

# A Survey of Modern Deep Learning based Object Detection Models

Syed Sahil Abbas Zaidi, Mohammad Samar Ansari, Asra Aslam,  
Nadia Kanwal, Mamoona Asghar, and Brian Lee

**Abstract**—Object Detection is the task of classification and localization of objects in an image or video. It has gained prominence in recent years due to its widespread applications. This article surveys recent developments in deep learning based object detectors. Concise overview of benchmark datasets and evaluation metrics used in detection is also provided along with some of the prominent backbone architectures used in recognition tasks. It also covers contemporary lightweight classification models used on edge devices. Lastly, we compare the performances of these architectures on multiple metrics.

**Index Terms**—Object detection and recognition, convolutional neural networks (CNN), lightweight networks, deep learning

## I. INTRODUCTION

Object detection is a trivial task for humans. A few months old child can start recognizing common objects, however teaching it to the computer has been an uphill task until the turn of the last decade. It entails identifying and localizing all instances of an object (like cars, humans, street signs, etc.) within the field of view. Similarly, other tasks like classification, segmentation, motion estimation, scene understanding, etc, have been the fundamental problems in computer vision.

Early object detection models were built as an ensemble of hand-crafted feature extractors such as Viola-Jones detector [1], Histogram of Oriented Gradients (HOG) [2] etc. These models were slow, inaccurate and performed poorly on unfamiliar datasets. The re-introduction of convolutional neural network (CNNs) and deep learning for image classification changed the landscape of visual perception. Its use in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) 2012 challenge by AlexNet [3] inspired further research of its application in computer vision. Today, object detection finds application from self-driving cars and identity detection to security and medical uses. In recent years, it has seen exponential growth with rapid development of new tools and techniques.

This survey provides a comprehensive review of deep learning based object detectors and lightweight classification architectures. While existing reviews are quite thorough [4]–[7], most of them lack new developments in the domain. The main contributions of this paper are as follows:

S.S.A. Zaidi, N. Kanwal, M. Asghar and B. Lee are with the Athlone Institute of Technology, Ireland. M.S. Ansari is with the Aligarh Muslim University, India. A. Aslam is with the Insight Center for Data Analytics, National University of Ireland, Galway. (Emails: sahilzaidi78@gmail.com, samar.ansari@znect.ac.in, asra.aslam@insight-centre.org, nkanwal@ait.ie, masghar@ait.ie, blee@ait.ie)

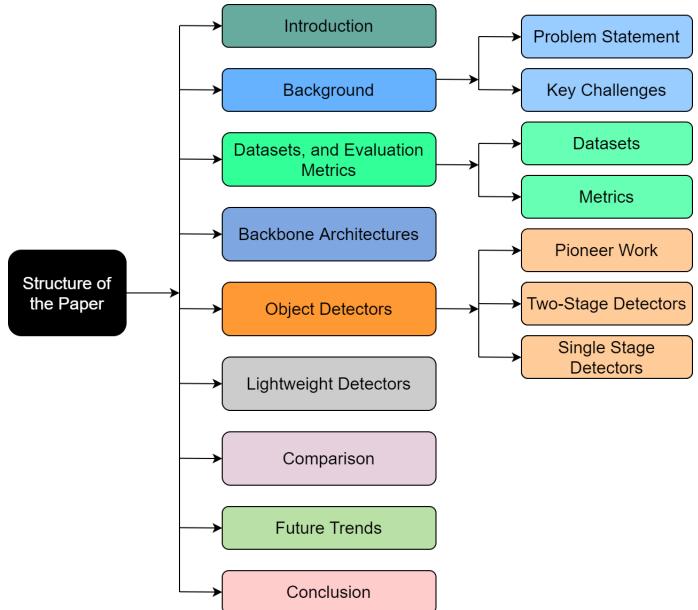


Fig. 1: Structure of the paper.

- 1) This paper provides an in-depth analysis of major object detectors in both categories – single and two stage detectors. Furthermore, we take historic look at the evolution of these methods.
- 2) We present a detailed evaluation of the landmark backbone architectures and lightweight models. We could not find any paper which provides a broad overview of both these topics.

In this paper, we have systematically reviewed various object detection architectures and its associated technologies, as illustrated in figure 1. Rest of this paper is organized as follows. In section II, the problem of object detection and its associated challenges are discussed. Various benchmark datasets and evaluation metrics are listed in Section III. In Section IV, several milestone backbone architectures used in modern object detectors are examined. Section V is divided into three major sub-section, each studying a different category of object detectors. This is followed by the analysis of a special classification of object detectors, called lightweight networks in section VI and a comparative analysis in Section VII. The future trends are mentioned in Section VIII while the paper is concluded in Section IX.



Fig. 2: Sample images from different datasets.

## II. BACKGROUND

### A. Problem Statement

The object detection is the natural extension of object classification, which aims only at recognizing the object in the image. The goal of the object detection is to detect all instances of the predefined classes and provide its coarse localization in the image by axis-aligned boxes. The detector should be able to identify all instances of the object classes and draw bounding box around it. It is generally seen as a supervised learning problem. Modern object detection models have access to large sets of labelled images for training and are evaluated on various canonical benchmarks.

### B. Key challenges in Object Detection

Computer vision has come a long way in the past decade, however it still has some major challenges to overcome. Some of these key challenges faced by the networks in real life applications are:

- *Intra class variation* : Intra class variation between the instances of same object is relatively common in nature. This variation could be due to various reasons like occlusion, illumination, pose, viewpoint, etc. These unconstrained external can have dramatic effect of the object appearance [5]. It is expected that the objects could have non-rigid deformation or be rotated, scaled or blurry. Some objects could have inconspicuous surroundings, making the extraction difficult.
- *Number of categories*: The sheer number of object classes available to classify makes it a challenging problem to solve. It also requires more high-quality annotated data, which is hard to come by. Using fewer examples for training a detector is an open research question.
- *Efficiency*: Present day models need high computation resources to generate accurate detection results. With mobile and edge devices becoming common place, efficient object detectors are crucial for further development in the field of computer vision.

## III. DATASETS AND EVALUATION METRICS

### A. Datasets

This section presents an overview of the datasets that are available, and have been most commonly used for object detection tasks.

1) *PASCAL VOC 07/12*: The Pascal Visual Object Classes (VOC) challenge was a multiyear effort to accelerate the development in the field of visual perception. It started in 2005 with classification and detection tasks on four object classes [8], but two versions of this challenges are mostly used as a standard benchmark. While the VOC07 challenge had 5k training images and more than 12k labelled objects [9], the VOC12 challenge increased them to 11k training images and more than 27k labelled objects [10]. Object classes was expanded to 20 categories and the tasks like segmentation and action detection were included as well. Pascal VOC introduced the mean Average Precision (*mAP*) at 0.5 IoU (Intersection over Union) to evaluate the performance of the models. Figure 3 depicts the distribution of the number of images w.r.t. to the different classes in the Pascal VOC dataset.

2) *ILSVRC*: The ImageNet Large Scale Visual Recognition Challenge (ILSVRC) [11] was an annual challenge running from 2010 to 2017 and became a benchmark for evaluating algorithm performance. The dataset size was scaled up to more than a million images consisting of 1000 object classification classes. 200 of these classes were hand-picked for object detection task, constitute of more than 500k images. Various sources including ImageNet [12] and Flickr, were used to construct detection dataset. ILSVRC also updated the evaluation metric by loosening the IoU threshold to help include smaller object detection. Figure 4 depicts the distribution of the number of images w.r.t. to the different classes in the ImageNet dataset.

3) *MS-COCO*: The Microsoft Common Objects in Context (MS-COCO) [13] is one of the most challenging datasets available. It has 91 common objects found in their natural context which a 4-year-old human can easily recognize. It was launched in 2015 and its popularity has only increased since then. It has more than two million instances and an

average of 3.5 categories per images. Furthermore, it contains 7.7 instances per image, comfortably more than other popular datasets. MS COCO comprises of images from varied viewpoints as well. It also introduced a more stringent method to measure the performance of the detector. Unlike the Pascal VOC and ILSVCR, it calculates the IoU from 0.5 to 0.95 in steps of 0.5, then using a combination of these 10 values as final metric, called Average Precision (AP). Apart from this, it also utilizes AP for small, medium and large objects separately to compare performance at different scales. Figure 5 depicts the distribution of the number of images w.r.t. to the different classes in the MS-COCO dataset.

4) *Open Image*: Google's Open Images [14] dataset is composed of 9.2 million images, annotated with image-level labels, object bounding boxes, and segmentation masks, among others. It was launched in 2017 and has received six updates. For object detection, Open Images has 16 million bounding boxes for 600 categories on 1.9 million images, which makes it the largest dataset of object localization. Its creators took extra care to choose interesting, complex and diverse images, having 8.3 object categories per image. Several changes were made to the AP introduced in Pascal VOC like ignoring un-annotated class, detection requirement for class and its subclass, etc. Figure 6 depicts the distribution of the number of images w.r.t. to the different classes in the Open Images dataset.

5) *Issues of Data Skew/Bias*: While observing Fig. 3 through Fig. 6, an alert reader would certainly notice that the number of images for difference classes vary significantly in all the datasets [15]. Three (Pascal VOC, MS-COCO, and Open Images Dataset) of the four datasets discussed above have a very significant drop in the number of images beyond the top-5 most frequent classes. As can be readily observed for Fig. 3, there are 13775 images which contain a 'person' and then 2829 images which contain a 'car'. The number of images for the remaining 18 classes in this dataset almost fall linearly to the 55 images of 'sheep'. Similarly, for the MS-COCO dataset, the class 'person' has 262465 images, and the next most-frequent class 'car' has 43867 images. The downward trend continues till there are only 198 images for the class 'hair drier'. A similar phenomenon is also observed in the Open Images Dataset, wherein the class 'Man' is the most frequent with 378077 images, and the class 'Paper Cutter' has only 3 images. This clearly represents a skew in the datasets and is bound to create a bias in the training process of any object detection model. Therefore, an object detection model trained on these skewed datasets will in all probability show better detection performance for the classes with more number of images in the training data. Although still present, this issue is slightly less pronounced in the ImageNet dataset, as can be observed from Fig. 4 from where it can be seen that the most frequent class i.e. 'koala' has 2469 images, and the least frequent class i.e. 'cart' has 624 images. However, this leads to another point of concern in the ImageNet dataset: the most frequent class is for 'koala' and the next most-appearing class is 'computer keyboard', which are clearly not the most sought after objects in a real-world object detection scenario (where person, cars, traffic signs, etc. are of higher concern).

## B. Metrics

Object detectors use multiple criteria to measure the performance of the detectors viz., frames per second (FPS), precision and recall. However, mean Average Precision (*mAP*) is the most common evaluation metric. Precision is derived from Intersection over Union (IoU), which is the ratio of the area of overlap and the area of union between the ground truth and the predicted bounding box. A threshold is set to determine if the detection is correct. If the IoU is more than the threshold, it is classified as True Positive while an IoU below it is classified as False Positive. If the model fails to detect an object present in the ground truth, it is termed as False Negative. Precision measures the percentage of correct predictions while the recall measure the correct predictions with respect to the ground truth.

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} = \frac{\text{True Positive}}{\text{All Observations}} \quad (1)$$

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} = \frac{\text{True Positive}}{\text{All Ground Truth}} \quad (2)$$

Based on the above equation, average precision is computed separately for each class. To compare performance between the detectors, the mean of average precision of all classes, called mean average precision (*mAP*) is used, which acts as a single metric for final evaluation.

## IV. BACKBONE ARCHITECTURES

Backbone architectures are one of the most important component of the object detector. These networks extract feature from the input image used by the model. Here, we have discussed some milestone backbone architectures used in modern detectors:

### A. *AlexNet*

Krizhevsky et al. proposed AlexNet [3], a convolutional neural network based architecture for image classification, and won the ImageNet Large-Scale Visual Recognition Challenge (ILSVRC) 2012 challenge. It achieved a considerably higher accuracy (more than 26%) than the contemporary models. AlexNet is composed of eight learnable layers - five convolutional and three fully connected layers. The last layer of the fully connected layer is connected to an *N*-way (*N*: number of classes) softmax classifier. It uses multiple convolutional kernels throughout the network to obtain features from the image. It also uses dropout and ReLU for regularization and faster training convergence respectively. The convolutional neural networks were given a new life by its reintroduction in AlexNet and it soon became the go-to technique in processing imaging data.







Fig. 6: (This image is best viewed in PDF form with magnification) Number of images for different classes annotated in the Open Images dataset [15]

#### F. CSPNet

Existing neural networks have shown incredible results in achieving high accuracy in computer vision tasks; however, they rely on excessive computational resources. Wang et al. believe that heavy inference computations can be reduced by cutting down the duplicate gradient information in the network. They proposed CSPNet [25] which creates different paths for the gradient flow within the network. CSPNet separates feature maps at the base layer into two parts. One part is passed through the partial convolution network block (e.g., Dense and Transition block in DenseNet [23] or Res(X) block in ResNeXt [24]) while the other part is combined with its outputs at a later stage. This reduces the number of parameters, increases the utilization of computation units and eases memory footprint. It is easy to implement and general enough to be applicable on other architectures like ResNet [21], ResNeXt [24], DenseNet [23], Scaled-YOLOv4 [26] etc. Applying CSPNet on these networks reduced computations from 10% to 20%, while the accuracy remained constant or improved. Memory cost and computational bottleneck is also reduced significantly with this method. It is leveraged in many state of the art detector models, while also being used for mobile and edge devices.

#### G. EfficientNet

Tan et al. systematically studied network scaling and its effects on the model performance. They summarized how altering network parameters like depth, width and resolution influence its accuracy. Scaling any parameter individually comes with an associated cost. Increasing depth of a network can help in capturing richer and more complex features, but they are difficult to train due to vanishing gradient problem. Similarly, scaling network width will make it easier to capture fine grained features but have difficulty in obtaining high level features. Gains from increasing the image resolution, like depth and width, saturate as model scales. In the paper [27], Tan et al. proposed the use of a compound coefficient that can uniformly scale all three dimensions. Each model parameter has an associated constant, which is found by fixing the coefficient as 1 and performing a grid search on a baseline network. The baseline architecture, inspired by their previous work [28], is developed by neural architecture search on a search target while optimizing accuracy and computations. EfficientNet is a simple and efficient architecture. It outperformed existing models in accuracy and speed while being considerably smaller. By providing a monumental increase in efficiency, it could potentially open a new era in the field of

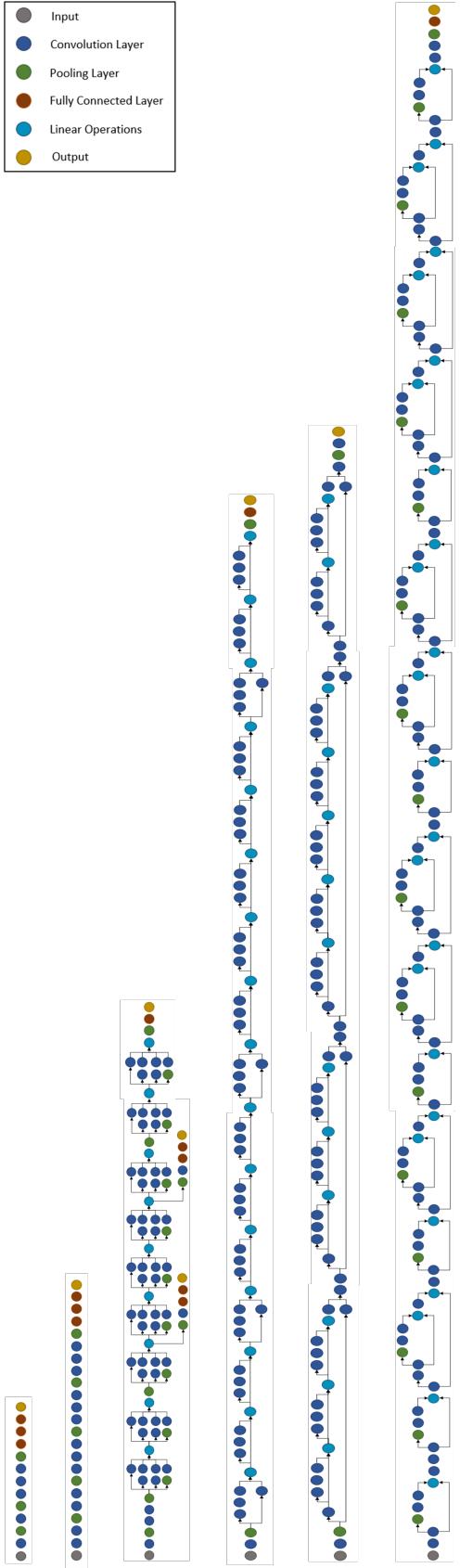


Fig. 7: Visualization of CNN Architectures<sup>1</sup>. Left to Right: AlexNet, VGG-16, GoogLeNet, ResNet-50, CSPResNeXt-50, EfficientNet-B4.

efficient networks.

## V. OBJECT DETECTORS

We have divided this review based on the two types of detectors — two-stage and single-stage detectors. However, we also discussed the pioneer work, where we briefly examine a few traditional object detectors. A network which has a separate module to generate region proposals is termed as a two-stage detector. These models try to find an arbitrary number of objects proposals in an image during the first stage and then classify and localize them in the second. As these systems have two separate steps, they generally take longer to generate proposals, have complicated architecture and lacks global context. Single-stage detectors classify and localize semantic objects in a single shot using dense sampling. They use predefined boxes/keypoints of various scale and aspect ratio to localize objects. It edges two-stage detectors in real-time performance and simpler design.

### A. Pioneer Work

1) *Viola-Jones*: Primarily designed for face detection, Viola-Jones object detector [1], proposed in 2001, was an accurate and powerful detector. It combined multiple techniques like Haar-like features, integral image, Adaboost and cascading classifier. First step is to search for Haar-like features by sliding a window on the input image and uses integral image to calculate. It then uses a trained Adaboost to find the classifier of each haar feature and cascades them. Viola Jones algorithm is still used in small devices as it is very efficient and fast.

2) *HOG Detector*: In 2005, Dalal and Triggs proposed the Histogram of Oriented Gradients (HOG) [2] feature descriptor used to extract features for object detection. It was an improvement over other detectors like [29]–[32]. HOG extracts gradient and its orientation of the edges to create a feature table. The image is divided into grids and the feature table is then used to create histogram for each cell in the grid. HOG features are generated for the region of interest and fed into a linear SVM classifier for detection. The detector was proposed for pedestrian detection; however, it could be trained to detect various classes.

3) *DPM*: Deformable Parts Model (DPM) [33] was introduced by Felzenszwalb et al. and was the winner Pascal VOC challenge in 2009. It used individual “part” of the object for detection and achieved higher accuracy than HOG. It follows the philosophy of *divide and rule*; parts of the object are individually detected during inference time and a probable arrangement of them is marked as detection. For example, a human body can be considered as a collection of parts like head, arms, legs and torso. One model will be assigned to capture one of the parts in the whole image and the process is repeated for all such parts. A model then removes improbable configurations of the combination of these parts to produce detection. DPM based models [34], [35] were one of the most successful algorithms before the era of deep learning.

<sup>1</sup>Tool Used: <https://netron.app/>

## B. Two-Stage Detectors

**1) R-CNN:** The Region-based Convolutional Neural Network (R-CNN) [36] was the first paper in the R-CNN family, and demonstrated how CNNs can be used to immensely improve the detection performance. R-CNN use a class agnostic region proposal module with CNNs to convert detection into classification and localization problem. A mean-subtracted input image is first passed through the region proposal module, which produces 2000 object candidates. This module find parts of the image which has a higher probability of finding an object using Selective Search [37]. These candidates are then warped and propagated through a CNN network, which extracts a 4096-dimension feature vector for each proposal. Girshick et al. used AlexNet [3] as the backbone architecture of the detector. The feature vectors are then passed to the trained, class-specific Support Vector Machines (SVMs) to obtain confidence scores. Non-maximum suppression (NMS) is later applied to the scored regions, based on its IoU and class. Once the class has been identified, the algorithm predicts its bounding box using a trained bounding-box regressor, which predicts four parameters i.e., center coordinates of box along with its width and height.

R-CNN has a complicated multistage training process. The first stage is pre-training the CNN with a large classification dataset. It is then fine-tuned for detection using domain-specific images (mean-subtracted, warped proposals) by replacing of the classification layer with a randomly initialized  $N+1$ -way classifier,  $N$  being the number of classes, using stochastic gradient descent (SGD) [38]. One liner SVM and bounding box regressor is trained for each class.

R-CNN ushered a new wave in the field of object detection, but it was slow (47 sec per image) and expensive in time and space [39]. It had complex training process and took days to train on small datasets even when some of the computations were shared.

**2) SPP-Net:** He et al. proposed the use of Spatial Pyramid Pooling (SPP) layer [40] to process image of arbitrary size or aspect ratio. They realized that only the fully connected part of the CNN required a fixed input. SPP-net [41] merely shifted the convolution layers of CNN before the region proposal module and added a pooling layer, thereby making the network independent of size/aspect ratio and reducing the computations. The selective search [37] algorithm is used to generate candidate windows. Feature maps are obtained by passing the input image through the convolution layers of a ZF-5 [16] network. The candidate windows are then mapped on to the feature maps, which are subsequently converted into fixed length representations by spatial bins of a pyramidal pooling layer. This vector is passed to the fully connected layer and ultimately, to SVM classifiers to predict class and score. Similar to R-CNN [36], SPP-net has as post processing layer to improve localization by bounding box regression. It also uses the same multistage training process, except that the fine tuning is done only on the fully connected layers.

SPP-Net is considerably faster than the R-CNN model with comparable accuracy. It can process images of any shape/aspect ratio and thus, avoid object deformation due to

input warping. However, as its architecture is analogous to R-CNN, it shared R-CNN's disadvantages too like multistage training, computationally expensive and training time as well.

**3) Fast R-CNN:** One of the major issues with R-CNN/SPP-Net was the need to train multiple systems separately. Fast R-CNN [39] solved this by creating a single end-to-end trainable system. The network takes as input an image and its object proposals. The image is passed through a set of convolution layers and the object proposals are mapped to the obtained feature maps. Girshick replaced pyramidal structure of pooling layers from SPP-net [41] with a single spatial bin, called RoI pooling layer. This layer is connected to 2 fully connected layer and then branches out into a  $N+1$ -class SoftMax layer and a bounding box regressor layer, which has a fully connected layer as well. The model also changed the loss function of bounding box regressor from L2 to smooth L1 to better performance, while introducing a multi-task loss to train the network.

The authors used modified version of existing state-of-art pre-trained models like [3], [17] and [42] as backbone. The network was trained in a single step by stochastic gradient descent (SGD) and a mini-batch of 2 images. This helped the network converge faster as the back-propagation shared computations among the RoIs from the two images.

Fast R-CNN was introduced as an improvement in speed (146x on R-CNN) while the increase in accuracy was supplementary. It simplified training procedure, removed pyramidal pooling and introduces a new loss function. The object detector, without the region proposal network, reported near real time speed with considerable accuracy.

**4) Faster R-CNN:** Even though Fast R-CNN inched closer to real time object detection, its region proposal generation was still an order of magnitude slower (2 sec per image compared to 0.2 sec per image). Ren et al. suggested a fully convoluted network [43] as a region proposal network (RPN) in [44] that takes an arbitrary input image and outputs a set of candidate windows. Each such window has an associated *objectness score* which determines likelihood of an object. Unlike its predecessors like [21], [34], [39] which used image pyramids to solve size variance of objects, RPN introduces Anchor boxes. It used multiple bounding boxes of different aspect ratios and regressed over them to localize object. The input image is first passed through the CNN to obtain a set of feature maps. These are forwarded to the RPN, which produces bounding boxes and their classification. Selected proposals are then mapped back to the feature maps obtained from previous CNN layer in RoI pooling layer, and ultimately fed to fully connected layer, which is sent to classifier and bounding box regressor. Faster R-CNN is essentially Fast R-CNN with RPN as region proposal module.

Training of Faster R-CNN is more convoluted, due to the presence of shared layers between two models which perform very different tasks. Firstly, RPN is pre-trained on ImageNet dataset [12] and fine-tuned on PASCAL VOC dataset [8]. A Fast R-CNN is trained from the region proposals of RPN from first step. Till this point, the networks do not have shared convolution layer. Now, we fix the convolution layers of the detector and fine-tune the unique layers in RPN. And finally,

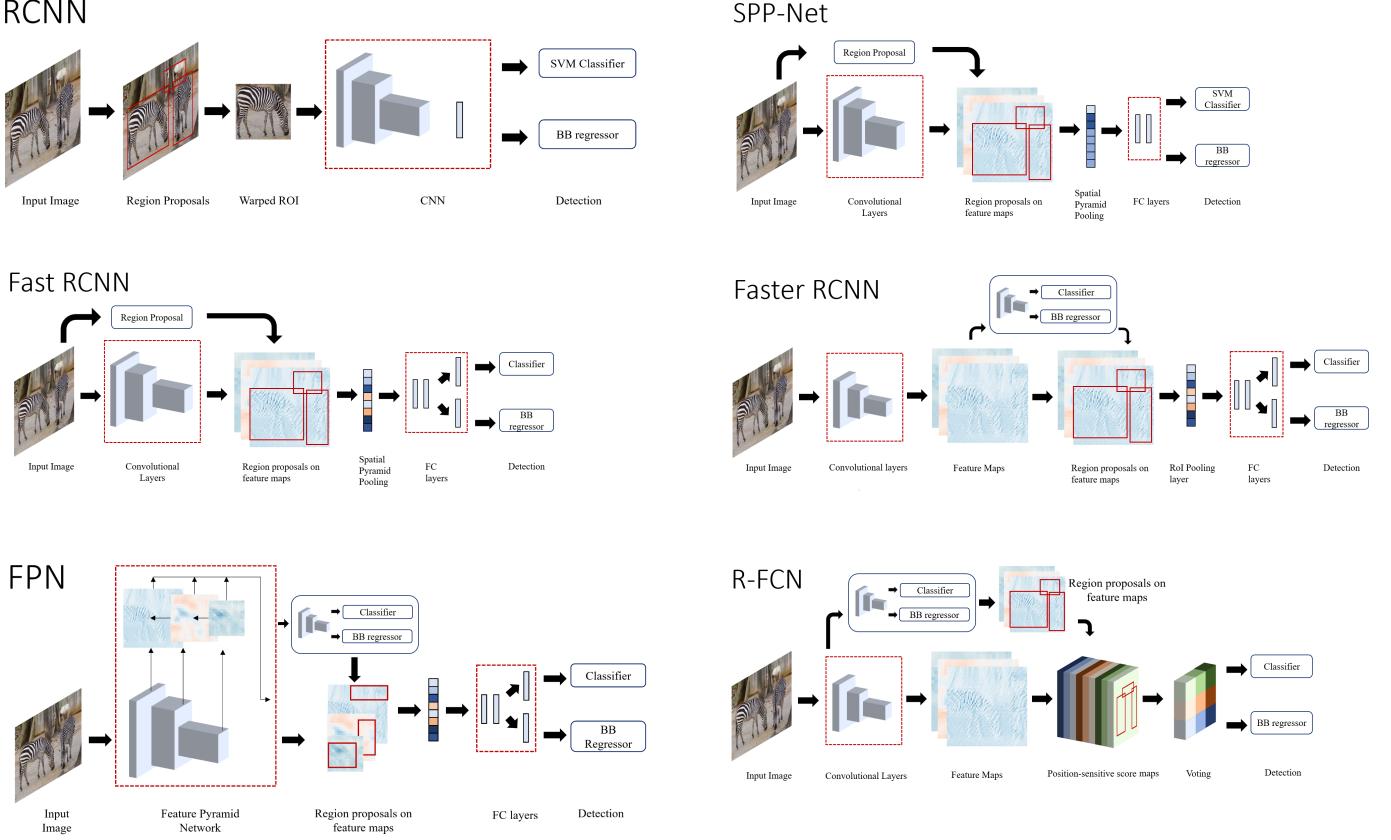


Fig. 8: Illustration of the internal architecture of different two stage object detectors<sup>2</sup>.

Fast R-CNN is fine-tuned from the updated RPN.

Faster R-CNN improved the detection accuracy over the previous state-of-art [39] by more than 3% and decreased inference time by an order of magnitude. It fixed the bottleneck of slow region proposal and ran in near real time at 5 frames per second. Another advantage of having a CNN in region proposal was that it could learn to produce better proposals and thereby increase accuracy.

5) **FPN:** Use of image pyramid to obtain feature pyramid (or *featurized image pyramids*) at multiple levels is a common method to increase detection of small objects. Even though it increases Average Precision of the detector, the increase in the inference time is substantial. Lin et al. proposed the Feature Pyramid Network (FPN) [45], which has a top-down architecture with lateral connections to build high-level semantic features at different scales. The FPN has two pathways, a bottom-up pathway which is a ConvNet computing feature hierarchy at several scales and a top-down pathway which upsamples coarse feature maps from higher level into high-resolution features. These pathways are connected by lateral connection by a  $1 \times 1$  convolution operation to enhance the semantic information in the features. FPN is used as a region proposal network (RPN) of a ResNet-101 [21] based Faster R-CNN here.

FPN could provide high-level semantics at all scales, which reduced the error rate in detection. It became a standard building block in future detections models and improved accuracy their accuracy across the table. It also lead to development of

other improved networks like PANet [46], NAS-FPN [47] and EfficientNet [27], which is current state of art detector.

6) **R-FCN:** Dai et al. proposed Region-based Fully Convolutional Network (R-FCN) [48] that shared almost all computations within the network, unlike previous two stage detectors which applied resource intensive techniques on each proposal. They argued against the use of fully connected layers and instead used convolutional layers. However, deeper layers in the convolutional network are translation-invariant, making them ineffective for localization tasks. The authors proposed the use of position-sensitive score maps to remedy it. These sensitive score maps encode relative spatial information of the subject and are later pooled to identify exact localization. R-FCN does it by dividing the region of interest into  $k \times k$  grid and scoring the likeliness of each cell with the detection class feature map. These scores are later averaged and used to predict the object class. R-FCN detector is a combination of four convolutional networks. The input image is first passed through the ResNet-101 [21] to get feature maps. An intermediate output (Conv4 layer) is passed to a Region Proposal Network (RPN) to identify ROI proposals while the final output is further processed through a convolutional layer and is input to classifier and regressor. The classification layer combines the generated the position-sensitive map with the ROI proposals to generate predictions while the regression network outputs the bounding box details. R-FCN is trained in a similar 4 step fashion as Faster-RCNN [44] whilst using a combined cross-entropy and box regression loss. It also adopts

online hard example mining (OHEM) [49] during the training.

Dai et al. offered a novel method to solve the problem of translation invariance in convolutional neural networks. R-FCN combines Faster R-CNN and FCN to achieve a fast, more accurate detector. Even though it did not improve accuracy by much, but it was 2.5-20 times faster than its counterpart.

**7) Mask R-CNN:** Mask R-CNN [50] extends on the Faster R-CNN by adding another branch in parallel for pixel-level object instance segmentation. The branch is a fully connected network applied on RoIs to classify each pixel into segments with little overall computation cost. It uses similar basic Faster R-CNN architecture for object proposal, but adds a mask head parallel to classification and bounding box regressor head. One major difference was the use of ROIAlign layer, instead of RoIPool layer, to avoid pixel level misalignment due to spatial quantization. The authors chose the ResNeXt-101 [24] as its backbone along with the feature Pyramid Network (FPN) for better accuracy and speed. The loss function of Faster R-CNN is updated with the mask loss and as in FPN, it uses 5 anchor boxes with 3 aspect ratio. Overall training of Mask R-CNN is similar to faster R-CNN.

Mask R-CNN performed better than the existing state of the art single-model architectures, added an extra functionality of instance segmentation with little overhead computations. It is simple to train, flexible and generalizes well in applications like keypoint detection, human pose estimation, etc. However, it was still below the real time performance ( $>30$  fps).

**8) DetectoRS:** Many contemporary two stage detectors like [44], [51], [52] use the mechanism of looking and thinking twice i.e. calculating object proposals first and using them to extract features to detect objects. DetectoRS [53] applies this mechanism at both macro and micro level of the network. At macro level, they propose Recursive Feature Pyramid (RFP), formed by stacking multiple feature pyramid network (FPN) with extra feedback connection from the top-down level path in FPN to the bottom-up layer. The output of the FPN is processed by the Atrous Spatial Pyramid Pooling layer (ASPP) [54] before passing it to the next FPN layer. A Fusion module is used to combine FPN outputs from different modules by creating an attention map. At micro level, Qiao et al. presented the Switchable Atrous Convolution (SAC) to regulate the dilation rate of convolution. An average pooling layer with  $5 \times 5$  filter and a  $1 \times 1$  convolution is used as a switch function to decide the rate of atrous convolution [55], helping the backbone detect objects at various scale on the fly. They also packed the SAC in between two global context modules [56] as it helps in making more stable switching. The combination of these two techniques, Recursive Feature Pyramid and Switchable Atrous Convolution results in DetectoRS. The authors incorporated the above techniques with the Hybrid Task Cascade (HTC) [51] as the baseline model and a ResNext-101 backbone.

DetectoRS combined multiple systems to improve performance of the detector and sets the state-of-the-art for the two stage detectors. Its RFP and SAC modules are well generalized and can be used in other detection models. However, it is not suitable for real time detections as it can only process about 4 frames per second.

### C. Single Stage Detectors

**1) YOLO:** Two stage detectors solve the object detection as a classification problem, a module presents some candidates which the network classifies as either an object or background. However, YOLO or You Only Look Once [57] reframed it as a regression problem, directly predicting the image pixels as objects and its bounding box attributes. In YOLO, the input image is divided into a  $S \times S$  grid and the cell where the object's center falls is responsible for detecting it. A grid cell predicts multiple bounding boxes, and each prediction array consists of 5 elements: center of bounding box – x and y, dimensions of the box – w and h, and the confidence score.

YOLO was inspired from the GoogLeNet model for image classification [18], which uses cascaded modules of smaller convolution networks [58]. It is pre-trained on ImageNet data [12] till the model achieves high accuracy and then modified by adding randomly initialized convolution and fully connected layers. At training time, grid cells predict only one class as it converges better, but it is increased during the inference time. Multitask loss, combined loss of all predicted components, is used to optimize the model. Non maximum suppression (NMS) removes class-specific multiple detections.

YOLO surpassed its contemporary single stage real time models by a huge margin in both accuracy and speed. However, it had significant shortcomings as well. Localization accuracy for small or clustered objects and limitation to number of objects per cell were its major drawbacks. These issues were fixed in later versions of YOLO [59]–[61].

**2) SSD:** Single Shot MultiBox Detector (SSD) [62] was the first single stage detector that matched accuracy of contemporary two stage detectors like Faster R-CNN [44], while maintaining real time speed. SSD was built on VGG-16 [17], with additional auxiliary structures to improve performance. These auxiliary convolution layers, added to the end of the model, decrease progressively in size. SSD detects smaller objects earlier in the network when the image features are not too crude, while the deeper layers were responsible for offset of the default boxes and aspect ratios [63].

During training, SSD match each ground truth box with the default boxes with the best jaccard overlap and train the network accordingly, similar to Multibox [63]. They also used hard negative mining and heavy data augmentation. Similar to DPM [33], it utilized weighted sum of the localization and confidence loss to train the model. Final output is obtained by performing non maximum suppression.

Even though SSD was significantly faster and more accurate than both state-of-art networks like YOLO and Faster R-CNN, it had difficulty in detecting small objects. This issue was later solved by using better backbone architectures like ResNet and other small fixes.

**3) YOLOv2 and YOLO9000:** YOLOv2 [59], an improvement on the YOLO [57], offered an easy tradeoff between speed and accuracy while the YOLO9000 model could predict 9000 object classes in real time. They replaced the backbone architecture of GoogLeNet [18] with DarkNet-19 [64]. It incorporated many impressive techniques like Batch Normalization [65] to improve convergence, joint training of classification and detection systems to increase detection

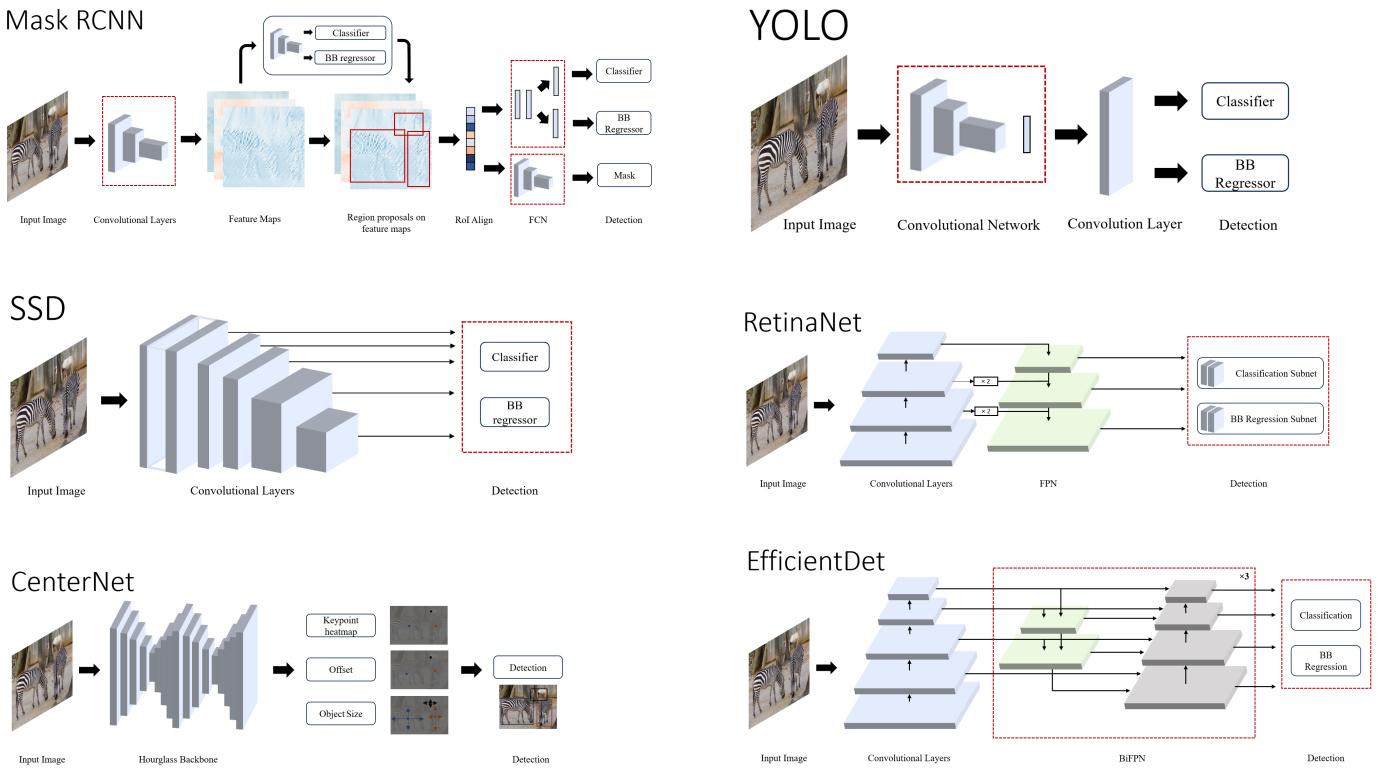


Fig. 9: Illustration of the internal architecture of different two and single stage object detectors<sup>2</sup>.

classes, removing fully connected layers to increase speed and using learnt anchor boxes to improve recall and have better priors. Redmon et al. also combined the classification and detection datasets in hierarchical structure using WordNet [66]. This WordTree can be used to predict a higher conditional probability of hypernym, even when the hyponym is not classified correctly, thereby increasing the overall performance of the system.

YOLOv2 provided better flexibility to choose the model on speed and accuracy, and the new architecture had fewer parameters. As the title of the paper suggests, it was “*better, faster and stronger*” [59].

4) **RetinaNet:** Given the difference between the accuracies of single and two stage detectors, Lin et al. suggested that the reason single stage detectors lag is the “extreme foreground-background class imbalance” [67]. They proposed a reshaped cross entropy loss, called Focal loss as the means to remedy the imbalance. Focal loss parameter reduces the loss contribution from easy examples. The authors demonstrate its efficacy with the help of a simple, single stage detector, called RetinaNet [67], which predicts objects by dense sampling of the input image in location, scale and aspect ratio. It uses ResNet [21] augmented by Feature Pyramid Network (FPN) [45] as the backbone and two similar subnets - classification and bounding box regressor. Each layer from the FPN is passed to the subnets, enabling it to detect objects at various scales. The classification subnet predicts the object score for each location while the box regression subnet regresses the offset for each anchor to the ground truth. Both subnets are small FCN and share parameters across the individual networks. Unlike most

previous works, the authors employ a class-agnostic bounding box regressor and found them to be equally effective.

RetinaNet is simple to train, converges faster and easy to implement. It achieved better performance in accuracy and run time than the two stage detectors. RetinaNet also pushed the envelope in advancing the ways object detectors are optimized by the introduction of a new loss function.

5) **YOLOv3:** YOLOv3 had “incremental improvements” from the previous YOLO versions [57], [59]. Redmon et al. replaced the feature extractor network with a larger Darknet-53 network [64]. They also incorporated various techniques like data augmentation, multi-scale training, batch normalization, among others. Softmax in classifier layer was replaced by a logistical classifier.

Even though YOLOv3 was faster than YOLOv2 [59], it lacked any ground breaking change from its predecessor. It even had lesser accuracy than an year old state-of-the-art detector [67].

6) **CenterNet:** Zhou et al. in [68] takes a very different approach of modelling objects as points, instead of the conventional bounding box representation. CenterNet predicts the object as a single point at the center of the bounding box. The input image is passed through the FCN that generates a heatmap, whose peaks correspond to center of detected object. It uses a ImageNet pretrained stacked Hourglass-101 [69] as the feature extractor network and has 3 heads – heatmap head to determine the object center, dimension head to estimate size of object and offset head to correct offset of object point. Multitask loss of all three heads is back propagated to feature

<sup>2</sup>Features created using: <https://poloclub.github.io/cnn-explainer/>

extractor while training. During inference, the output from offset head is used to determine the object point and finally a box is generated. As the predictions, not the result, are points and not bounding boxes, non-maximum suppression (NMS) is not required for post-processing.

CenterNet brings a fresh perspective and set aside years of progress in the field of object detection. It is more accurate and has lesser inference time than its predecessors. It has high precision for multiple tasks like 3D object detection, keypoint estimation, pose, instance segmentation, orientation detection and others. However, it requires different backbone architectures as general architectures that work well with other detectors give poor performance with it and vice-versa.

**7) EfficientDet:** EfficientDet [70] builds towards the idea of scalable detector with higher accuracy and efficiency. It introduces efficient multi-scale features, BiFPN and model scaling. BiFPN is bi-directional feature pyramid network with learnable weights for cross connection of input features at different scales. It improves on NAS-FPN [47], which required heavy training and had complex network, by removing one-input nodes and adding an extra lateral connection. This eliminates less efficient nodes and enhances high-level feature fusion. Unlike existing detectors which scale up with bigger, deeper backbone or stacking FPN layers, EfficientDet introduces a compounding coefficient which can be used to “jointly scale up all dimensions of backbone network, BiFPN network, class/box network and resolution” [70]. EfficientDet utilizes EfficientNet [27] as the backbone network with multiple sets of BiFPN layers stacked in series as feature extraction network. Each output from the final BiFPN layer is sent to class and box prediction network. The model is trained using SGD optimizer along with synchronized batch normalization and uses swish activation [71], instead of the standard ReLU activation, which is differentiable, more efficient and has better performance.

EfficientDet achieves better efficiency and accuracy than previous detectors while being smaller and computationally cheaper. It is easy to scale, generalizes well for other tasks and is the current state-of-the-art model for single-stage object detection.

**8) YOLOv4:** YOLOv4 [61] incorporated a lot of exciting ideas to design a fast and easy to train object detector that could work in existing production systems. It utilizes “bag of freebies” i.e., methods that only increase training time and do not affect the inference time. YOLOv4 utilizes data augmentation techniques, regularization methods, class label smoothing, CIoU-loss [72], Cross mini-Batch Normalization (CmBN), Self-adversarial training, Cosine annealing scheduler [73] and other tricks to improve training. Methods that only affect the inference time, called “Bag of Specials”, are also added to the network, including Mish activation [74], Cross-stage partial connections (CSP) [25], SPP-Block [41], PAN path aggregated block [46], Multi input weighted residual connections (Mi-WRC), etc. It also used genetic algorithm for searching hyperparameter. It has an ImageNet pre-trained CSPNetDarknet-53 backbone, SPP and PAN block neck and YOLOv3 as detection head.

Most existing detection algorithms require multiple GPUs to train model, but YOLOv4 can be easily trained on a single

GPU. It is twice as fast as EfficientDet with comparable performance. It is the state-of-the-art for real time single stage detectors.

**9) Swin Transformer:** Transformers [75] have had a profound impact in the Natural Language Processing (NLP) domain since its inception. Its application in language models like BERT (Bidirectional Encoder Representation from Transformers) [76], GPT (Generative Pre-trained Transformer) [77], T5 (Text-To-Text Transfer Transformer) [78] etc. have pushed the state of the art in the field. Transformers [75] uses the attention model to establish dependencies among the elements of the sequence and can attend to longer context than other sequential architectures. The success of transformers in NLP sparked interest in its application in computer vision. While CNNs have been the backbone on advancement in vision, they have some inherent shortcomings like the lack of importance of global context, fixed post-training weights [79] etc.

Swin Transformer [80] seeks to provide a transformer based backbone for computer vision tasks. It splits the input images in multiple, non-overlapping patches and converts them into embeddings. Numerous Swin Transformer blocks are then applied to the patches in 4 stages, with each successive stage reducing the number of patches to maintain hierarchical representation. The Swin Transformer block is composed of local multi-headed self-attention (MSA) modules, based on alternating shifted patch window in successive blocks. Computation complexity becomes linear with image size in local self-attention while shifted window enables cross-window connection. [80] also shows how shifted windows increase detection accuracy with little overhead.

Transformers present a paradigm shift from the CNN based neural networks. While its application in vision is still in a nascent stage, its potential to replace convolution from these tasks is very real. Swin Transformer achieved the state-of-the-art on MS COCO dataset, but utilises comparatively higher parameters than convolutional models.

## VI. LIGHTWEIGHT NETWORKS

A new branch of research has shaped up in recent years, aimed at designing small and efficient networks for resource constrained environments as is common in Internet of Things (IoT) deployments [81]–[84]. This trend has percolated to the design of potent object detectors too. It is seen that although a large number of object detectors achieve excellent accuracy and perform inference in real-time, a majority of these models require excessive computing resources and therefore cannot be deployed on edge devices.

Many different approaches have shown exciting results in the past. Utilization of efficient components and compression techniques like pruning ([85], [86]), quantization ([87], [88]), hashing [89], etc. have improved the efficiency of deep learning models. Use of trained large network to train smaller models, called distillation [90], has also shown interesting results. However in this section, we explore some prominent examples of efficient neural network design for achieving high performance on edge devices.

### A. SqueezeNet

Recent advances in the field of CNNs had mostly focused on improving the state-of-the-art accuracy on the benchmark datasets, which led to an explosion of model size and their parameters. But in 2016, Iandola et al. proposed a smaller, smarter network called SqueezeNet [91], which reduced the parameters while maintaining the performance. They achieved it by employing three main design strategies viz. using smaller filters, decreasing the number of input channels to  $3 \times 3$  filters and placing downsampling layers later in the network. The first two strategies decrease the number of parameters while attempting to preserve the accuracy and the third strategy increases the accuracy of the network. The building block of SqueezeNet is called a fire module, which consist of two layers: a squeeze layer and an expand layer, each with a ReLU activation. The squeeze layer is made up of multiple  $1 \times 1$  filters while the expand layer is a mix of  $1 \times 1$  and  $3 \times 3$  filters, thereby limiting the number of input channels. The SqueezeNet architecture is composed of a stack of 8 Fire modules squashed in between the convolution layers. Inspired by ResNet [21], SqueezeNet with residual connections was also proposed which increased the accuracy over the vanilla model. The authors also experimented with Deep Compression [87] and achieved  $510\times$  reduction in model size compared to AlexNet, while maintaining the baseline accuracy. SqueezeNet presented a good candidate for improving the hardware efficiency of the neural network architectures.

### B. MobileNets

MobileNet [92] moved away from the conventional methods of small models like shrinking, pruning, quantization or compressing, and instead used efficient network architecture. The network used depthwise separable convolution, which factorizes a standard convolution into a depthwise convolution and a  $1 \times 1$  pointwise convolution. A standard convolution uses kernels on all input channels and combines them in one step while the depthwise convolution uses different kernels for each input channel and uses pointwise convolution to combine inputs. This separation of filtering and combining of features reduces the computation cost and model size. MobileNet consists of 28 separate convolutional layers, each followed by batch normalization and ReLU activation function. Howard et al. also introduced the two model shrinking hyperparameters: width and resolution multiplier, in order to further improve speed and reduce size of the model. The width multiplier manipulates the width of the network uniformly by reducing the input and output channels while the resolution multiplier influences the size of the input image and its representations throughout the network. MobileNet achieves comparable accuracy to some full-fledged models while being a fraction of their size. Howard et al. also showed how it could generalize over various applications like face attribution, geolocalization and object detection. However, it was too simple and linear like VGG and therefore had fewer avenues for gradient flow. These were fixed in later iterations of this model [93], [94].

### C. ShuffleNet

In 2017, Zhang et al. introduced ShuffleNet [95], an extremely computationally efficient neural network architecture, specifically designed for mobile devices. They recognized that many efficient networks become less effective as they scale down and purported it to be caused by expensive  $1 \times 1$  convolutions. In conjunction with channel shuffle, they proposed the use of group convolution to circumvent its drawback of limited information flow. ShuffleNet consists mainly of a standard convolution followed by stacks of ShuffleNet units grouped in three stages. The ShuffleNet unit is similar to the ResNet block where they use depthwise convolution in the  $3 \times 3$  layer and replace the  $1 \times 1$  layer with pointwise group convolution. The depthwise convolution layer is preceded by a channel shuffle operation. The computation cost of the ShuffleNet can be administered by two hyperparameters: group number to control the connection sparsity and scaling factor to manipulate the model size. As group numbers become large, the error rate saturates as the input channels to each group decreases and therefore may reduce the representational capabilities. ShuffleNet outperformed contemporary models ([3], [18], [91], [92]) while having considerably smaller size. As the only advancement in ShuffleNet was channel shuffle, there isn't any improvement in inference speed of the model.

### D. MobileNetv2

Improving on MobileNetv1 [92], Sandler et al. proposed MobileNetv2 [93] in 2018. It introduced the inverted residual with linear bottleneck, a novel layer module to reduce computation and improve accuracy. The module expands a low-dimensional representation of the input into high dimension, filters with a depthwise convolution and then projects it back to low dimension, unlike the common residual block which performs compression, convolution and then expansion operations. The MobileNetv2 contains a convolution layer followed by 19 residual bottleneck modules and subsequently two convolutional layers. The residual bottleneck module has a shortcut connection only when the stride is 1. For higher stride, the shortcut is not used because of the difference in dimensions. They also employed ReLU6 as the non-linearity function, instead of simple ReLU, to limit computations. For object detection, the authors used MobileNetv2 as the feature extractor of a computationally efficient variant of the SSD [62]. This model, called SSDLite, claimed to have 8x fewer parameters than the original SSD while achieving competitive accuracy. It generalizes well over on other datasets, is easy to implement and hence, was well-received by the community.

### E. PeleeNet

Existing lightweight deep learning models like [92], [93], [95] relied heavily on depthwise separable convolution, which lacked efficient implementation. Wang et al. proposed a novel efficient architecture based on conventional convolution, named PeleeNet [96], using an assortment of computation conserving techniques. PeleeNet was centered around the DenseNet [23] but looked at many other models for inspiration. It introduced two-way dense layers, stem block, dynamic









