

Rank of Experts: Detection Network Ensemble

Seung-Hwan Bae¹, Youngwan Lee², Youngjoo Jo², Yuseok Bae², Joong-won Hwang²

¹Computer Vision Laboratory, Incheon National University, South Korea

² Electronics and Telecommunications Research Institute, Daejeon, South Korea

shbae@inu.ac.kr, {yw.lee, run.youngjoo, baeys, jwhwang}@etri.re.kr

Abstract

The recent advances of convolutional detectors show impressive performance improvement for large scale object detection. However, in general, the detection performance usually decreases as the object classes to be detected increases, and it is a practically challenging problem to train a dominant model for all classes due to the limitations of detection models and datasets. In most cases, therefore, there are distinct performance differences of the modern convolutional detectors for each object class detection. In this paper, in order to build an ensemble detector for large scale object detection, we present a conceptually simple but very effective class-wise ensemble detection which is named as Rank of Experts. We first decompose an intractable problem of finding the best detections for all object classes into small subproblems of finding the best ones for each object class. We then solve the detection problem by ranking detectors in order of the average precision rate for each class, and then aggregate the responses of the top ranked detectors (i.e. experts) for class-wise ensemble detection. The main benefit of our method is easy to implement and does not require any joint training of experts for ensemble. Based on the proposed Rank of Experts, we won the 2nd place in the ILSVRC 2017 object detection competition.

1. Introduction

An object detection problem is to predict object hypotheses of one or more classes given an image. In general, an object detection process consists of feature extraction, region search and object classification, and classical object detectors are constructed by combing them. However, the recent advances of deep learning change the detection paradigm. In particular, the region-based convolutional detection approach [41, 8, 36] enables us to detect large classes of objects with a single convolutional neural network (CNN).

In the recent years, in order to detect large classes of objects, deep CNNs [38, 13, 39] with high classification rate trained on the ImageNet dataset [35] are used for ex-

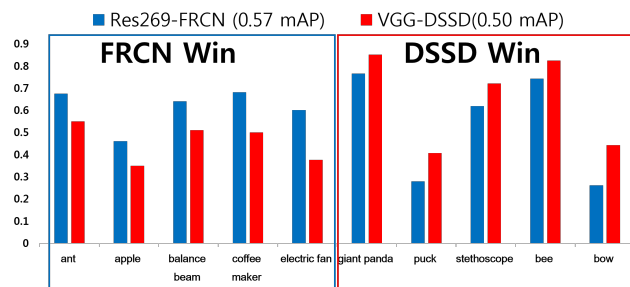


Figure 1. Performance comparison between FRCN [34] using ResNet [13] with 269 layers and DSSD [13] using VGG [38] with 16 layers on the ILSVRC 2017 val2 set. Even though the FRCN has the better mAP, the DSSD shows the better average precision rate for some object classes.

tracting object features and combined with several meta-architectures (i.e. detectors) such as Faster RCNN (FRCN) [34], SSD [27], DSSD [7] and R-FCN[3]. However, it is still challenging to train a single dominant model for all object classes due to the limitations of models and datasets. Also, in general, the detection performance is degraded as object classes increase.

To resolve this problem, we therefore propose an ensemble detection in order to combine different types of detectors for large scale object detection. In this paper we named the proposed ensemble detection as Rank of Experts since it can determine best detectors (or experts) for each class object detection. The motivation of our ensemble detection is that detectors with different network structures have different detection rate for object classes in general. In addition, we presume that a detector with the low mean average precision (mAP) can detect object instances of some classes better than a detector with high mAP. To support our assumption, we compare the average precision between the FRCN and DSSD. As shown in Fig. 1, for some object classes the DSSD shows the better detection rates even though the FRCN has the higher mAP than the DSSD.

Based on the assumption, in this study our main idea is to develop the class-wise ensemble detection which can determine the best experts and combine their detections ef-

fectively for detecting instances of the specific object class. To this end, we first generate a pool of detectors. In order to enhance the diversity of the detectors, we trained a set of detectors with various feature extractors and meta-architectures. Once the detector pool is generated, we generate a ranking matrix which represents the ranking of the detectors according to their average precision on each class. We then combine responses of top rankers of the pool for detecting the instances of the corresponding class.

Although our method is rather simple, there are obvious benefits of the proposed ensemble detection. Most of all, it can be easily implemented and does not require any extra joint training of detectors for ensemble. In addition, in most cases, the performance of the ensemble model is better than that of best single detector since it can selectively combine the good detection results of multiple detectors.

Recently, the 1st place winners at ILSVRC [35] & MSCOCO [26] 2015/2016 competitions improve the performance further using the model ensemble. MSRA [12] generates a union set of region proposals of models and the set is processed by an ensemble of per-region classifiers. However, the computational complexity increases significantly as the number of models increases since region proposals of whole models should be processed by each individual detector. For more efficient network ensemble, Google [14] combines models with complementary strength only. To this end, they greedily select models with the highest mAP on a validation set, but prune away a model which shows high similarity for the category AP with the selected ones. However, mAP is an indirect metric to select models for class-wise ensemble since high mAP does not always ensure the superiority on class-wise performance as shown in Fig. 1. On the other hand, we use the more direct and meaningful average precision metric for class-wise ensemble. As a result, we can select and combine the best experts for each category detection, and it can improve the class-wise performance.

To prove the effects of our ensemble detection, we generate a detector pool consisting of different types of detectors. Each detector was implemented based on recently developed feature extractors (ResNet101/152/269 [12], VGGNet [38] and WRI(wide-residual-inception)Net [24]) and meta-architectures (FRCN [34], SSD [27] and DSSD [7]). By combining them based on our ensemble method, we greatly improve the mAP up to about 4% and 5% on PASCAL VOC 2007/2012 (07/12) and ILSVRC 2017 datasets, respectively, as shown in Table 3 and 5. Remarkably, we won the 2nd place in the ILSVRC 2017 detection competition.

To summarize, our main contributions are as follows: (1) proposition of a new ensemble method to combine different detectors effectively without extra joint training and improve the class-wise performance directly, (2) extensive implementation of modern convolutional detectors with vari-

ous feature extractors and meta-architectures to generate a detector pool with high diversity, (3) thorough comparison of the proposed methods with recent detection and ensemble methods, and the details of our approach for winning the object detection challenge.

2. Related Work

We discuss previous works on deep networks and network ensemble for detection, which are related to our work.

Detection network: Since Krizhevsky et al. [20] achieved remarkable progress using CNN in object classification, a series of detection methods based on CNNs have been flourished. In RCNN [9], CNN is used to extract feature maps from warped regions of interest (RoI) generated by selective search [42]. Spatial pyramid pooling network [18] improved RCNN by sharing feature maps of images instead of extracting a CNN feature for each RoI.

For end-to-end training, Fast RCNN [8] introduces a single-stage training algorithm which can jointly learn to classify object proposals and refine their locations. In addition, Faster RCNN [34] further improves the speed and accuracy by embedding a region proposal network into the Fast RCNN.

In addition, many unified detection frameworks such as YOLO [33], SSD [27], and R-FCN [3] trained with a multi-task loss function for object classification and regression have been developed for the better speed. On the other hand, feature pyramid network [25] improves the accuracy and reduces computational complexity by constructing feature pyramid instead of image pyramid. Deformable convolutional network [16] proposes deformable convolution to enhance the robustness of the CNN against the geometric transformation. Note that this paper mainly focuses on combining the modern detectors effectively for better detection rather than improving their performance.

Network ensemble: The idea of ensemble learning, which combines multiple models, has been widely used due to two important benefits. The first one is that it can usually enhance the generalization performance. In fact, even models with similar performance may show different generalization ability. By combining several models with majority voting [10] for the classification or weighted averaging [6, 28, 31] for regression, the overall risk of making a particularly poor selection can be reduced [32].

When creating an ensemble system, increasing the diversity of the models can reduce the generalization error [32, 21] in general. Therefore, it is important to make models produce different errors to represent different regions of input space. For achieving this, training models with data resampling [10, 21, 15, 17] is one of the most popular methods. Training CNNs with different initial states and parameters or the order of mini-batches can also generate different outputs. AlexNet [20] increases the diversity by training the

same CNNs with different initial parameters. VGGNet [38] selects two best performing models among five models with distinct architecture. GoogleNet [40] achieves the diversity using the resampling.

Moreover, the other benefit is that ensemble learning can handle a complex task based on a divide-and-conquer approach. It divides the task into small subproblems and solve them with multiple learners. Mixture of experts [15, 37] that controls the activation of experts with a gate function on a given subtask is representative of this approach.

For combining models, the methods [31, 21, 17, 45] using the weighted averages of network outputs are most well-known. A new layer [30] or gate networks [22, 37, 4] are trained to select appropriate models for a given task. Recently, He et al. [12] collect region proposals of several models and forwards them to each RCNN. Huang et al. [14] select dissimilar models by using the cosine similarity.

In this work, for improving the generality and diversity we build a detector pool with modern detectors and effectively aggregate their outputs by using our Rank of Experts.

3. Detector Pool

To enhance the diversity of our ensemble detection, we generate a detector pool which contains convolutional detectors with various feature extractors and meta-architectures. In particular, we have implemented several modern convolutional detectors. We have combined them with the recent feature extractors, and trained them with the multi-task loss function [8] for object classification and localization. We provide the details of the implementations and evaluation results of each detector on PASCAL VOC & ILSVRC datasets as in Table 3 and 5, respectively.

3.1. Feature Extractor

For feature extraction, we use the ResNet [13], VGGNet [38] and WRINet [24]. A ResNet uses identity mapping with shortcut connection between feature maps to propagate signals effectively in a deep network. We implement several ResNets with different layers (101/152/269) by stacking different number of residual blocks. We also use the modified atrous VGG16 for reducing network parameters [27]. Compared to original VGG16 [38], fully connected layers are replaced with convolutional layers or removed to reduce the computational cost and memory usage, but its accuracy is similar to [38].

A WRINet [24] has a shallow and wide structure for fast feature extraction with low memory. For accurate feature extraction for objects of various sizes, it is designed based on the combination of residual [13] and inception [40] modules. In specific, it consists of the basic residual blocks with two consecutive 3x3 convolutional layers and identity mapping, and inception-residual blocks with 1x1, 3x3 and 5x5 convolutional layers which are concatenated and followed

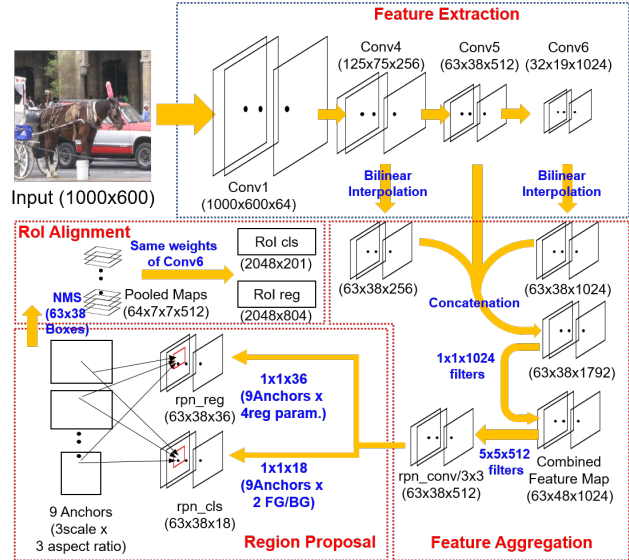


Figure 2. The proposed FRCN-Type2 detector with feature aggregation and RoI Alignment.

by identity mapping. It has achieved higher accuracy than the ResNet-101 with three times faster speed on the KITTI and CIFAR-10/100 datasets.

Before combining extractors with meta-architecture, it is possible to train them on ILSVRC [35] or other datasets to enhance their generalization ability. However, we use the pre-trained ResNet [44, 1], VGGNet [27] and WRINet [24] models only opened to the public because the focus of our work is not on improving the feature extractors.

3.2. Meta-Architecture

Based on the modern convolutional detectors (FRCN [34] and SSD [27]), we implement the original and extended versions using feature map aggregation [19], deconvolution [7] and RoI alignment [11]. The implemented detectors are listed in Table 3 and 5. All the detectors are designed based on a single CNN, and have in common the box (or region) proposal, box classification and refinement components. The parameters of the components are trained with confidences and locations of predicted and ground truth boxes with a multi-task loss function through end-to-end learning, where the loss is a weighted sum of the classification L_{cls} and localization L_{reg} losses as $L = L_{cls} + L_{reg}$ (see [8] for more details of L_{reg} and L_{reg}).

For improving the performance, we implement FRCN-Type1 with the following modifications : (1) we make the ratio of positive and negative samples in each min-batch equal, and (2) we allow region proposals of the small sizes (≥ 0) to be used for training in order to detect small objects. These effects are proved in Table 2.

In FRCN, a feature map of the fixed size (e.g. 7x7) is extracted by the RoIPooling operation from each proposal,

and the feature map is fed to higher layers for classification and box regression. We also implement the variant of FRCN by replacing the RoIPooling layer with the RoIAlignment layer [11], which uses bilinear interpolation instead of the hard quantization. In addition, since exploiting multiple feature maps enhances the FRCN accuracy as shown [19], we implement FRCN-Type2 by aggregating multi-feature of several layers and use the aggregated feature for region/RoI proposal and classifications as shown in Fig. 2. However, our implementation is different from [19]: we aggregate feature maps of different layers (Conv4/5/6) using bilinear interpolation to resize them instead of aggregating feature maps of layers (Conv1/3/5) using pooling and deconvolution. As a result, we can achieve the similar accuracy to [19] but improve the speed by avoiding the expensive pooling and deconvolution operations, and using features of the small sizes for aggregation.

We also implement SSD [27] and DSSD [7]. DSSD improves the SSD for the small object detection by merging encoded and decoded feature maps. Remarkably, as discussed in [7, 14], SSD and DSSD show lower accuracies for the small objects but higher those for the large objects than FRCN. Thus, we apply those detectors for ensemble detection to improve the robustness against object scale changes.

4. Rank of Experts

As mentioned, the main benefits of network ensemble are to improve the generalization performance and handle a difficult task using the divide-and-conquer approach. For large scale object detection, we also use the network ensemble for improving the generalization ability by combining multiple different detectors and dividing a difficult detection problem into small subproblems, where we reduce an intractable problem of finding the best detection responses for all the object classes to small subproblems of finding the best ones for each object class.

In the most existing ensemble methods, a joint training between multiple networks is usually required. The weights of the networks (*i.e.* beliefs in networks) are learned jointly by minimizing the generalize loss of an ensemble model [15, 45, 32], where the loss is evaluated with the residual between desired outputs and the networks' combined outputs with weights. A new layer is trained to control the amount of feature maps between the networks [30]. To allocate appropriate experts for a given task, they [15, 17, 14, 22] train an additional gating network together.

However, we found that it is not effective to apply the traditional ensemble learning approach directly for large scale object detection. The main reason is that learning several deep networks (*e.g.* ResNet and GoogleNet) together is not practical because of huge GPU memory usage and expensive computational complexity. In addition, the diversity between networks is not much when training several detec-

tors using only data resampling as done in [15, 17, 21, 45].

Unlike the previous works, the proposed Rank of Experts is designed for ensemble appropriate for large scale object detection. This allows us to combine deep convolutional detectors without the expensive joint training between them. By exploiting the proposed ensemble method, we can determine best convolutional detectors and combine their outputs for the specific object class detection. In the proposed object detection, we further improve their diversity by training detectors with various feature extractors and meta-architectures as well as data resampling and augmentation as discussed in Sec. 3 and 5. Our Rank of Experts is written in Python.

4.1. Ensemble Detection Formulation

Given an input image, a trained detector T_i generates a set of detections $D_i = \{\mathbf{d}_k | k = 1, \dots, M_i\}$, where M_i is the number of detection boxes of T_i . Each detection is represented as $\mathbf{d}_k = (\mathbf{p}_k, s_k, l_k)$, where \mathbf{p}_k , s_k and l_k are the bounding box left top and right bottom position, confidence score and object class. When trained N detectors are applied for an image, all detections of N detectors can be defined as $\mathbb{D} = \bigcup_{i=1}^N D_i$. Similarly, we denote a set of ground truth boxes in the image as \mathbb{G} . Then, a detection problem can be formulated to find the optimal \mathbb{D} by minimizing the loss function as

$$\hat{\mathbb{D}} = \underset{\mathbb{D}}{\operatorname{argmin}} L_{cls}(\mathbb{D}, \mathbb{G}) + L_{reg}(\mathbb{D}, \mathbb{G}) \quad (1)$$

where the same classification L_{cls} and regression loss L_{reg} are used as defined in [8]. Since the ground truth is not available, the solving Eq. (1) is not feasible in practice. To resolve this problem, we divide the problem Eq. (1) of finding \mathbb{D} for all classes to small subproblems of finding \mathbb{D}^j for class j . Then, the problem can be solved in the following two phases: For each class j , we first evaluate the average precision of N detectors and rank them according to its precision in a network ranking phase, and then collect outputs of the top ranked detectors to generate \mathbb{D}^j in an inference phase.

4.2. Network Ranking

In order to rank detectors, we can exploit K -fold cross validation. We divide a training set into K subsets. Then, we train a detector T_i with $K - 1$ sets and evaluate its average precision $AP^k(T_i^j)$ for object class j with the held-out set at round k . An average precision matrix for N detectors for C object classes is denoted as

$$R = [r_{ji}]_{C \times N}, r_{ji} = \frac{1}{K} \sum_{k=1}^K AP^k(T_i^j), \quad (2)$$

Then, we generate a ranking matrix $R^* = [r_{ji}^*]$ by sorting each of its rows in descending order, where $r_{ji}^* \geq r_{j(i+1)}^*$, $i = 1, \dots, N - 1$. Here, it is worthy of notice that we rank the detectors for each class with the AP measure

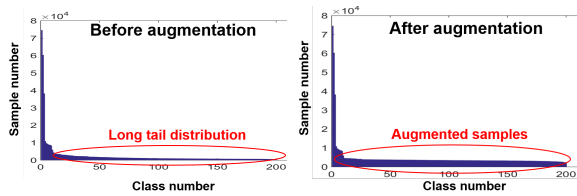


Figure 3. Class distribution before/after augmentation on the ILSVRC 2017 training dataset

instead of using the mAP for all classes. This makes detection results for each class decoupled. As a result, we can select and combine best experts for each class regardless of the AP rates for other classes, and improve the class-wise performance directly.

4.3. Ensemble Inference

Once a ranking matrix is generated, we select detectors with high APs among trained detectors for each class. Let denote the number of selected ones for object class j as N^j . To determine N^j for class j , we evaluate a threshold V^j by subtracting δ from r_{j1}^* (*i.e.* the maximum AP of N detectors’ APs for class j), and then select T_i^* when $r_{ji}^* \geq V^j$.

For prediction, we stack the outputs of selected detectors as $\mathbb{D}^j = \bigcup_{i=1}^{N^j} D_i^j$, where D_i^j and \mathbb{D}^j are the detection boxes of T_i^* and an union set of detections for class j . We then apply soft non-maximum suppression (Soft-NMS) [2] for \mathbb{D}^j , which decays s_k of an overlapped box \mathbf{d}_k with a box having a maximum score rather than removing \mathbf{d}_k , because it shows the better AP scores than conventional NMS for several benchmark datasets. After the Soft-NMS, we use the remaining boxes as final detections of the object.

Note that before generating \mathbb{D}^j applying the Soft-NMS for \mathbb{D}_i^j could reduce the detection performance. The reason is that some false positive samples could not be suppressed since the detections of other detectors, which can remove the false ones, are eliminated by the previous Soft-NMS.

5. Experiments

5.1. Implementation Details

Data augmentation: For the ILSVRC 2017 detection dataset, we augment the training set to mitigate a data imbalanced problem between object classes. As shown in Fig. 3, the training dataset has long tail distribution, and fine-tuning models with it may degrade performance [29].

The additional benefit is that we can generate different detectors by training them with/without the augmented samples as in Table 5. As a result, we can increase the diversity without using bagging which reduces a training set.

For data augmentation, we use 2D rotation, 2D translation, color histogram normalization, and noise addition. In 2D rotation, a whole image is rotated, but we restrict the

Table 1. The comparison of GPU memory usage with (w/t) and without (w/o) Caffe memory optimization. We test it on a IBM Minsky sever with 16280 (Mib) GPU memory.

Network	w/o optimization (Mib)	w/t optimization (Mib)
Res101-FRCN	9731	5969
Res152-FRCN	12889	7771
Res269-FRCN	Out of memory	13581

Table 2. Comparison results of the Res-FRCN models with different FG and TR on PASCAL 2007 test set.

Network	FG	TR	mAP	Network	FG	TR	mAP
Res101-FRCN	0.25	0	76.31%	Res101-FRCN	0.25	16	76.75%
Res101-FRCN	0.5	0	76.82%	Res101-FRCN	0.5	16	76.52%
Res152-FRCN	0.25	0	75.71%	Res152-FRCN	0.25	16	76.76%
Res152-FRCN	0.5	0	76.92%	Res152-FRCN	0.5	16	75.99%
Res269-FRCN	0.25	0	77.31%	Res269-FRCN	0.25	16	76.75%
Res269-FRCN	0.5	0	77.40%	Res269-FRCN	0.5	16	76.90%

rotation degree in a small range (-6° , 6°) to maintain box annotation quality. For 2D translation, we crop a (randomly selected) border of each image. In FRCN the size of the feature map (*e.g.* the last block layer of the “conv4” block in ResNet to be extracted by RoIPooling) is 16 times smaller (by max pooling and 2 stride in “conv1/3/4” block) than the input image size. For preventing the large translation of the feature map, we therefore cut $\{4, 8, 12\}$ pixels in the border after resizing. In the histogram normalization, histogram equalization, color casting, and color balancing are selectively applied. We also generate a noise image with additive white Gaussian noise with variance $[0.1, 0.001]$.

We randomly select and apply one of those methods for each original image. The number of augmented images from each original image varies 0 to 9. More samples are augmented when the sample number of the class is smaller. In our experiment, we generate 713,294 images with the 333,474 original images. Figure 3 shows the change of class distribution via data augmentation.

Multi-scale testing: For more accurate detection, we also use multi-scale testing as used in [18, 8, 12]. We apply the multi-scale testing only for the FRCN models and use four image scales $s \in \{400, 600, 800, 900\}$, where the s is the length of an image’s shorter side. We cap the longer size l of an image’s side at 1500 pixels for fitting the Res269-FRCN in our GPU memory during testing, and keep the image’s aspect ratio when resizing. However, we have not performed multi-scale training (*c.f.* [18, 8]). All the FRCN models are trained with $s = 600$ and $l = 1000$. Table 5 shows the detection results of the trained detectors at each scale. Since we trained those detectors with $s = 600$, they have the best rates at $s = 600$. We use detection results of the models for all the scales for ensemble detection.

Batch normalization and memory optimization: We use optimized Caffe [43] for improving the speed of batch

normalization. We have inserted the batch normalization layer in the original FRCN code [34]. For reducing GPU memory usage during network training and testing, we also apply the memory optimization code [43] to [34]. It can reduce the memory usage by reusing the memory allocated by previous blobs in a network. Table 1 compares the memory usages when training different Res-FRCN models with a mini-batch with 128 samples.

5.2. Performance Evaluation

Experimental setup: We compare our Rank of Experts method with the modern detection and ensemble methods. For evaluation, we use PASCAL VOC (2007&2012) [5] and ILSVRC datasets [35]. For more evaluation of our Rank of Experts, we participated in ILSVRC 2017 detection competition. We use the common metrics consisting of AP (averaged over IoU thresholds) and mAP (averaged over classes). When evaluating our methods on each set, we do not use other additional dataset (e.g. MSCOCO) for training.

Training parameters: All the parameters of implemented detectors have been determined experimentally and remained unchanged for each dataset evaluation. The most parameters are same to those given in official codes of FRCN [34] and SSD [27]. From extensive evaluation, we find that most parameters do not affect the overall performance significantly. We set the learning rates to 10^{-3} and 10^{-4} for PASCAL VOC and ILSVRC datasets, respectively. However, we found that the fraction FG of labeled foreground samples in a mini-batch and threshold TR , which limits the sizes of proposal boxes, can affect the performance. Therefore, both parameters are determined from evaluations as discussed in the next section. Since the explicit validation sets are available in both datasets, we evaluate the ranking matrix with the provided validation sets that are not exploited for training a detector.

Implemented detectors: As shown in Table 5, we implemented 19 detectors with different feature extractors, meta-architectures, and training sets. All the details of those can be found in Sec. 3. In Table 3-6, the performance of our detectors are highlighted.

PASCAL VOC evaluation: For VOC 2007 and 2012 (VOC07/12) tests, we train detectors with VOC07trainval (5011 images) and VOC07/VOC12trainval (21503 images) sets as also done in [8, 34, 27, 33]. We train detectors with $\mu = 10^{-3}$ learning rate for 50k iterations, and continue the training for 30k iteration with $\mu = 10^{-4}$ for VOC07. For VOC12, we train the detectors during 70k and 50k iterations with the same learning rates.

In Table 3, we compare the implemented detectors with the modern convolutional detectors. For the VOC07 test set, all the implemented detectors show the better results than the original ones. Main reason is that we more successfully train detectors by using feature aggregation in Sec.

Table 3. Performance comparison with other detectors on the PASCAL VOC07/12 test sets. The more results can be found in our supplementary material and [the PASCAL VOC 2012 website](#).

Contributors	Feature extractor	Meta architecture	RoI warping	mAP	Speed (fps)
Training set: VOC07 trainval / Test set: VOC07 test					
[8]	VGG16	Fast	Pooling	66.9%	-
[34]	VGG16	FRCN	Pooling	69.9%	-
[27]	VGG16	SSD300	-	68.0%	-
[27]	VGG16	SSD512	-	71.6%	-
[33]	VGG16	YOLO	-	57.9%	-
Ours	Res269	FRCN-Type1	Pooling	77.4%	3.0
Ours	Res101	FRCN-Type2	Alignment	75.8%	6.5
Ours	Res152	FRCN-Type1	Alignment	76.1%	5.2
Ours	Res152	FRCN-Type2	Alignment	76.2%	5.6
Ours	WRI	SSD	-	77.0%	12.0
Ours	VGG	SSD	-	74.2%	14.8
Ours	VGG	DSSD	-	74.5%	7.6
Rank of Experts (for 7 detectors)				81.2%	0.87
Training set: VOC07/12trainval / Test set: VOC12 test					
[7]	Res101	DSSD513	-	80.0%	-
[23]	VGG16	RUN-3WAY	-	79.8%	-
[27]	VGG16	SSD512	-	78.5%	-
[34]	Res101	FRCN	Pooling	76.8%	-
Ours	Res269	FRCN-Type1	Pooling	79.3%	3.0
Ours	Res101	FRCN-Type2	Alignment	76.9%	6.5
Ours	Res152	FRCN-Type1	Alignment	77.9%	5.2
Ours	Res152	FRCN-Type2	Alignment	77.2%	5.6
Ours	WRI	SSD	-	78.4%	12.0
Ours	VGG	SSD	-	76.6%	14.8
Ours	VGG	DSSD	-	77.6%	7.6
Rank of Experts (for 7 detectors)				82.2%	0.87

Table 4. ILSVRC 2017 competition results. More detection results can be found [in the ILSVRC 2017 website](#).

Team name	Categories Won	mAP
BDAT	85	73.13%
NUS-Qihoo-DPNs	9	65.69%
KAISTNIA-ETRI	0	61.02%
QINIUI-ATLAB+CAS-SAR	0	52.66%
FACEALL-BUPT	0	49.58%
PaulForDream	0	49.42%
DeepView (Rank of Experts)	10	59.30%

3.2 and tuning important parameters well as shown in Table 2. In particular, we obtain the best detection rates using the Res269-FRCN models. Exploiting feature extractors with low classification errors also improves the detection rates when comparing Res101/Res152/Res269-FRCN models. In addition, SSD and DSSD detectors show the comparable performance with the FRCN detectors. We can improve the detection rate about 4% using ensemble of 7 detectors.

We further evaluate our methods in the PASCAL VOC 2012 challenge. In this evaluation, we only use the VOC dataset (without using COCO dataset) for a fair comparison. The comparison results are given in Table 3. Compared to the state-of-the-art detectors, the implemented detectors show the competitive performance, and we achieve the best mAP by using the proposed ensemble detection.

Table 5. Evaluation results of different single and ensemble detectors on the ILSVRC 2017 val2 set. For each image resolution, **the best detection rates are marked with red color** and **the number of the most selected detectors are marked with blue color**.

Network	Feature extractor	Meta architecture	RoI warping	Training Dataset	mAP (# of selection) per image resolution				
					400	512	600	800	900
D1	Res101	FRCN-Type1	Pooling	train	50.07% (0)	-	53.13% (13)	53.25% (16)	52.08% (10)
D2	Res101	FRCN-Type1	Pooling	train+val1	49.79% (2)	-	53.57% (12)	53.32% (7)	51.96% (5)
D3	Res152	FRCN-Type1	Pooling	train+val1	52.16% (11)	-	55.77% (32)	54.94% (20)	53.59% (19)
D4	Res152	FRCN-Type1	Pooling	train	52.59% (13)	-	55.71% (23)	55.27% (26)	-
D5	Res152	FRCN-Type1	Pooling	train+val1	51.41% (6)	-	55.35% (32)	54.21% (19)	52.91% (11)
D6	Res152	FRCN-Type1	Pooling	train	51.92% (5)	-	56.00% (28)	55.41% (18)	54.38% (24)
D7	Res152	FRCN-Type1	Pooling	train+val1	52.41% (8)	-	56.19% (32)	55.54% (25)	54.34% (22)
D8	Res269	FRCN-Type1	Pooling	train+val1	-	-	56.92% (49)	-	-
D9	Res269	FRCN-Type1	Pooling	train+val1	54.09% (24)	-	57.65% (69)	56.29% (38)	54.94% (19)
D10	Res269	FRCN-Type1	Pooling	trainval+val1+aug	53.21% (17)	-	56.76% (43)	55.84% (34)	54.49% (20)
D11	Res269	FRCN-Type1	Pooling	trainval1	53.98% (23)	-	57.72% (76)	56.64% (45)	55.41% (26)
D12	Res269	FRCN-Type1	Pooling	train	53.59% (17)	-	57.34% (63)	56.56% (33)	55.53% (29)
D13	Res152	FRCN-Type2	Alignment	train+val1	48.91% (6)	-	54.65% (46)	54.53% (30)	54.49% (42)
D14	Res152	FRCN-Type2	Alignment	trainval+val1+aug	-	-	54.57% (40)	53.92% (29)	-
D15	Res152	FRCN-Type2	Alignment	train+val1	51.95% (8)	-	56.15% (35)	55.41% (25)	54.73% (20)
D16	Res152	FRCN-Type2	Alignment	trainval+val1+aug	43.42% (1)	-	51.39% (10)	49.00% (5)	47.40% (2)
D17	VGG	SSD	-	trainval+val1+aug	-	50.48% (14)	-	-	-
D18	VGG	DSSD	-	trainval+val1+aug	-	49.98% (15)	-	-	-
D19	WRI	SSD	-	trainval+val1+aug	-	49.21% (8)	-	-	-
Rank of Experts (for 19 detectors)					62.54%				

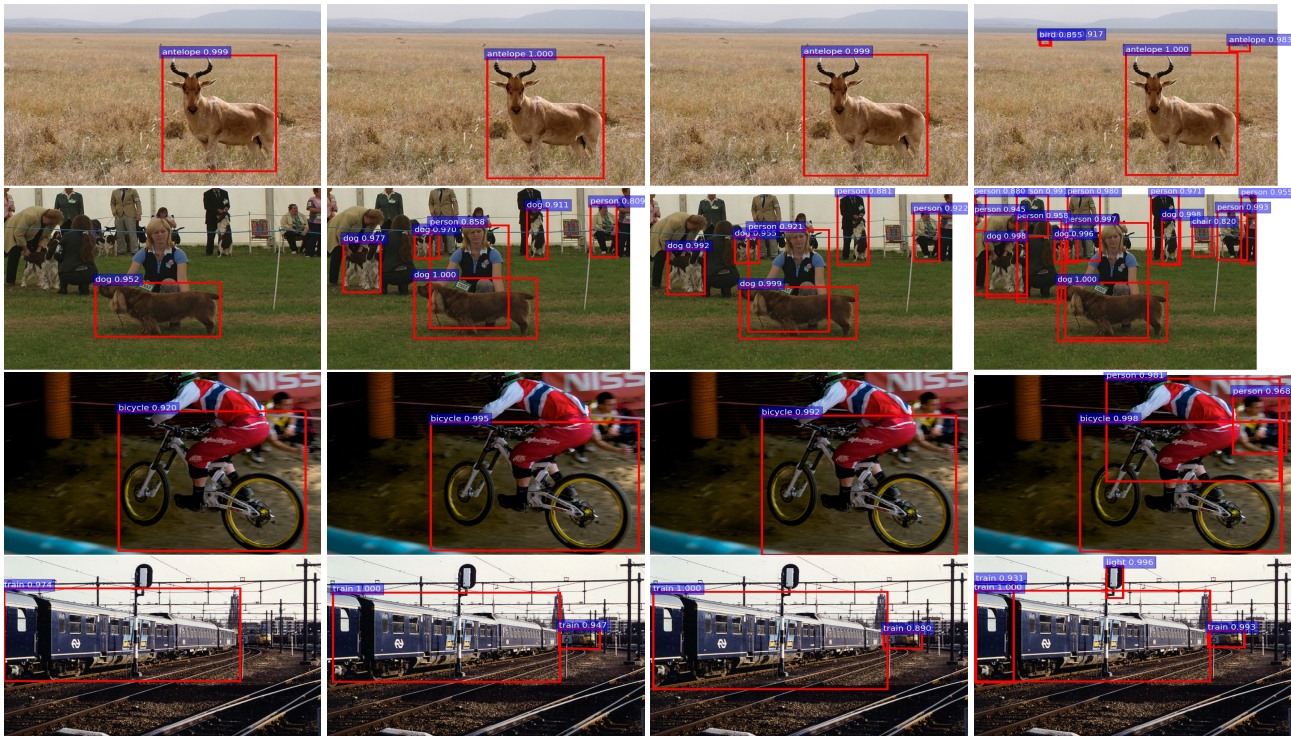


Figure 4. Detection results of VGG-DSSD, Res152-FRCN, Res269-FRCN and Rank of Experts are shown in each column.

ILSVRC 2017 competition: For this competition, we implement 19 detectors with different structures, training sets and initial states as shown in Table 5. We train the detectors with the $\mu = 10^{-4}$ and $\mu = 10^{-5}$ for 1M and 400k iterations. We also provide the performance of each detector. Among them D11 using RES269-FRCN produces the best mAP. Moreover, the mAP difference between detectors is much bigger than on PACAL VOC evaluation. The more

comparison results are given in Fig. 4. Our Rank of Experts shows the better results than other single detectors.

Table 4 shows the comparison results with detectors of other teams in ILSVRC 2017 competition. Although three teams show the better mAPs than ours, we won the 2nd place for object category won. These results prove that the proposed Rank of Expert can improve class-wise detection rates significantly.

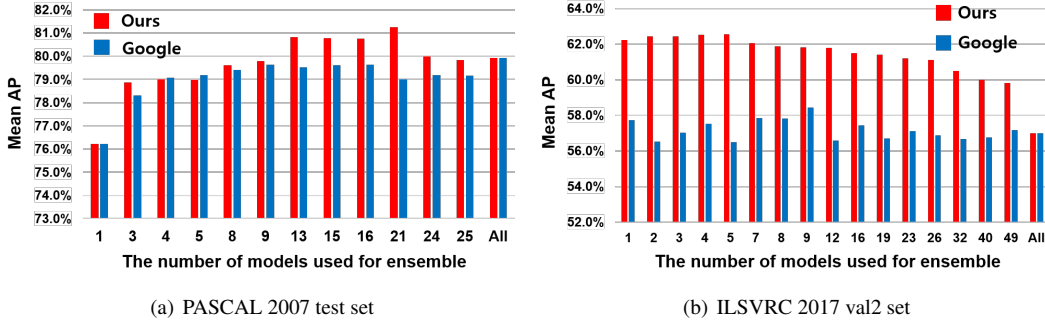


Figure 5. Performance comparison with other ensemble method [14] on PASCAL and ILSVRC datasets.

Table 6. Mean AP(mAP) improvement on ILSVRC 2017 val2 set.

Method	mAP	Method	mAP
Soft-NMS [2]	~ 1%	Data augmentation	1 ~ 2%
Multi-scale test	~ 1%	Rank of experts	4 ~ 5 %

Ablation experiment: In Table 6, we show the performance improvement when applying each method. We evaluate the averaged mAP improvement of D10, D14 and D16-19 detectors. We found that the mAP is greatly enhanced by using the proposed Rank of Experts.

Comparison with other ensemble method: To show effectiveness of our Rank of Experts, we compare our method with the ensemble method of Google [14] which won 1st places on the MSCOCO 2016 competition. To this end, we evaluate a similarity matrix by computing cosine distances between the AP score vectors of all the detectors. We then greedily select the best models with highest mAP, but prune away a model when the similarity score between selected ones and the candidate exceeds to a threshold. By changing the threshold, we can vary the number of detectors used for ensemble.

Figure 5 shows the comparison results of both ensemble methods on PASCAL VOC 2007 and ILSVRC 2017 datasets. In the most cases, we achieve the better mAPs on both datasets. On average the proposed method improves the mAPs to 0.31% and 4.15% for PASCAL and ILSVRC datasets, respectively. This indicates that our method is more effective especially for large scale object detection.

Discussion: For our Rank of Experts, the most important parameter is $\delta = [0, 1]$ that determines N_j of how many detectors will be used for each class detection. If $\delta = 0$, we can use one detector with the best mAP only, but use all the detectors if $\delta = 1$. The more detectors can be used for ensemble when δ increases.

To show the effects of δ , we change the δ from 0 to 0.1 with interval 0.05 on PASCAL and ILSVRC datasets. For this evaluation, 7 detectors (PASCAL) and 19 detectors (ILSVRC) with multi-scale testing as described in Table 3 and 5 were used. We obtain the best mAP 81.24% and 62.54% with $\delta = 0.03$ for each dataset. However, the num-

ber of detectors used for ensemble on each dataset is different. More detectors were used on PASCAL evaluation since the mAP difference of detectors is not significant. Note that the performance difference of ensemble detection using different δ is negligible. There are only ~ 1% and ~ 2% mAP differences for the PASCAL and ILSVRC datasets. Thus, we verify that δ does not affects the overall performance significantly and the robustness of our Rank of Experts.

In Table 5, we provide the number of selected detectors for ensemble. In general, detectors with high mAPs are selected more frequently. However, all the detectors can be used for ensemble detection by using the proposed Rank of Experts although some detectors (*i.e.* WRI-SSD and VGG-DSSD) have much lower mAPs than other detectors.

Running time: The speed of each detector depends on types of feature extractors and meta-architectures. We present the speed of our implemented detectors in Table 3. The running time of our Rank of Experts relies on the number of detection boxes. For 162K and 3.66M boxes¹ on PASCAL and ILSVRC dataset evaluation, the speeds are ~ 140(fps) and ~ 9.37(fps), respectively. The main bottleneck occurs during Soft-NMS. The speed can be greatly improved with the GPU programming.

6. Conclusion

In the recent object detection challenges, the top-ranked teams significantly improve the detection rates using ensemble learning. However, the conventional methods requires additional training stages with huge GPU memory usage. In this paper, we therefore propose an effective ensemble detection for combining different detectors without the expensive training. Furthermore, the additional benefit of our ensemble detection can improve the class-wise detection rate directly. In order to increase the diversity between detectors, we implement a variety of detectors with different feature extractors, meta-architectures and training sets. We have proved that the proposed method leads to the obvious performance improvement through extensive evaluations.

¹For PASCAL dataset, 27 (models) \times 20 (classes) \times 300 (maximum number of boxes per class). For ILSVRC dataset, 61 \times 200 \times 300.

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org. **3**
- [2] N. Bodla, B. Singh, R. Chellappa, and L. S. Davis. Improving object detection with one line of code. *CoRR*, abs/1704.04503, 2017. **5, 8**
- [3] J. Dai, Y. Li, K. He, and J. Sun. R-FCN: object detection via region-based fully convolutional networks. *Neural Information Processing Systems*. **1, 2**
- [4] M. R. David Eigen and I. Sutskever. Learning factored representations in a deep mixture of experts. *ICLR Workshops*, 2014. **3**
- [5] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>. **6**
- [6] Y. Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121(2):256 – 285, 1995. **2**
- [7] C. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg. DSSD : Deconvolutional single shot detector. *CoRR*, abs/1701.06659, 2017. **1, 2, 3, 4, 6**
- [8] R. B. Girshick. Fast r-cnn. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 1440–1448, 2015. **1, 2, 3, 4, 5, 6**
- [9] R. B. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, pages 580–587, 2014. **2**
- [10] L. K. Hansen and P. Salamon. Neural network ensembles. *IEEE Trans. Pattern Anal. Mach. Intell.*, 12(10):993–1001, 1990. **2**
- [11] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick. Mask R-CNN. *CoRR*, abs/1703.06870, 2017. **3, 4**
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778, 2016. **2, 3, 5**
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV*, pages 630–645, 2016. **1, 3**
- [14] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, and K. Murphy. Speed/accuracy trade-offs for modern convolutional object detectors. *Computer Vision and Pattern Recognition*, abs/1611.10012, 2017. **2, 3, 4, 8**
- [15] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton. Adaptive mixtures of local experts. *Neural Computation*, 3(1):79–87, 1991. **2, 3, 4**
- [16] Y. X. Y. L. G. Z. H. H. Y. W. Jifeng Dai, Haozhi Qi. Deformable convolutional networks. *International Conference on Computer Vision*, 2017. **2**
- [17] M. I. Jordan and R. A. Jacobs. Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, 6(2):181–214, 1994. **2, 3, 4**
- [18] H. Kaiming, Z. Xiangyu, R. Shaoqing, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In *European Conference on Computer Vision*, 2014. **2, 5**
- [19] T. Kong, A. Yao, Y. Chen, and F. Sun. Hypernet: Towards accurate region proposal generation and joint object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 845–853, 2016. **3, 4**
- [20] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States.*, pages 1106–1114, 2012. **2**
- [21] A. Krogh and J. Vedelsby. Neural network ensembles, cross validation, and active learning. In *Advances in Neural Information Processing Systems 7, [NIPS Conference, Denver, Colorado, USA, 1994]*, pages 231–238, 1994. **2, 3, 4**
- [22] T. Lasota, B. Londzin, Z. Telec, and B. Trawinski. Comparison of ensemble approaches: Mixture of experts and adaboost for a regression problem. In *Intelligent Information and Database Systems - 6th Asian Conference, ACHIDS 2014, Bangkok, Thailand, April 7-9, 2014, Proceedings, Part II*, pages 100–109, 2014. **3, 4**
- [23] K. Lee, J. Choi, J. Jeong, and N. Kwak. Residual features and unified prediction network for single stage detection. *CoRR*, abs/1707.05031, 2017. **6**
- [24] Y. Lee, H. Kim, E. Park, X. Cui, and H. Kim. Wide-residual-inception networks for real-time object detection. *CoRR*, abs/1702.01243, 2017. **2, 3**
- [25] T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie. Feature pyramid networks for object detection. *Computer Vision and Pattern Recognition*. **2**
- [26] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. **2**
- [27] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg. SSD: single shot multibox detector. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I*, pages 21–37, 2016. **1, 2, 3, 4, 6**

- [28] D. Opitz and R. Maclin. Popular ensemble methods: An empirical study. *Journal of Artificial Intelligence Research*, 11:169–198, 1999. [2](#)
- [29] W. Ouyang, X. Wang, C. Zhang, and X. Yang. Factors in finetuning deep model for object detection with long-tail distribution. pages 864–873, 06 2016. [5](#)
- [30] E. Park, X. Han, T. L. Berg, and A. C. Berg. Combining multiple sources of knowledge in deep cnns for action recognition. In *2016 IEEE Winter Conference on Applications of Computer Vision, WACV 2016, Lake Placid, NY, USA, March 7-10, 2016*, pages 1–8, 2016. [3](#), [4](#)
- [31] M. P. Perrone and L. N. Cooper. When networks disagree: Ensemble methods for hybrid neural networks. pages 126–142. Chapman and Hall, 1993. [2](#), [3](#)
- [32] R. Polikar. Ensemble based systems in decision making. *IEEE Circuits and Systems Magazine*, 6(3):21–45, 2006. [2](#), [4](#)
- [33] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 779–788, 2016. [2](#), [6](#)
- [34] S. Ren, K. He, R. B. Girshick, and J. Sun. Faster R-CNN: towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 91–99, 2015. [1](#), [2](#), [3](#), [6](#)
- [35] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. [1](#), [2](#), [3](#), [6](#)
- [36] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *CoRR*, abs/1312.6229, 2013. [1](#)
- [37] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *CoRR*, 2017. [3](#)
- [38] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. [1](#), [2](#), [3](#)
- [39] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA.*, pages 4278–4284, 2017. [1](#)
- [40] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Computer Vision and Pattern Recognition (CVPR)*, 2015. [3](#)
- [41] C. Szegedy, A. Toshev, and D. Erhan. Deep neural networks for object detection. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 2553–2561, 2013. [1](#)
- [42] J. R. R. Uijlings, K. E. A. van de Sande, T. Gevers, and A. W. M. Smeulders. Selective search for object recognition. *International Journal of Computer Vision*, 104(2):154–171, 2013. [2](#)
- [43] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Val Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. [5](#), [6](#)
- [44] X. Zeng, W. Ouyang, J. Yan, H. Li, T. Xiao, K. Wang, Y. Liu, Y. Zhou, B. Yang, Z. Wang, H. Zhou, and X. Wang. Crafting gbd-net for object detection. *CoRR*, abs/1610.02579, 2016. [3](#)
- [45] Z. Zhou, J. Wu, and W. Tang. Ensembling neural networks: Many could be better than all. *Artif. Intell.*, 137(1-2):239–263, 2002. [3](#), [4](#)