feature maps” (NoCs). The proposed method can classify the object and practically apply on region-wise classifier networks that are custom pooled and region-independent convolutional features <sup>[30]</sup>. NoCs discovered the deep and convolutional per-region classifier significance for the detection of the object. The proposed method experimentally validate as compared to ResNets and Faster R-CNN schemes. NoCs admitted in the ImageNet, and MS COCO challenged 2015 and achieved the first position.

*TDM*. The authors represented the feed-forward ConvNet with a top-down modulation (TDM) network as complements of the bottom-up approach, which links through the adjacent link <sup>[31]</sup>. The top-down layer network responsible for the managed assortment and amalgamation of background information and low-level-features. The proposed framework delivered a notable boost on the COCO dataset and attained 28.6 AP for VGG-16 and 35.2 AP for ResNet-101 networks. The proposed method achieves 37.3 AP as compared to Inception ResNet v2.

*FPN*. In this investigation, authors concentrated on Feature Pyramid Network (FPN), which used an intrinsic multi-scale and pyramidal hierarchy of deep convolutional networks designed for feature pyramids thru the minimal spare budget <sup>[32]</sup>. FPN designed by basic Faster R-CNN structure. The proposed approach achieved an effective and efficient outcome on COCO detection dataset. The proposed method became titleholder in COCO 2016 challenge for single model entries. The proposed method is configured and run at 5 FPS on a GPU and thus is an experimental and precise elucidation of multi-scale object detection.

*RON*. The authors focused on novel Reverse-Connection-with-Objectness-Prior Networks (RON) for Object-detection <sup>[33]</sup>, which is fast and efficient. The authors joined the 02 approaches energetically named as region-based and region-free. For fully convolution framework of RON, there are 02 main difficulties-first, multi-scale localization and second, negative sample searching. The solution to the 02 main difficulties is obtained by construction of reverse-connection and objectness-prior. The authors used 03 standard datasets for validation of the proposed method with extensive experiments, i.e. PASCAL VOC 2007, PASCALVOC 2012 and MSCOCO. With  $384 \times 384$  pixel size image and VGG16, the PASCALVOC 2007 achieved 81.3% in mAP, and PASCALVOC 2012 achieved 80.7% in mAP. When applied with the huge dataset and complexity, the proposed method gets superior performance on the MS COCO dataset.

*RSA*. The authors suggested the novel recurrent scale approximation (RSA) unit <sup>[34]</sup>. This method validates the

deep-CNN feature via a major level to minor level. RSA is superior for face detection and attains identical results in the detection of an object. RSA is easy and fast for detection. The authors also construct the innovative retracing network to emphasize global & local scale data to improve forecasting. The extensive experiment demonstrates the effectiveness of the proposed technique and summed up the proposed method to beat state-of-the-art approaches.

*DCN*. For improved transformation modelling competence of CNNs, the authors presented 02 novel segments, i.e. deformable-convolution-networks and deformable RoIpool [35]. All segments are grounded on the notion of augmenting spatial specimen positions via extra-offsets (the detailed framework is given below). The proposed method substitutes the previous one and trained by back-propagation. The proposed method outperformed other state-of-the-art techniques.

In *DCN*, Feature maps and convolution are considered as 3D. The 02 sections-deformable convolutions and RoI-pooling control the 2D spatial-domain. The operation remains the same across the channel dimension. For simplicity, the section is described in 2D.

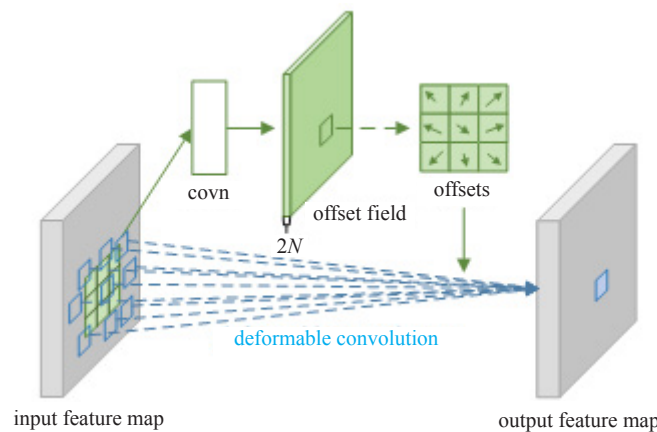


Figure 5.  $3 \times 3$  deformable convolution [35]

### (a). Working of deformable convolution

The 2D convolution includes 02 phases: First, sampling uses a regular-grid  $\mathfrak{R}$  over the incoming feature-map  $x$ ; Second, the summation of all specimen values weighted via  $w$ . The grid  $\mathfrak{R}$  describes the receptive-field dimension and expansion. For instance,

$$R = \{(-1, -1), (-1, 0), \dots, (0, 1), (1, 1)\}$$

Describes a  $3 \times 3$  kernel with expansion one. For each location  $p_0$  on the outcome feature map  $y$ ,

$$y(p_0) = \sum_{p_n \in R} w(p_n) \cdot x(p_0 + p_n) \quad (13)$$

Where  $p_n$  count the position in  $R$ .

In deformable-convolution, the regular grid  $\mathfrak{R}$  is improved with offsets

$$\{\Delta p_n | n = 1, \dots, N\}, \text{ where } N = |R|$$

Eq. (13) becomes

$$y(p_0) = \sum_{p_n \in R} w(p_n) \cdot x(p_0 + p_n + \Delta p_n) \quad (14)$$

Now, the sampling is on the uneven and offset position  $p_n + \Delta p_n$ . As the offset  $\Delta p_n$  is usually fractional. Eq. (14) is executed via bilinear interpolation as

$$x(p) = \sum_q G(q, p).x(q), \quad (15)$$

Where  $p$  denotes a random (fractional) position ( $p = p_0 + p_n + \Delta p_n$  for Eq. (15)),  $q$  show all integral-spatial position in the feature map  $x$ , and  $G(\cdot, \cdot)$  is the bilinear interpolation kernel. Note that  $G$  is 02-dimensional. It divides 01-dimensional kernels as

$$G = (q, p) = g(q_x, p_x).g(q_y, p_y), \quad (16)$$

Where  $g(a, b) = \max(0, 1 - |a - b|)$ . Eq. (15) is efficient to compute as  $G(q, p)$  is non-zero only.

As shown in Figure 6, the offsets are achieved by applying a convolutional-layer over the identical input-feature map.

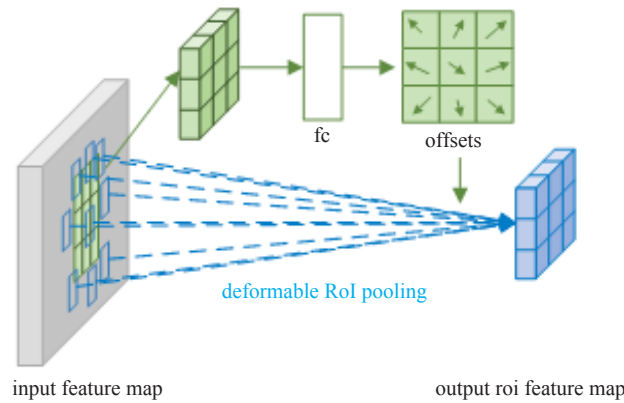


Figure 6. 3 × 3 deformable RoI pooling <sup>[35]</sup>

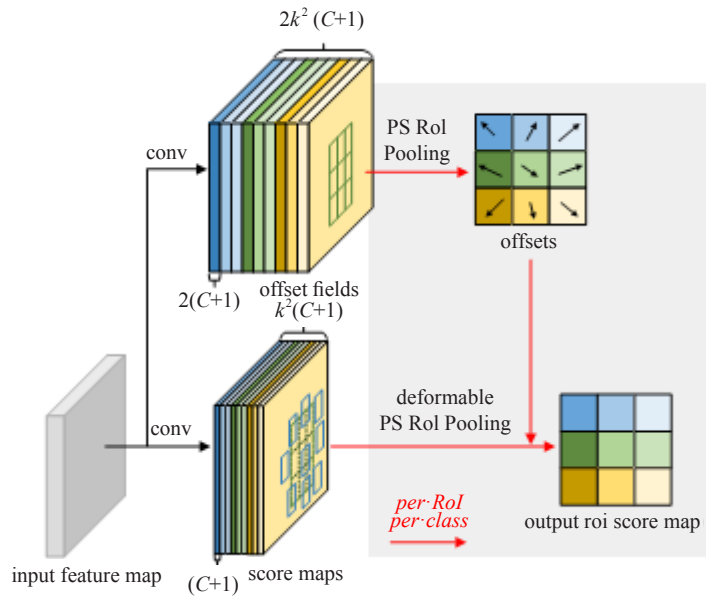


Figure 7. 3 × 3 deformable PS RoI pooling <sup>[35]</sup>

The convolution-kernel of the identical spatial-resolution & expansion is the same as the current convolutional-layer (e.g., 3 × 3 with expansion one in Figure 5). The outcome offset areas have the exact spatial-resolution with the incoming feature map. The channel dimension is similar to the 2D-offsets. Throughout the training, the convolutional-kernels produce the outcome features and offsets are learned together. For offsets learning, the gradients are back-propagated through Eq. (15) and Eq. (16).

### (b). Working of deformable RoI pooling

In the region proposal, the object detection method RoI pooling is used. It changes an input rectangular region of random size into fixed-size features.

RoI Pooling, given the input feature map  $x$  and a RoI of size  $w \times h$  and top-left corner  $p_0$ , RoI pooling splits the RoI into  $k \times k$  ( $k$  is a free-parameter) bins and outcome a  $k \times k$  feature-map  $y$ . For  $(i, j)$ -th bin ( $0 \leq i, j < k$ ),

$$y(i, j) = \sum_{p \in \text{bin}(i, j)} x(p_0 + p) / n_{ij}, \quad (17)$$

Where  $n_{ij}$  is the no. of pixels in the bin. The  $(i, j)$ -th bin spans  $\lfloor i \frac{w}{k} \rfloor \leq p_x < \lceil (i+1) \frac{w}{k} \rceil$  and  $\lfloor j \frac{h}{k} \rfloor \leq p_y < \lceil (j+1) \frac{h}{k} \rceil$ .

In the same way, as in Eq. (13), in deformable-RoI-pooling, offsets  $\{\Delta p_{ij} | 0 \leq i, j < k\}$  are joined to the spatial binning locations. Eq.(16) becomes

$$y(i, j) = \sum_{p \in \text{bin}(i, j)} x(p_0 + p + \Delta p_{ij}) / n_{ij} \quad (18)$$

Typically,  $\Delta p_{ij}$  is fractional. Eq. (18) is implemented by bilinear interpolation via Eq. (15) and (16). Figure 7 shows how to get offsets. Firstly, RoI pooling (Eq. (17)) produces the pooled-feature-maps. From the maps, a fc layer creates the normalized offsets, which are then transformed to the offsets  $\Delta p_{ij}$  in Eq. (18) thru element-wise product with the RoI's width & height, as  $\Delta p_{ij} = \gamma \cdot \hat{p}_{ij} \circ (w, h)$ . Here  $\gamma$  is a pre-defined scalar to control the magnitude of the offsets. By default, it set to  $\gamma = 0.1$ . The offset normalization is essential to make the offset learning invariant to RoI-size.

### (c). Deformable ConvNets for enhanced object detection

Both deformable convolutions and RoI-pooling sections have identical input and outcome as their basic-versions. Hence, they can readily replace their basic counterparts in current CNNs. In training, these integrate conv and fc layers for offset learning. Their learning-rates are fixed to  $\beta$  times ( $\beta = 1$  by default, and  $\beta = 0.01$  for the fc layer in Faster R-CNN) of the learning-rate for the current layers. They trained via back-propagation through the bilinear-interpolation operations in Eq. (15) and Eq. (16). The resulting CNNs are called deformable-ConvNets. for the incorporation of deformable ConvNets with the state-of-the-art CNN framework, this framework contains 02-stages. First, a deep-fully-convolutional net produces feature maps over the entire input image. Second, a shallow task-specific net produces an outcome from the feature-maps.

*DeNet*. The authors described the image detection of huge sparse bounding box distribution. The authors also classified the sparse distribution assessment structure, sparse specimen and CNN detection approach. For this purpose, the authors studied and analyzed the prevalent method as evaluation speed to decrease, the labour-intensive worked. They presented 02 uniquenesses, corner based region-of-interest estimator + deconvolution based CNN prototypical<sup>[36]</sup>. The proposed model automates and worked on an actual environment with 03 datasets named as MSCOCO, PascalVOC 2007 and PascalVOC 2012. The wide-ranging experiments demonstrated the efficiency of the proposed method.

*CoupleNet*. Connect the whole construction via an accessible feature segment for detection of an object, and the authors proposed a new fully convolutional network known as CoupleNet<sup>[37]</sup>. Region-Proposal-Network (RPN) give the object proposal and tightly add the connected unit which contains 02 undergrowths, first one extracts the feature segment info of the object from position sensitive RoI (PSRoI) pooling whereas the other one encrypts the whole context info from RoI pooling. The authors integrate these 02 aspects into one framework. The wide-ranging experiments validate the effectiveness of the proposed method as compared on 03 datasets, i.e. PASCAL VOC 2007 achieved 82.7% in mAP, PASCALVOC 2012 achieved 80.4% in mAP and COCO achieved 34.4% in mAP.

*RetinaNet*. The problem encountered by a class imbalance in foreground-background during the training of dense object. So for this purpose, authors suggested class imbalance thru restructuring cross-entropy loss. During the training, the focal loss centers around the sparse routine and also handle the negatives detector. To consider the efficacy of loss, authors construct an easy dense object detector known as RetinaNet<sup>[38]</sup>. RetinaNet outperformed as a prevalent method.

*Mask R-CNN*. In this paper, object instance segmentation structure proposed by the authors and known as Mask R-CNN<sup>[39]</sup>. This method is practically easy, efficient and able to detect objects and produce a segmentation mask for each instance. Mask R-CNN successor of Faster R-CNN for forecasting. Mask R-CNN is a winner of COCO 2016 challenge and showed the superior outcome in COCO set of contests with instance segmentation, bounding box and key point detection. The proposed method beat the other prevalent method given in the literature.

*SMN*. The authors<sup>[40]</sup> proposed Spatial Memory Network (SMN) framework. This method is theoretically modest and influential. The spatial memory accumulates the object instances posteriorly hooked on a self-styled "image" depiction.

The SMN network outperformed, and it delivered a 2.2% enhancement over Faster RCNN on the COCO dataset.

*Light-Head R-CNN.* The authors proposed 02 phase detector methods known as Light-Head RCNN <sup>[41]</sup>. For the construction of Light-Head RCNN, the method needs a primary network which is easily adjustable thru highly feature map & pooling and single fully connected layer (R-CNN subnet). The proposed ResNet-101 grounded on light-head R-CNN overtakes the other method. The extensive experiment showed LightHead R-CNN acquired 30.7 mAP at 102 FPS on COCO and beat the YOLO and SSD on rapidity and accurateness by single-stage and fast detector.

*Soft-NMS.* The novel method Soft-NMS introduced by the authors <sup>[42]</sup>, Soft-NMS is very easy to implement and no need for further training. This method also accumulates any detection framework. Soft-NMS method separates the detection rate as a continuous function, and never remove the object. The proposed method experiment validates and gets continuous enhancements for the coco style mAP metric-datasets, i.e. PASCALVOC 2007 and MS-COCO. The wide-ranging investigations applied on PASCALVOC 2007, and MS-COCO outperformed as compared 1.7% for both RFCN and Faster-RCNN & 1.3% for R-FCN and 1.1% for Faster-RCNN. The proposed method improves the state-of-the-art object-detection from 39.8%-40.9% via one model.

*ZIP.* The new network zoom-out-and-in constructed with superior features with superior pixelization & combined with low features to detect objects <sup>[43]</sup>. The special construction is based on a different scale, size and location. The feature representation of different stages is used to detect the object of small, large and medium scale. The authors proposed the iterative framework to sequentially degenerate the region proposals at the time of training the phase to check whether iterative regression at the testing phase is the same or not. The wide-ranging experiments demonstrated the effectiveness of the proposed method on 02 datasets, named as ILSVRC DET and MS COCO, whereas the proposed method outperforms than the state-of-the-art methods in numerous evaluation metrics. The proposed technique enhances the average accuracy by 2% in the object-detection system.

*SIN.* The proposed methods worked on 02 tasks, i.e. contextual scene info and single image object extraction, and define that solved this with the help of cognition and reasoning. Edges and nodes represent the problem framed by graph structure inference. The authors presented Structure Inference Network (SIN) <sup>[44]</sup>. The graphical model integrates into a typical object detection scheme. The extensive practical on 02 datasets named as PASCAL VOC and MS COCO designated that the 02 tasks improved the effectiveness of the detection of an object by precise and fast outcomes.

*STDN.* The authors focused on multi-scale object detection by Scale-Transferrable Detection Network (STDN) <sup>[45]</sup>. The proposed network integrates with the scale transfer layer, which discovers the inter-scale steadiness manifold environment. Scale-transfer component combined with leading base network smartly, and dense-convolutional-network (DenseNet) harvest a one-phase object-detector. Comprehensive experiments on 02 datasets named PASCALVOC 2007 and MSCOCO is applied by the proposed technique and demonstrated the effectiveness over the equivalent state-of-the-art detection techniques.

*RefineDet.* The authors suggested a novel RefineDet <sup>[46]</sup> single-shot detection approach, which is fast, effective and efficient. The recommended approach is superior in the accurateness as comparable 02 stage approaches. RefineDet, arranged with 02 internal parts, i.e. anchor-refinement part and the object-detection part. The suggested method works on the concept of noise removal from the anchor and decreases the computational complexity for the classifier. This method also regulates the anchor to delivered superior representation for the subsequent regressor. The authors also construct the new segment, i.e. a transfer construction, transfer the feature map in the anchor-refinement part to forecast the objects. The experiments on 03 datasets, i.e. PASCALVOC 2007, PASCALVOC 2012, and MSCOCO validates the proposed method to achieve state-of-the-art detection accurateness with superior efficiency.

*MegDet.* The authors, considered a fast effective and efficient method Large Mini-Batch Object Detector (MegDet) <sup>[47]</sup>. MegDet excellently configured on 128 GPUs to reduce the training the period. The authors proposed a warmup learning rate strategy and Cross-GPU-Batch-Normalization that integrate and permit a mini-batch detector a shorten period and attain superior accuracy (see Algorithm 5). MegDet achieved first place in the competition of object detection tasks in COCO 2017. MegDet earned 52.5% nmAP in COCO 2017.

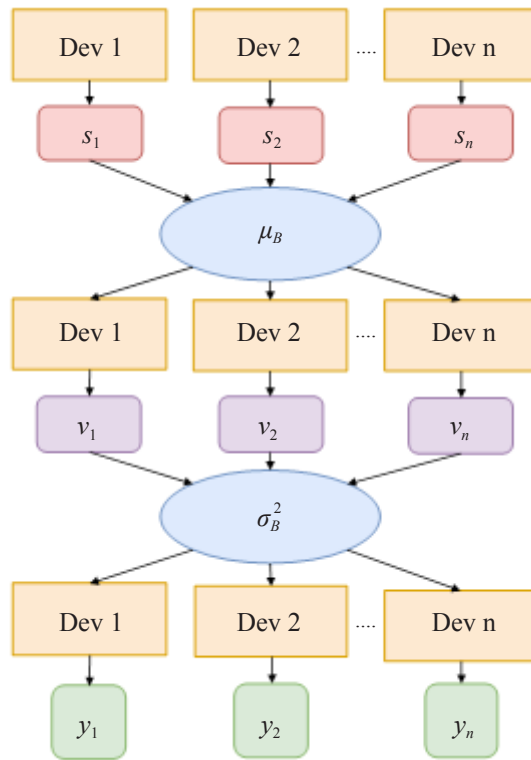


Figure 8. Execution of Cross-GPU batch normalization

The execution of Cross-GPU Batch-Normalization is sketched in Figure 8. Given  $n$  GPU devices in total, the sum value  $s_k$  is first calculated based on the training examples assigned to the device  $k$ . By averaging the sum-values from whole devices, the mean value  $\mu_\beta$  is obtained for current mini batch. This step needs an AllReduce operation. The  $\sigma_\beta^2$  value is obtained for the calculated variance. After distribution  $\sigma_\beta^2$  to each device, the standard normalization is obtained by  $y = \gamma \frac{x - \mu_\beta}{\sqrt{\sigma_\beta^2 + \epsilon}} + \beta$ . Algorithm 5 gives the complete flow. In the implementation, the NVIDIA-Collective-Communication-Library (NCCL) efficiently perform AllReduce operation for achieving and broadcasting.

Algorithm 5: Cross-GPU Batch Normalization over a mini batch  $\beta$ .

Input: Values of input  $x$  on multiple devices

In a minibatch:  $\beta = \cup_{i=1}^n \beta_i, \beta_i = \{x_{i1} \dots i_n\}$

BN parameters:  $\gamma, \beta$

Output:  $y = CGBN(x)$

1. For  $i = 1, \dots, n$  do
2.     Compute the device sum  $s_i$  overset  $B_i$
3. End for
4. Reduce the set  $s_1, \dots, n$  to *mi nibatchmean*  $\mu_\beta$
5. Broadcast  $\mu_\beta$  to each device
6. For  $i = 1, \dots, n$  do
7.     Compute the device variance sum  $v_i$  overset  $B_i$
8. End for
9. Reduce the set  $s_1, \dots, n$  to *mi nibatchvarianve*  $\sigma_\beta^2$
10. Broadcast  $\sigma_\beta^2$  to each device
11. Compute the output:  $y = \gamma \frac{x - \mu_\beta}{\sqrt{\sigma_\beta^2 + \epsilon}} + \beta$  over devices

*DA Faster R-CNN*. This method handles the domain shift in 02 ways, i.e. image level and instance level shift. Image level is used for illumination, elegance, etc. whereas the instance level used for the presence of an object, scale, etc. The authors enhanced the cross-domain rapidity of object detection. The authors proposed a Faster R-CNN scheme that works on image and instance level, and both are integrated with the divergence theory and trained by domain classifier learner<sup>[48]</sup>. This method learns province invariant region-Proposal-network (RPN) thru FASTER R-CNN. The extensive experiment

on various datasets such as Cityscapes, KITTI, SIM10K, etc. validates the success of the proposed method.

### 2.2.2 Regression or classification based approach

Region proposal based approach is time-consuming and has no real-time applications, and these drawbacks are eliminated by the regression or classification-based approaches.

*MultiBox*. In this paper, the authors <sup>[49]</sup> introduced the saliency enthused neural network scheme for object-detection, which detects the region-of-interest from the image via class-agnostic bounding-boxes and a particular score for the respective box. The technique automates the number of variables for each group and permits for cross-class over-simplification at the uppermost levels of the network. The proposed method outperforms on 02 datasets, namely VOC 2007 and ILSVRC 2012, although considering only topmost locations in an image.

*AttentionNet*. The novel method for object detection given by the authors known as AttentionNet <sup>[50]</sup>, which built with CNN. The authors performed iterative classification problem as object detection, and CNN solved that type of classification. AttentionNet forecast the exact boundary box of an image because of its nature of a unified network. There are 02 datasets used for experimental evaluation PASCALVOC 2007 and PASCALVOC 2012 and achieved 65% (AP) on human detection scheme by eight layered frameworks.

*YOLO v1*. The authors <sup>[51]</sup> introduced YOLO for object detection. YOLO used detached bounding boxes and accompanying class likelihoods. Each NN forecast bounding box and class likelihoods openly in one cycle. YOLO optimized endwise framework arrangements, so the entire detection process under the one network. YOLO method works fast, and the progression of the image is 45 fps in actual. As an experimental validation YOLO, less forecast the false positives on contextual but with extra localization errors. YOLO trained for a simple demonstration of an object. YOLO overtake as compared to the DPM and R-CNN.

*G-CNN*. G-CNN introduced by the authors <sup>[52]</sup>, and this method is based on CNN without proposal procedures. The multi-scale network of fixed bounding boxes is a key point of G-CNN. G-CNN strongly bind the neighbouring objects as an immovable network. The proposed method used 180 squares in a multi-scale network and achieved Fast R-CNN. G-CNN detect more rapidly and decreases the number of squares to be treated and eliminate the object proposal phase (see Algorithm 6 & 7).

#### Algorithm 6: G-CNN Training Algorithm

1. Procedure TRAIN GCNN
2. For  $1 \leq c \leq S_{train}$  do
3.      $TrainTuples \leftarrow \{ \}$
4.     For  $1 \leq s \leq c$  do
5.         If  $s = 1$  then
6.              $B^1 \leftarrow$  Spatial pyramid grid of boxes
7.              $G^* \leftarrow A(B^1)$
8.         Else
9.              $B^s \leftarrow T^{s-1}$
10.         End if
11.          $T^s \leftarrow \mathcal{O}(B^s, G^*, s)$
12.          $\Delta^s \leftarrow \Delta(B^s, T^s)$
13.         Add  $(B^s, \Delta^s)$  to TrainTuples
14.     End for
15.     Train G-CNN  $N_{iter}$  iterations with TrainTuples
16.     End for
17. End procedure

#### Algorithm 7: Detection algorithm

1. Let  $f(\cdot)$  be the feed-forward G-CNN regression network
2. Let  $c(\cdot)$  be the classifier function
3. Procedure DETECT
4.      $B^1 \leftarrow$  Spatial pyramid grid of boxes
5.     For  $1 \leq s \leq S_{test}$  do
6.          $l \leftarrow c(B^s)$
7.          $\delta_l^s \leftarrow f(B^s)$
8.          $B^{s+1} \leftarrow B^s + \Delta^{-1}(\delta_l^s)$

9. End for
10. Output  $B^{S_{test}+1}$
11. End procedure

*SSD*. In this investigation, the authors used a single deep NN, namely SSD<sup>[53]</sup>, which discretizes the resultant interstellar of bounding-boxes. In this structure, the network produces the score of each object and yields correction in the forecast period. It handles the multiple features thru different pixel values that vary in size. The proposed method is easy to implement that needs object proposals. So it was cool in the training period and easily accommodated into the environment that necessity a detection part. The experiment performed on 03 datasets, namely PASCALVOC, COCO, and ILSVRC. The extensive investigation proved that the proposed approach is faster and accurate as compared to Faster R-CNN. Designed for  $300 \times 300$  resolutions, SSD attains 74.3% mAP on VOC 2007 assessment at 59 FPS on an Nvidia-TitanX and designed for  $512 \times 512$  resolutions, SSD attains 76.9% mAP.

*DSSD*. The authors<sup>[54]</sup> proposed a DSSD model, designed for accumulated context. The proposed method constructs the object detection model and validates its efficiency on a standard benchmark. The performance is increased via the feature arrangement of encoder and decoder. DSSD outperformed as compared to prevalent SSD method in the literature. The proposed prototypically quietly attains state-of-the-art detection outcomes on PASCALVOC and COCO.

*YOLO v2*. The authors introduced the novel method YOLO9000 for an actual environment that can detect 9000+ categories of objects<sup>[55]</sup>. The authors also showed the version of YOLO, i.e. YOLO v2 efficient multi-scale model and is able to execute on different sizes and demonstrated effectively. Extensive experiments demonstrated that the effectiveness of the YOLO v2, on PASCAL VOC 2007 achieved a 76.8 % mAP at 67FPS and PASCAL VOC 2007 achieved 78.6 % mAP at 40 FPS. YOLOv2 outperformed as compared to Faster RCNN with ResNet and SSD with the specific condition. Lastly, the authors proposed the method to train object detection and classification, trained via 02 datasets named as COCO detection and the ImageNet classification dataset. YOLO9000 can detect the unlabeled object. Extensive experiments demonstrated the effectiveness of the YOLO9000 with ImageNet detection scheme and achieved 19.7% mAP for 44 out of 200 classes and worked on the actual environment.

*DSOD*. The authors considered Deeply Supervised Object Detector (DSOD) scheme to detect the object from scratch<sup>[56]</sup>. Deep supervision is permitted through dense layer knowledge which is essential in learning. The outcome of the DSOD is a single-shot detection structure. DSOD tackle the two main problems, i.e. loss function and category distribution among detection & classification job. The extensive experiments on 03 datasets named PASCALVOC 2007, 2012 and MSCOCO demonstrated the effectiveness of the proposed method. The proposed method overtakes the other prevalent methods given in the literature.

*YOLO v3*. The authors<sup>[57]</sup> introduced the upgraded version of YOLO with some minor modifications, named as YOLOv3. The introduced version is fast, effective and more efficient, and train the network smartly. YOLO v3 speed 22 ms at 28.2 mAP in  $320 \times 320$  resolution image. YOLO v3 outperform than SSD and  $03 \times$  faster. The authors tested YOLO v3 for object detection and obtained .5IOU mAP object detection metric. The proposed method obtains 57.9 AP50 in 51 ms configured on a TitanX as compared to 57.5 AP50 in 198 ms thru RetinaNet, and the received performance is  $3.8 \times$  quicker. Several standard datasets are available that play an essential role in deep learning-based object detection methods, such as VOC 2007<sup>[58]</sup>, VOC 2008<sup>[59]</sup>, VOC 2009<sup>[59]</sup>, VOC 2010<sup>[59]</sup>, VOC 2011<sup>[59]</sup>, VOC 2012<sup>[59]</sup>, ILSVRC13<sup>[60]</sup>, ILSVRC14<sup>[61]</sup>, ILSVRC15<sup>[61]</sup>, ILSVRC16<sup>[61]</sup>, ILSVRC17<sup>[61]</sup>, MSCOCO15<sup>[62]</sup>, MSCOCO16<sup>[62]</sup>, MSCOCO17<sup>[62]</sup>, MSCOCO18<sup>[62]</sup> and OID18<sup>[63]</sup>, which contains a large number of images and videos.

### 2.3 Other methods of object-detection

The authors proposed an innovative training arrangement namely Scale-Normalization-for-Image-Pyramids (SNIP) which selectively back-propagates the gradients of object-instances of diverse sizes as a utility of the image gauge<sup>[64]</sup>, Relation-Network operates on feature appearance and geometry, and they permit modeling with relations<sup>[65]</sup>, Cascade R-CNN Framework used for superior quality detection of an object by which this method eliminates the overfitting and quality mismatch problem during training time<sup>[66]</sup>, MLKP used for low-dimensional polynomial kernel approximation for intelligible computation on a modified multi-scale feature representation<sup>[67]</sup>, Fitness-NMS, Constructed to categorize only an adequately precise bounding box, rather than the best available one<sup>[68]</sup>, RFBNet, Find the connection of the dimension and peculiarity of receptive fields and achieved superior feature discriminability and robustness<sup>[69]</sup>, CornerNet, through this technique bounding box of an object form of key point pair is detected. The top-left and bottom-right position is used by single Convolution Neural Network<sup>[70]</sup>, PFPNet, The feature pyramid (FP) is designed by the broadening net breadth in place of the cumulative network. Thus spatial pyramid pool has some extra features to produce a pool that differs in size<sup>[71]</sup>, Pelee, Fast and effective structure construct with conventional as a substitute<sup>[72]</sup>, Hybrid Knowledge



Routed Modules integrates the cognitive path via O2 forms, i.e. implicit and explicit. Explicit phase is constructed through linguistic information & structured restriction whereas implicit phase designated via implicit restriction [73], M2Det is fast and accumulated with SSD framework for superior detection [74], R-DAD, in this method the object region divided into numerous segments-Extract the CNN feature in whole region, segments region and learn the semantic relation by use of high level semantic feature for classification and localization [75], ScratchDet, this method discover and train the detector from scratch [76], Libra R-CNN, is used for balanced training of an object detection and consist of 03 main parts-IoU-balanced sampler, sensible feature-pyramid, and sensible L1 loss for decreasing the disparity [77], Reasoning-RCNN, Delivered any detection net and the competence of adaptive reasoning over object-regions by developing diverse human commonsensical acquaintance [78], FSAF, is used for single shot object detection and integrate the feature pyramid scheme [79], AmoebaNet + NAS-FPN, learn by feature-pyramid net for object-detection and neural framework are able to search and find out the novel feature-pyramid framework in a new accessible search space casing across-scale joints [80], Cascade-RetinaNet, Used for reduction of misalignments and it contains sequential pipelining thru enhanced IoU thresholds design [81], TridentNet, this method produce scale-specific feature maps through a sequentially representational control [82], DAFS, Used dynamic feature selection scheme to select the novel pels from a feature map, and pels are designated on anchor arbitrary view and size [83], Auto-FPN, Framework is constructed for further detection and classify the foundation, and its uses O2 auto search, auto fusion is used to search multi-level features and auto head search is used for classification and bounding-box-regression [84], FCOS, in this method anchor box free and proposal free detector is used that removes the complexity during training [85], FreeAnchor, train the supernet with tight detector training roster via ImageNet pre-training [86], DetNAS, is a backbone net for object detection and this method contains extra phases for classification purpose. This method also has high-spatial-resolution maintained by deeper layer [87], NATS, Discover the framework space on the ground of prevalent net and reused weights, worked on channel level in place of path level [88], AmoebaNet + NAS-FPN + AA: -Demonstrated the influence of data augmentation on the detection and data augmentation affect the generalization performance for detection model [89], EfficientDet, Demonstrate the weighted bi-directional-feature-pyramid-network (BiFPN), also use the effective multi-scale feature fusion method is used to integrate the scaling approach for uniformity [90].

### 3. Dataset and performance

Object-detection is an essential aspect of cognitive science technology, and it allows the computer system to identify the object in the given picture or image. The object detection method also determines the position and size of an object with respect to the other object. Nowadays, state-of-the-art object detection representations are mechanically built by deep learning. Figure 9 demonstrates the prevalent object detection method and its performance statistics.

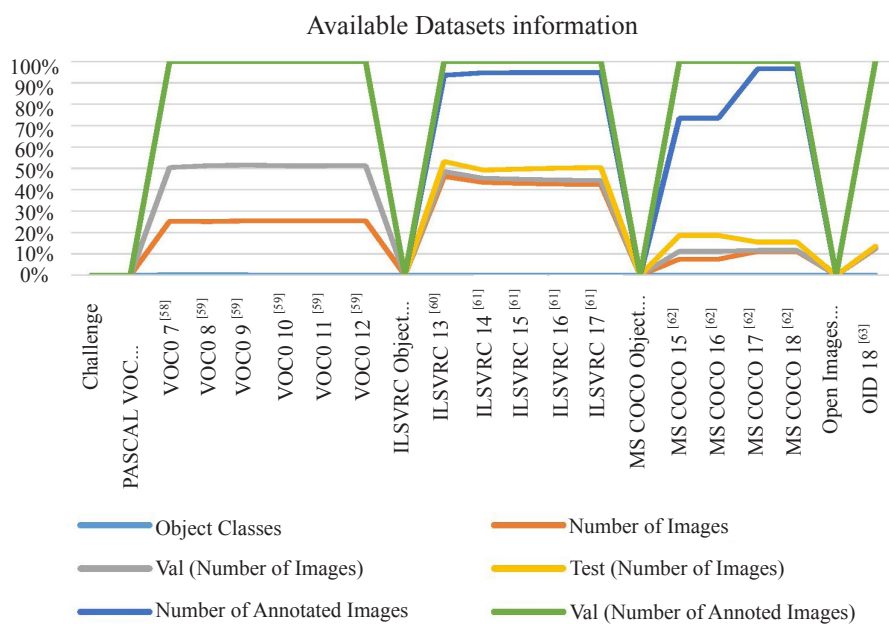


Figure 9. Available datasets for object detection in literature

As the object detection methods configured and tested on a different platform, so it is not easy to treat each method equally, based on performance, speed and metric, etc. Some method is tested on high configuration and some on the low-configurations system. This configuration depends on hardware compatibility and environment, so the performance of the method is greatly affected by the outcome. As we know, there are lots of methods given in literature so to measure effectiveness, and their accuracy is a very crucial and tedious job. In the below figure, we show the performance of the various method on 03 different datasets (see in Figure 10 and Figure 11).

Numerous Method tested on VOC07 (mAP@IoU = 0.5)

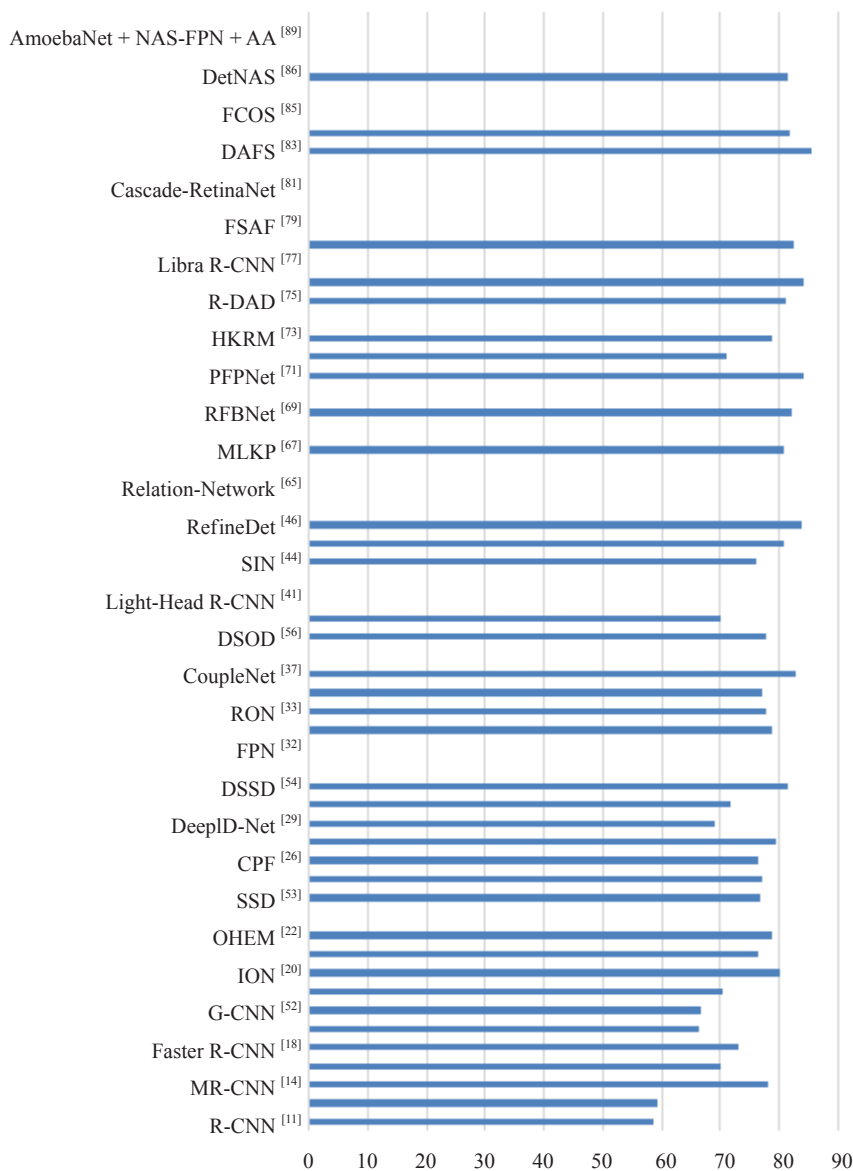


Figure 10. Numerous method tested on VOC07 (mAP@IoU = 0.5)

### Numerous Method tested on VOC12 and COCO

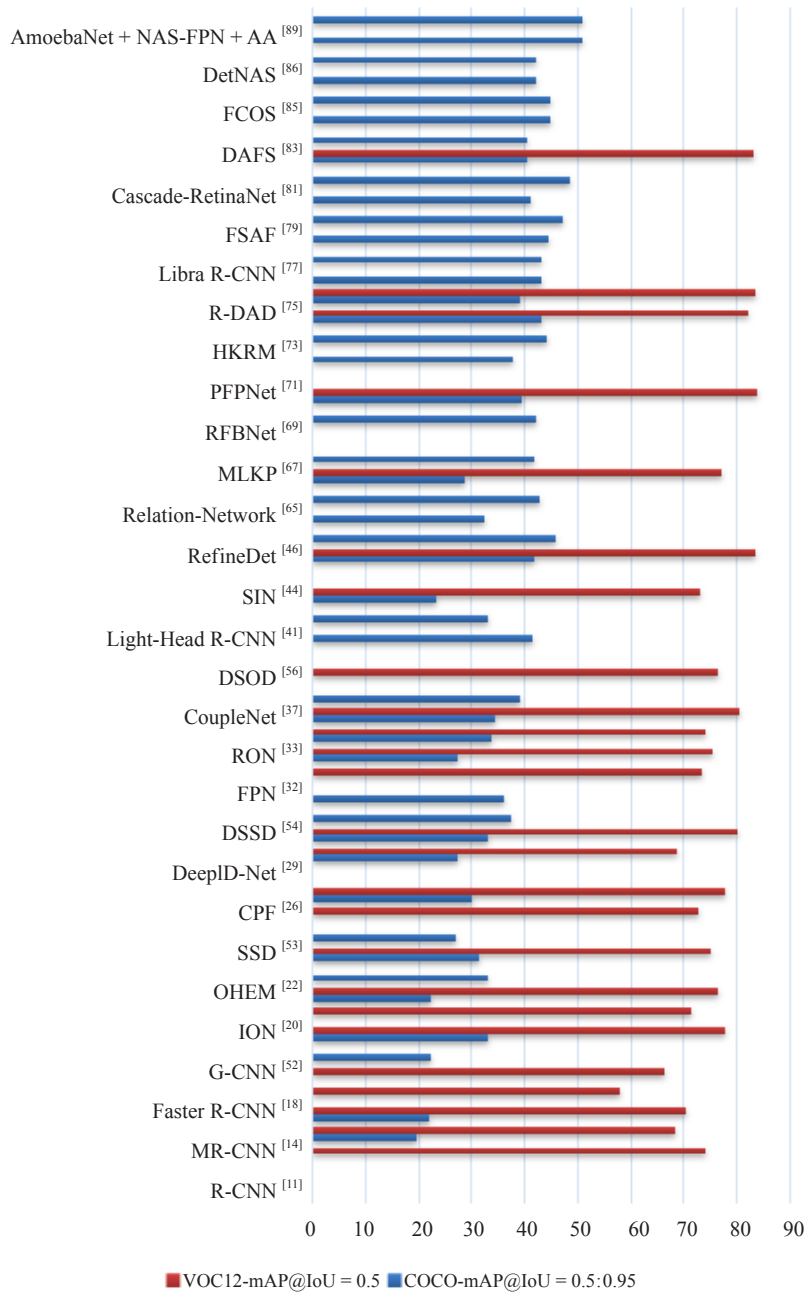


Figure 11. Numerous method tested on VOC12-mAP@IoU = 0.5 and COCO-mAP@IoU = 0.5:0.95

## 4. Conclusion

Object detection is a crucial and challenging task in computer vision and has obtained significant attention. Object detection has received great success in the last few years, but there exists a considerable gap between accuracy and speed. This research paper comprises all the major object detection methods such as CNN, R-CNN, YOLO, SIN, SNIP, Viola-Jones, HOG, etc. In this research paper, we have also discussed the modernization of the object-detection methods in recent years and described various techniques of object-detection according to their usage. Finally, we have shown the performance of all the existing methods and well-known datasets based on their evaluation criteria.

## References

---

- [1] Jiao, L., Zhang, F., Liu, F., et al. A survey of deep learning-based object detection. *IEEE Access*. 2019; 7: 128837-128868.
- [2] Liu, L., Ouyang, W., Wang, X., et al. Deep learning for generic object detection: A survey. *International Journal of Computer Vision*. 2019; 128: 261-318.
- [3] Fu, K., Zhang, T., Zhang, Y., et al. Meta-SSD: Towards fast adaptation for few-shot object detection with meta-learning. *IEEE Access*. 2019; 7: 77597-77606.
- [4] Amit, Yali, Felzenszwalb, P. Object detection. *Computer Vision*. 2014. Available from: [https://doi.org/10.1007/978-3-030-03243-2\\_660](https://doi.org/10.1007/978-3-030-03243-2_660).
- [5] Alganci, U., Soydas, M., Sertel, E. Comparative research on deep learning approaches for airplane detection from very high-resolution satellite images. *Remote. Sens*. 2020; 12: 458.
- [6] Viola, P., Jones, M. Robust real-time face detection. *International Journal of Computer Vision*. 2004; 57: 137-154.
- [7] Tony Lindeberg. Scale invariant feature transform. *Scholarpedia*. 2012; 7(5): 10491.
- [8] Kim, K., Cheon, Y., Hong, S., et al. PVANET: Deep but lightweight neural networks for real-time object detection. 2006; 1-7. Available from: [abs/1608.08021](https://arxiv.org/abs/1608.08021).
- [9] Muralidharan, R., Chandrasekar, C. Object recognition using support vector machine augmented by RST invariants. *International Journal of Computer Science*. 2011; 8(5): 280-286.
- [10] Sermanet, P., Eigen, D., Zhang, X., et al. OverFeat: Integrated recognition, localization and detection using convolutional networks. 2014; 1-16. Available from: [abs/1312.6229](https://arxiv.org/abs/1312.6229).
- [11] Girshick, R.B., Donahue, J., Darrell, T., et al. Rich feature hierarchies for accurate object detection and semantic segmentation. *IEEE Conference on Computer Vision and Pattern Recognition*. 2014. p. 580-587.
- [12] He, K., Zhang, X., Ren, S., et al. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2015; 37: 1904-1916.
- [13] Zhang, Y., Sohn, K., Villegas, R., et al. Improving object detection with deep convolutional networks via bayesian optimization and structured prediction. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2015. p. 249-258.
- [14] Gidaris, S., Komodakis, N. Object detection via a multi-region and semantic segmentation-aware CNN model. *IEEE International Conference on Computer Vision (ICCV)*. 2015. p. 1134-1142.
- [15] Kuo, W., Hariharan, B., Malik, J. DeepBox: Learning objectness with convolutional networks. *IEEE International Conference on Computer Vision (ICCV)*. 2015. p. 2479-2487.
- [16] Girshick, R.B. Fast R-CNN. *IEEE International Conference on Computer Vision (ICCV)*. 2015. p. 1440-1448.
- [17] Ghodrati, A.H., Diba, A., Pedersoli, M., et al. DeepProposal: Hunting objects by cascading deep convolutional layers. *IEEE International Conference on Computer Vision (ICCV)*. 2015. p. 2578-2586.
- [18] Ren, S., He, K., Girshick, R.B., et al. Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2015; 39: 1137-1149.
- [19] Lu, Y., Javidi, T., Lazebnik, S. Adaptive object detection using adjacency and zoom prediction. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016. p. 2351-2359.
- [20] Bell, S., Zitnick, C.L., Bala, K., et al. Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016. p. 2874-2883.
- [21] Kong, T., Yao, A., Chen, Y., et al. HyperNet: Towards accurate region proposal generation and joint object detection. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016. p. 845-853.
- [22] Shrivastava, A., Gupta, A., Girshick, R.B. Training region-based object detectors with online hard example mining. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016. p. 761-769.
- [23] Yang, B., Yan, J., Lei, Z., et al. CRAFT objects from Images. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016. p. 6043-6051.
- [24] Zagoruyko, S., Lerer, A., Lin, T., et al. A multipath network for object detection. 2016; 1-14. Available from: [abs/1604.02135](https://arxiv.org/abs/1604.02135).
- [25] Zeng, X., Ouyang, W., Yan, J., et al. Crafting GBD-Net for object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2018; 40: 2109-2123.
- [26] Shrivastava, A., Gupta, A. Contextual priming and feedback for Faster R-CNN. *European Conference on Computer Vision*. 2016. p. 330-348.
- [27] Cai, Z., Fan, Q., Feris, R.S., et al. A unified multiscale deep convolutional neural network for fast object detection. 2016; 1-16. Available from: [abs/1607.07155](https://arxiv.org/abs/1607.07155).
- [28] Dai, J., Li, Y., He, K., et al. R-FCN: Object detection via region-based fully convolutional networks. 2016; 1-11. Available from: [abs/1605.06409](https://arxiv.org/abs/1605.06409).

- [29] Ouyang, W., Wang, X., Zeng, X., et al. DeepID-Net: Deformable deep convolutional neural networks for object detection. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2015. p. 2403-2412.
- [30] Ren, S., He, K., Girshick, R.B., et al. Object detection networks on convolutional feature maps. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2017; 39: 1476-1481.
- [31] Shrivastava, A., Sukthankar, R., Malik, J., et al. Beyond skip connections: Top-down modulation for object detection. 2016; 1-11. Available from: abs/1612.06851.
- [32] Lin, T., Dollár, P., Girshick, R.B., et al. Feature pyramid networks for object detection. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2017. p. 936-944.
- [33] Kong, T., Sun, F., Yao, A., et al. RON: Reverse connection with objectness prior networks for object detection. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2017. p. 5244-5252.
- [34] Liu, Y., Li, H., Yan, J., et al. Recurrent scale approximation for object detection in CNN. *IEEE International Conference on Computer Vision (ICCV)*. 2017. p. 571-579.
- [35] Dai, J., Qi, H., Xiong, Y., et al. Deformable convolutional networks. *IEEE International Conference on Computer Vision (ICCV)*. 2017. p. 764-773.
- [36] Tychsen-Smith, L., Petersson, L. DeNet: Scalable real-time object detection with directed sparse sampling. *IEEE International Conference on Computer Vision (ICCV)*. 2017. p. 428-436.
- [37] Zhu, Y., Zhao, C., Wang, J., et al. CoupleNet: Coupling global structure with local parts for object detection. *IEEE International Conference on Computer Vision (ICCV)*. 2017. p. 4146-4154.
- [38] Lin, T., Goyal, P., Girshick, R.B., et al. Focal loss for dense object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2020; 42: 318-327.
- [39] He, K., Gkioxari, G., Dollár, P., et al. Mask R-CNN. *IEEE International Conference on Computer Vision (ICCV)*. 2017. p. 2980-2988.
- [40] Chen, X., Gupta, A. Spatial memory for context reasoning in object detection. *IEEE International Conference on Computer Vision (ICCV)*. 2017. p. 4106-4116.
- [41] Li, Z., Peng, C., Yu, G., et al. Light-head R-CNN: In defense of two-stage object detector. 2017; 1-9. Available from: abs/1711.07264.
- [42] Bodla, N., Singh, B., Chellappa, R., et al. Soft-NMS-improving object detection with one line of code. *IEEE International Conference on Computer Vision (ICCV)*. 2017. p. 5562-5570.
- [43] Li, H., Liu, Y., Ouyang, W., et al. Zoom out-and-in network with recursive training for object proposal. 2017; 1-9. Available from: abs/1702.05711.
- [44] Liu, Y., Wang, R., Shan, S., et al. Structure inference net: Object detection using scene-level context and instance-level relationships. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018. p. 6985-6994.
- [45] Zhou, P., Ni, B., Geng, C., et al. Scale-transferrable object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018. p. 528-537.
- [46] Zhang, S., Wen, L., Bian, X., et al. Single-shot refinement neural network for object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018. p. 4203-4212.
- [47] Peng, C., Xiao, T., Li, Z., et al. MegDet: A large mini-batch object detector. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018. p. 6181-6189.
- [48] Chen, Y., Li, W., Sakaridis, C., et al. Domain adaptive Faster R-CNN for object detection in the wild. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018. p. 3339-3348.
- [49] Erhan, D., Szegedy, C., Toshev, A., et al. Scalable object detection using deep neural networks. *2014 IEEE Conference on Computer Vision and Pattern Recognition*. 2014. p. 2155-2162.
- [50] Yoo, D., Park, S., Lee, J., et al. AttentionNet: Aggregating weak directions for accurate object detection. *IEEE International Conference on Computer Vision (ICCV)*. 2015. p. 2659-2667.
- [51] Redmon, J., Divvala, S.K., Girshick, R.B., et al. You only look once: Unified, real-time object detection. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016. p. 779-788.
- [52] Najibi, M., Rastegari, M., Davis, L.S. G-CNN: An iterative grid based object detector. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016. p. 2369-2377.
- [53] Liu, W., Anguelov, D., Erhan, D., et al. SSD: Single shot multibox detector. *European Conference on Computer Vision*. 2016. p. 21-37.
- [54] Fu, C., Liu, W., Ranga, A., et al. DSSD: Deconvolutional single shot detector. 2017; 1-11. Available from: abs/1701.06659.
- [55] Redmon, J., Farhadi, A. YOLO9000: Better, faster, stronger. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2017. p. 6517-6525.
- [56] Shen, Z., Liu, Z., Li, J., et al. DSOD: Learning deeply supervised object detectors from scratch. *IEEE International Conference on Computer Vision (ICCV)*. 2017. p. 1937-1945.
- [57] Redmon, J., Farhadi, A. YOLOv3: An incremental improvement. 2018; 1-6. Available from: abs/1804.02767.

- [58] Everingham, M., Gool, L.V., Williams, C.K., et al. The pascal visual object classes (VOC) challenge. *International Journal of Computer Vision*. 2009; 88: 303-338.
- [59] Everingham, M., Eslami, SM, Gool, L.V., et al. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision*. 2014; 111: 98-136.
- [60] Deng, J., Dong, W., Socher, R., et al. ImageNet: A large-scale hierarchical image database. *IEEE Conference on Computer Vision and Pattern Recognition*. 2009. p. 248-255.
- [61] Russakovsky, O., Deng, J., Su, H., et al. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*. 2015; 115: 211-252.
- [62] Lin, T., Maire, M., Belongie, S.J., et al. Microsoft COCO: Common objects in context. 2014; 1-15. Available from: [abs/1405.0312](https://arxiv.org/abs/1405.0312).
- [63] Kuznetsova, A., Rom, H., Alldrin, N., et al. The open images dataset V4. *International Journal of Computer Vision*. 2020; 128: 1956-1981.
- [64] Singh, B., Davis, L.S. An analysis of scale invariance in object detection-SNIP. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018. p. 3578-3587.
- [65] Hu, H., Gu, J., Zhang, Z., et al. Relation networks for object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2017. p. 3588-3597.
- [66] Cai, Z., Vasconcelos, N. Cascade R-CNN: Delving into high quality object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018. p. 6154-6162.
- [67] Wang, H., Wang, Q., Gao, M., et al. Multi-scale location-aware kernel representation for object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018. p. 1248-1257.
- [68] Tychsen-Smith, L., Petersson, L. Improving object localization with fitness NMS and bounded IoU loss. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018. p. 6877-6885.
- [69] Liu, S., Huang, D., Wang, Y. Receptive field block net for accurate and fast object detection. *European Conference on Computer Vision*. 2017. p. 404-419.
- [70] Law, H., Deng, J.B. CornerNet: Detecting objects as paired keypoints. *International Journal of Computer Vision*. 2019; 128: 642-656.
- [71] Kim, S., Kook, H., Sun, J., et al. Parallel feature pyramid network for object detection. *European Conference on Computer Vision*. 2018. p. 239-256.
- [72] Wang, R.J., Li, X., Ao, S., et al. Pelee: A real-time object detection system on mobile devices. *NeurIPS*. 2018; 1-10. Available from: [arXiv:1804.06882](https://arxiv.org/abs/1804.06882).
- [73] Jiang, C., Xu, H., Liang, X., et al. Hybrid knowledge routed modules for large-scale object detection. *NeurIP*. 2018; 1-12. Available from: [arXiv:1810.12681v1](https://arxiv.org/abs/1810.12681v1).
- [74] Zhao, Q., Sheng, T., Wang, Y., et al. M2Det: A single-shot object detector based on multi-level feature pyramid network. 2018; 1-8. *AAAI*. Available from: [arXiv:1811.04533](https://arxiv.org/abs/1811.04533).
- [75] Bae, S. Object detection based on region decomposition and assembly. *AAAI*. 2019; 1-8. Available from: [arXiv:1901.08225](https://arxiv.org/abs/1901.08225).
- [76] Zhu, R., Zhang, S., Wang, X., et al. ScratchDet: Training single-shot object detectors from scratch. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019. p. 2263-2272.
- [77] Pang, J., Chen, K., Shi, J., et al. Libra R-CNN: Towards balanced learning for object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019. p. 821-830.
- [78] Xu, H., Jiang, C., Liang, X., et al. Reasoning-RCNN: Unifying adaptive global reasoning into large-scale object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019. p. 6412-6421.
- [79] Zhu, C., He, Y., Savvides, M. Feature selective anchor-free module for single-shot object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019. p. 840-849.
- [80] Ghiasi, G., Lin, T., Pang, R., et al. NAS-FPN: Learning scalable feature pyramid architecture for object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019. p. 7029-7038.
- [81] Zhang, H., Chang, H., Ma, B., et al. Cascade RetinaNet: Maintaining consistency for single-stage object detection. *BMVC*. 2019; 1-12. Available from: [arXiv:1907.06881](https://arxiv.org/abs/1907.06881).
- [82] Li, Y., Chen, Y., Wang, N., et al. Scale-aware trident networks for object detection. *IEEE/CVF International Conference on Computer Vision (ICCV)*. 2019. p. 6053-6062.
- [83] Li, S., Yang, L., Huang, J., et al. Dynamic anchor feature selection for single-shot object detection. *IEEE/CVF International Conference on Computer Vision (ICCV)*. 2019. p. 6608-6617.
- [84] Xu, H., Yao, L., Li, Z., et al. Auto-FPN: Automatic network architecture adaptation for object detection beyond classification. *IEEE/CVF International Conference on Computer Vision (ICCV)*. p. 6648-6657.
- [85] Tian, Z., Shen, C., Chen, H., et al. FCOS: Fully convolutional one-stage object detection. *IEEE/CVF International Conference on Computer Vision (ICCV)*. 2019. p. 9626-9635.
- [86] Zhang, X., Wan, F., Liu, C., et al. FreeAnchor: Learning to match anchors for visual object detection. *NeurIPS*. 2019;

1-9. Available from: arXiv:1909.02466v2.

- [87] Chen, Y., Yang, T., Zhang, X., et al. DetNAS: Backbone search for object detection. *NeurIPS*. 2019; 1-12. Available from: arXiv:1903.10979.
- [88] Peng, J., Sun, M., Zhang, Z., et al. Efficient neural architecture transformation searchin channel-level for object detection. *NeurIPS*. 2019; 1-11. Available from: arXiv:1909.02293v1.
- [89] Zoph, B., Cubuk, E.D., Ghiasi, G., et al. Learning data augmentation strategies for object detection. 2019; 1-13. Available from: abs/1906.11172.
- [90] Tan, M., Pang, R., Le, Q.V. EfficientDet: Scalable and efficient object detection. 2019; 1-10. Available from: abs/1911.09070.