

DSOD: Learning Deeply Supervised Object Detectors from Scratch

Zhiqiang Shen^{*1}, Zhuang Liu^{*2}, Jianguo Li³, Yu-Gang Jiang¹, Yurong Chen³, Xiangyang Xue¹
¹Fudan University, ²Tsinghua University, ³Intel Labs China

{zhiqiangshen13, ygj, yxxue}@fudan.edu.cn, liuzhuangthu@gmail.com

{jianguo.li, yurong.chen}@intel.com

Abstract

We present *Deeply Supervised Object Detector (DSOD)*, a framework that can learn object detectors from scratch. State-of-the-art object detectors rely heavily on the off-the-shelf networks pre-trained on large-scale classification datasets like ImageNet, which incurs learning bias due to the difference on both the loss functions and the category distributions between classification and detection tasks. Model fine-tuning for the detection task could alleviate this bias to some extent but not fundamentally. Besides, transferring pre-trained models from classification to detection between discrepant domains is even more difficult (e.g. RGB to depth images). A better solution to tackle these two critical problems is to train object detectors from scratch, which motivates our proposed DSOD. Previous efforts in this direction mostly failed due to much more complicated loss functions and limited training data in object detection. In DSOD, we contribute a set of design principles for training object detectors from scratch. One of the key findings is that deep supervision, enabled by dense layer-wise connections, plays a critical role in learning a good detector. Combining with several other principles, we develop DSOD following the single-shot detection (SSD) framework. Experiments on PASCAL VOC 2007, 2012 and MS COCO datasets demonstrate that DSOD can achieve better results than the state-of-the-art solutions with much more compact models. For instance, DSOD outperforms SSD on all three benchmarks with real-time detection speed, while requires only 1/2 parameters to SSD and 1/10 parameters to Faster RCNN.

1. Introduction

Convolutional Neural Networks (CNNs) have produced impressive performance improvements in many computer vision tasks, such as image classification [17, 28, 32, 9, 10], object detection [5, 4, 27, 19, 21, 25], image segmenta-

tion [23, 8, 2, 36], etc. In the past several years, many innovative CNN network structures have been proposed. Szegedy *et al.* [32] propose an “Inception” module which concatenates features maps produced by various sized filters. He *et al.* [9] propose residual learning blocks with skip connections, which enable training very deep networks with more than 100 layers. Huang *et al.* [10] propose DenseNets with dense layer-wise connections. Thanks to these excellent network structures, the accuracy of many vision tasks has been greatly improved. Among them, object detection is one of the fastest moving areas due to its wide applications in surveillance, autonomous driving, etc.

In order to achieve good performance, most of the advanced object detection systems fine-tune networks pre-trained on ImageNet [3]. This fine-tuning process is also viewed as transfer learning [24]. Fine-tuning from pre-trained models has at least two advantages. First, there are many state-of-the-art deep models publicly available. It is very convenient to reuse them for object detection. Second, fine-tuning can quickly generate the final model and requires much less instance-level annotated training data than the classification task.

However, there are also critical limitations when adopting the pre-trained networks in object detection: (1) *Limited structure design space*. The pre-trained network models are mostly from ImageNet-based classification task, which are usually very heavy — containing a huge number of parameters. Existing object detectors directly adopt the pre-trained networks, and as a result there is little flexibility to control/adjust the network structures (even for small changes of network structure). The requirement of computing resources is also bounded by the heavy network structures. (2) *Learning bias*. As both the loss functions and the category distributions between classification and detection tasks are different, we argue that this will lead to different searching/optimization spaces. Therefore, learning may be biased towards a local minimum which is not the best for detection task. (3) *Domain mismatch*. As is known, fine-tuning can mitigate the gap due to different target category distribution. However, it is still a severe problem when the source

^{*}indicates equal contribution. This work was done when Zhiqiang Shen and Zhuang Liu were interns at Intel Labs China. Jianguo Li is the corresponding author.

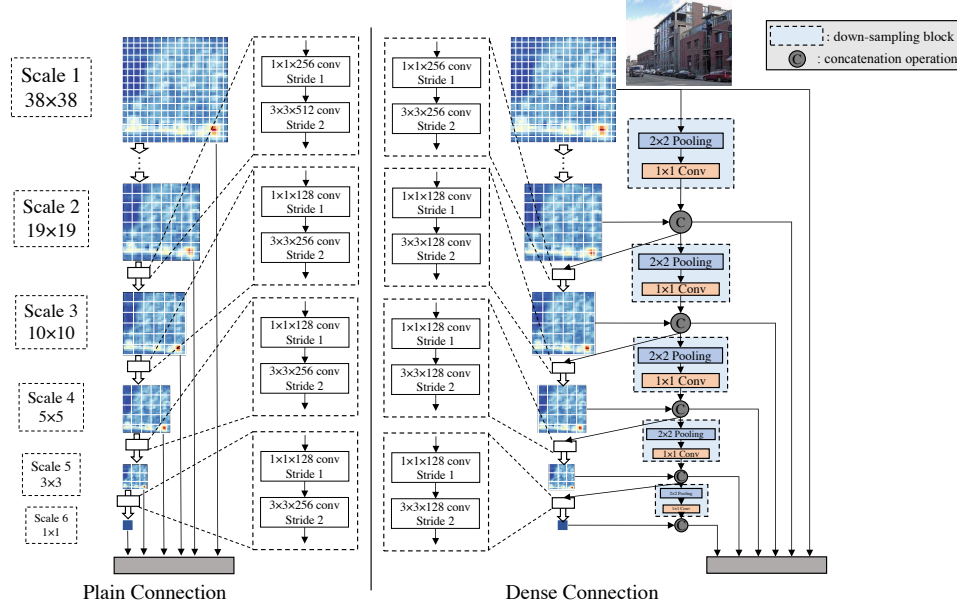


Figure 1: DSOD prediction layers with plain and dense structures (300×300 input). Plain structure is introduced by SSD [21]. See Sec. 3 for more details.

domain (ImageNet) has a huge mismatch to the target domain such as depth images, medical images, etc [7].

Our work is motivated by the following two questions. First, is it possible to train object detection networks **from scratch**? Second, if the first answer is positive, are there any principles to design a resource efficient network structure for object detection while keeping high detection accuracy? To meet this goal, we propose deeply supervised object detectors (DSOD), a simple yet efficient framework which could learn object detectors from scratch. DSOD is fairly flexible, so that we can tailor various network structures for different computing platforms such as server, desktop, mobile and even embedded devices.

We contribute a set of principles for designing DSOD. One key point is that *deep supervision* plays a critical role, which is motivated by the work of [18, 35]. In [35], Xie *et al.* proposed a holistically-nested structure for edge detection, which included the side-output layers to each conv-stage of base network for explicit deep supervision. Instead of using the multiple cut-in loss signals with side-output layers, this paper adopts *deep supervision* implicitly through the dense layer-wise connections proposed in DenseNet [10]. Dense structures are not only adopted in the backbone sub-network, but also in the front-end multi-scale prediction layers. Figure 1 illustrates the structure comparison in front-end prediction layers. The fusion and reuse of multi-resolution prediction-maps help keep or even improve the final accuracy while reducing model parameters to some extent.

Our main contributions are summarized as follows:

- (1) We present DSOD, to the best of our knowledge, the first framework that can train object detection networks from scratch with state-of-the-art performance,

even with limited training data.

- (2) We introduce and validate a set of principles to design efficient object detection networks from scratch through step-by-step ablation studies.
- (3) We show that our DSOD can achieve state-of-the-art performance on three standard benchmarks (PASCAL VOC 2007, 2012 and MS COCO datasets) with real-time processing speed and more compact models.

2. Related Work

Object Detection. State-of-the-art CNN based object detection methods can be divided into two groups: (i) **region proposal based methods** and (ii) **proposal-free methods**.

Proposal based methods include R-CNN [5], Fast R-CNN [4], Faster R-CNN [27] and R-FCN [19]. R-CNN uses **selective search** [34] to first generate potential object regions in an image and then perform classification on the proposed regions. R-CNN requires high computational costs since each region is processed by the CNN network separately. Fast R-CNN and Faster R-CNN improve the efficiency by sharing computation and using neural networks to generate the region proposals. R-FCN further improves speed and accuracy by removing fully-connected layers and adopting position-sensitive score maps for final detection.

Proposal-free methods like YOLO [25] and SSD [21] have recently been proposed for real-time detection. YOLO uses a single feed-forward convolutional network to directly predict object classes and locations. Comparing with the region-based methods, YOLO no longer requires a second per-region classification operation so that it is extremely fast. SSD improves YOLO in several aspects, including (1) using small convolutional filters to predict categories and

Layers		Output Size (Input $3 \times 300 \times 300$)	DSOD
Stem	Convolution	$64 \times 150 \times 150$	3×3 conv, stride 2
	Convolution	$64 \times 150 \times 150$	3×3 conv, stride 1
	Convolution	$128 \times 150 \times 150$	3×3 conv, stride 1
	Pooling	$128 \times 75 \times 75$	2×2 max pool, stride 2
Dense Block (1)		$416 \times 75 \times 75$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$
Transition Layer (1)		$416 \times 75 \times 75$	1×1 conv
		$416 \times 38 \times 38$	2×2 max pool, stride 2
Dense Block (2)		$800 \times 38 \times 38$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 8$
Transition Layer (2)		$800 \times 38 \times 38$	1×1 conv
		$800 \times 19 \times 19$	2×2 max pool, stride 2
Dense Block (3)		$1184 \times 19 \times 19$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 8$
Transition w/o Pooling Layer (1)		$1120 \times 19 \times 19$	1×1 conv
Dense Block (4)		$1568 \times 19 \times 19$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 8$
Transition w/o Pooling Layer (2)		$1568 \times 19 \times 19$	1×1 conv
DSOD Prediction Layers		–	Plain/Dense

Table 1: DSOD architecture (growth rate $k = 48$ in each dense block).

anchor offsets for bounding box locations; (2) using pyramid features for prediction at different scales; and (3) using default boxes and aspect ratios for adjusting varying object shapes. Our proposed DSOD is built upon the SSD framework and thus it inherits the speed and accuracy advantages of SSD, while produces smaller and more flexible models.

Network Architectures for Detection. Significant efforts have been devoted to the design of network architectures for image classification. Many different networks are emerged, such as AlexNet [17], VGGNet [28], GoogLeNet [32], ResNet [9] and DenseNet [10]. Meanwhile, several regularization techniques [29, 12] have also been proposed to further enhance the model capabilities. Most detection methods [5, 4, 27, 21] directly utilize pre-trained ImageNet models as the backbone network.

Some other works design specific backbone network structures for object detection, but still require pre-training the network on ImageNet classification dataset first. For instance, YOLO [25] defines a network with 24 convolutional layers followed by 2 fully connected layers. YOLO9000 [26] improves YOLO by proposing a new network named Darknet-19, which is a simplified version of VGGNet [28]. Kim *et al.* [15] proposes PVANet for object detection, which consists of the simplified “Inception” block from GoogLeNet. Huang *et al.* [11] investigated various combination of network structures and detection frameworks, and found that Faster R-CNN [27] with Inception-ResNet-v2 [31] achieved the highest performance. In this paper, we also consider network structures for generic object detection. However, the pre-training on ImageNet is no longer required by the proposed DSOD.

Learning Deep Models from Scratch. To the best of our knowledge, there are no works which train object detection networks from scratch. The proposed approach has very

appealing advantages over existing solutions. We will elaborate and validate the method in the following sections. In semantic segmentation, Jégou *et al.* [13] demonstrated that a well-designed network structure can outperform state-of-the-art solutions without using the pre-trained models. It extends DenseNets to fully convolutional networks by adding an upsampling path to recover the full input resolution.

3. DSOD

In this section, we first introduce our DSOD architecture and its components, and elaborate several important design principles. Then we describe the training settings.

3.1. DSOD Architecture

Overall Framework. The proposed DSOD method is a multi-scale proposal-free detection framework similar to SSD [21]. The network structure of DSOD can be divided into two parts: the backbone sub-network for feature extraction and the front-end sub-network for prediction over multi-scale response maps. The backbone sub-network is a variant of the deeply supervised DenseNets [10] structure, which is composed of a *stem block*, four *dense blocks*, two *transition layers* and two *transition w/o pooling layers*. The front-end subnetwork (or named *DSOD prediction layers*) fuses multi-scale prediction responses with an elaborated *dense structure*. Figure 1 illustrates the proposed DSOD prediction layers along with the plain structure of multi-scale predicting maps as used in SSD [21]. The full DSOD network architecture¹ is detailed in Table 1. We elaborate each component and the corresponding design principle in

¹The visualization of the complete network structure is available at: <http://ethereon.github.io/netscope/#/gist/b17d01f3131e2a60f9057b5d3eb9e04d>.

the following.

Principle 1: Proposal-free. We investigated all the state-of-the-art CNN based object detectors, and found that they could be divided into three categories. First, R-CNN and Fast R-CNN require external object proposal generators like selective search. Second, Faster R-CNN and R-FCN require integrated region-proposal-network (RPN) to generate relatively fewer region proposals. Third, YOLO and SSD are single-shot and proposal-free methods, which handle object location and bounding box coordinates as a regression problem. We observe that only the proposal-free method (the 3rd category) can converge successfully without the pre-trained models. We conjecture this is due to the RoI (Regions of Interest) pooling in the other two categories of methods — RoI pooling generates features for each region proposals, which hinders the gradients being smoothly back-propagated from region-level to convolutional feature maps. The proposal-based methods work well with pre-trained network models because the parameter initialization is good for those layers before RoI pooling, while this is not true for training from scratch.

Hence, we arrive at the first principle: training detection network from scratch requires a proposal-free framework. In practice, we derive a multi-scale proposal-free framework from the SSD framework [21], as it could reach state-of-the-art accuracy while offering fast processing speed.

Principle 2: Deep Supervision. The effectiveness of deeply supervised learning has been demonstrated in GoogLeNet [32], DSN [18], DeepID3 [30], etc. The central idea is to provide integrated objective function as direct supervision to the earlier hidden layers, rather than only at the output layer. These “companion” or “auxiliary” objective functions at multiple hidden layers can mitigate the “vanishing” gradients problem. The proposal-free detection framework contains both classification loss and localization loss. The explicit solution requires adding complex side-output layers to introduce “companion” objective at each hidden layer for the detection task, similar to [35]. Here we empower deep supervision with an elegant & implicit solution called dense layer-wise connection, as introduced in DenseNets [10]. A block is called *dense block* when all preceding layers in the block are connected to the current layer. Hence, earlier layers in DenseNet can receive additional supervision from the objective function through the skip connections. Although only a single loss function is required on top of the network, all layers including the earlier layers still can share the supervised signals unencumbered. We will verify the benefit of deep supervision in Section 4.1.2.

Transition w/o Pooling Layer. We introduce this layer in order to increase the number of dense blocks without reducing the final feature map resolution. In the original design of DenseNet, each transition layer contains a pooling operation to down-sample the feature maps. The number of

dense blocks is fixed (4 dense blocks in all DenseNet architectures) if one wants to maintain the same scale of outputs. The only way to increase network depth is adding layers inside each block for the original DenseNet. The transition w/o pooling layer eliminates this restriction of the number of dense blocks in our DSOD architecture, and can also be used in the standard DenseNet.

Principle 3: Stem Block. Motivated by Inception-v3 [33] and v4 [31], we define stem block as a stack of three 3×3 convolution layers followed by a 2×2 max pooling layer. The first conv-layer works with stride = 2 and the other two are with stride = 1. We find that adding this simple stem structure can evidently improve the detection performance in our experiments. We conjecture that, compared with the original design in DenseNet (7×7 conv-layer, stride = 2 followed by a 3×3 max pooling, stride = 2), the stem block can reduce the information loss from raw input images. We will show that the reward of this stem block is significant for detection performance in Section 4.1.2.

Principle 4: Dense Prediction Structure. Figure 1 illustrates the comparison of the plain structure (as in SSD) and our proposed dense structure in the front-end sub-network. SSD designs prediction-layers as an asymmetric hourglass structure. For 300×300 input images, six scales of feature maps are applied for predicting objects. The Scale-1 feature maps are from the middle layer of the backbone sub-network, which has the largest resolution (38×38) in order to handle the small objects in an image. The remaining five scales are on top of the backbone sub-network. Then, a plain transition layer with the *bottleneck* structure (a 1×1 conv-layer for reducing the number of feature maps plus a 3×3 conv-layer) [33, 9] is adopted between two contiguous scales of feature maps.

Learning Half and Reusing Half. In the plain structure as in SSD (see Figure 1), each later scale is directly transited from the adjacent previous scale. We propose the dense structure for prediction, which fuses multi-scale information for each scale. For simplicity, we restrict that each scale outputs the same number of channels for the prediction feature maps. In DSOD, in each scale (except scale-1), half of the feature maps are learned from the previous scale with a series of conv-layers, while the remaining half feature maps are directly down-sampled from the contiguous high-resolution feature maps. The down-sampling block consists of a 2×2 , stride = 2 max pooling layer followed by a 1×1 , stride = 1 conv-layer. The pooling layer aims to match resolution to current size during concatenation. The 1×1 conv-layer is used to reduce the number of channels to 50%. The pooling layer is placed before the 1×1 conv-layer for the consideration of reducing computing cost. This down-sampling block actually brings each scale with the multi-resolution feature maps from all of its preceding scales, which is essentially identical to the dense

	DSOD300							
transition w/o pooling?	✓	✓	✓	✓	✓	✓	✓	✓
hi-comp factor θ ?		✓	✓	✓	✓	✓	✓	✓
wide bottleneck?			✓	✓	✓	✓	✓	✓
wide 1st conv-layer?				✓	✓	✓	✓	✓
big growth rate?					✓	✓	✓	✓
stem block?						✓	✓	✓
dense pred-layers?							✓	✓
VOC 2007 mAP	59.9	61.6	64.5	68.6	69.7	74.5	77.3	77.7

Table 2: **Effectiveness of various designs on VOC 2007 test set.** Please refer to Table 3 and Section 4.1 for more details.

layer-wise connection introduced in DenseNets. For each scale, we only learn half of new feature maps and reuse the remaining half of the previous ones. This dense prediction structure can yield more accurate results with fewer parameters than the plain structure, as will be studied in Section 4.1.

3.2. Training Settings

We implement our detector based on the Caffe framework [14]. All our models are trained from scratch with SGD solver on NVidia TitanX GPU. Since **each scale of DSOD feature maps is concatenated** from multiple resolutions, we adopt the L2 normalization technique [22] to scale the feature norm to 20 on all outputs. Note that SSD only applies this normalization to scale-1. Most of our training strategies follow SSD, including data augmentation, scale and aspect ratios for default boxes and loss function (*e.g.*, smooth L1 loss for localization purpose and softmax loss for classification purpose), while we have our own learning rate scheduling and mini-batch size settings. Details will be given in the experimental section.

4. Experiments

We conduct experiments on the widely used PASCAL VOC 2007, 2012 and MS COCO datasets that have 20, 20, 80 object categories respectively. Object detection performance is measured by mean Average Precision (mAP).

4.1. Ablation Study on PASCAL VOC2007

We first investigate each component and design principle of our DSOD framework. The results are mainly summarized in Table 2 and Table 3. We design several controlled experiments on PASCAL VOC 2007 with our DSOD300 (with 300×300 inputs) for this ablation study. A consistent setting is imposed on all the experiments, unless when some components or structures are examined. In this study, we train the models with the combined training set from VOC 2007 `trainval` and 2012 `trainval` (“07+12”), and test on the VOC 2007 testset.

4.1.1 Configurations in Dense Blocks

We first investigate the impact of different configurations in dense blocks of the backbone sub-network.

Compression Factor in Transition Layers. We compare two compression factor values ($\theta = 0.5, 1$) in the transition layers of DenseNets. Results are shown in Table 3 (rows 2 and 3). Compression factor $\theta = 1$ means that there is no feature map reduction in the transition layer, while $\theta = 0.5$ means half of the feature maps are reduced. Results show that $\theta = 1$ yields 2.9% higher mAP than $\theta = 0.5$.

Channels in bottleneck layers. As shown in Table 3 (rows 3 and 4), we observe that wider bottleneck layers (with more channels of response maps) improve the performance greatly (4.1% mAP).

Channels in the 1st conv-layer We observe that a large number of channels in the first conv-layers is beneficial, which brings 1.1% mAP improvement (in Table 3 rows 4 and 5).

Growth rate. A large growth rate k is found to be much better. We observe 4.8% mAP improvement in Table 3 (rows 5 and 6) when increase k from 16 to 48 with $4k$ bottleneck channels.

4.1.2 Effectiveness of Design Principles

We now justify the effectiveness of the key design principles elaborated earlier.

Proposal-free Framework. We tried to train object detectors from scratch using the proposal-based framework such as Faster R-CNN and R-FCN. However, the training process failed to converge for all the network structures we attempted (VGGNet, ResNet, DenseNet). We further tried to train object detectors using the proposal-free framework SSD. The training converged successfully but gives much worse results (69.6% for VGG), compared with the case fine-tuning from pre-trained model (75.8%), as shown in Table 4. This experiment validates our design principle to choose a proposal-free framework.

Deep Supervision. We then tried to train object detectors from scratch with deep supervision. Our DSOD300 achieves 77.7% mAP, which is much better than the SSD300S that is trained from scratch using VGG16 (69.6%) without deep supervision. It is also much better than the fine-tuned results by SSD300 (75.8%). This validates the principle of deep supervision.

Transition w/o Pooling Layer. We compare the case without this designed layer (only 3 dense blocks) and the case with the designed layer (4 dense blocks in our design). The backbone network is DS/32-12-16-0.5. Results are shown in Table 3. The network structure with the *Transition w/o pooling layer* brings 1.7% performance gain, which validates the effectiveness of this layer.

Stem Block. As can be seen in Table 3 (rows 6 and 9), the stem block improves the performance from 74.5% to 77.3%. This validates our conjecture that using stem block can protect information loss from the raw input images.

Method	data	pre-train	transition w/o pooling	stem	backbone network	prediction Layer	# parameters	mAP
DSOD300	07+12	✗	✗	✗	DS/32-12-16-0.5	Plain	4.1M	59.9
DSOD300	07+12	✗	✓	✗	DS/32-12-16-0.5	Plain	4.2M	61.6
DSOD300	07+12	✗	✓	✗	DS/32-12-16-1	Plain	5.5M	64.5
DSOD300	07+12	✗	✓	✗	DS/32-64-16-1	Plain	6.1M	68.6
DSOD300	07+12	✗	✓	✗	DS/64-64-16-1	Plain	6.3M	69.7
DSOD300	07+12	✗	✓	✗	DS/64-192-48-1	Plain	18.0M	74.5
DSOD300	07+12	✗	✓	✓	DS/64-12-16-1	Plain	5.2M	70.7
DSOD300	07+12	✗	✓	✓	DS/64-36-48-1	Plain	12.5M	76.0
DSOD300	07+12	✗	✓	✓	DS/64-192-48-1	Plain	18.2M	77.3
DSOD300	07+12	✗	✓	✓	DS/64-64-16-1	Dense	5.9M	73.6
DSOD300	07+12	✗	✓	✓	DS/64-192-48-1	Dense	14.8M	77.7
DSOD300	07+12+COCO	✗	✓	✓	DS/64-192-48-1	Dense	14.8M	81.7

Table 3: **Ablation study on PASCAL VOC 2007 test set.** DS/A-B- k - θ describes our backbone network structure. A denotes the number of channels in the 1st conv-layer. B denotes the number of channels in each *bottleneck* layer (1×1 convolution). k is the growth rate in dense blocks. θ denotes the compression factor in transition layers. See Section 4.1 for more explanations.

Method	data	pre-train	backbone network	prediction layer	speed (fps)	# parameters	input size	mAP
Faster RCNN [27]	07+12	✓	VGGNet	-	7	134.7M	~ 600 × 1000	73.2
Faster RCNN [27]	07+12	✓	ResNet-101	-	2.4*	-	~ 600 × 1000	76.4
R-FCN [19]	07+12	✓	ResNet-50	-	11	31.9M	~ 600 × 1000	77.4
R-FCN [19]	07+12	✓	ResNet-101	-	9	50.9M	~ 600 × 1000	79.5
R-FCNmulti-sc [19]	07+12	✓	ResNet-101	-	9	50.9M	~ 600 × 1000	80.5
YOLOv2 [26]	07+12	✓	Darknet-19	-	81	-	352 × 352	73.7
SSD300 [21]	07+12	✓	VGGNet	Plain	46	26.3M	300 × 300	75.8
SSD300* [21]	07+12	✓	VGGNet	Plain	46	26.3M	300 × 300	77.2
Faster RCNN	07+12	✗	VGGNet/ResNet-101/DenseNet	Failed				
R-FCN	07+12	✗	VGGNet/ResNet-101/DenseNet	Failed				
SSD300S [†]	07+12	✗	ResNet-101	Plain	12.1	52.8M	300 × 300	63.8*
SSD300S [†]	07+12	✗	VGGNet	Plain	46	26.3M	300 × 300	69.6
SSD300S [†]	07+12	✗	VGGNet	Dense	37	26.0M	300 × 300	70.4
DSOD300	07+12	✗	DS/64-192-48-1	Plain	20.6	18.2M	300 × 300	77.3
DSOD300	07+12	✗	DS/64-192-48-1	Dense	17.4	14.8M	300 × 300	77.7
DSOD300	07+12+COCO	✗	DS/64-192-48-1	Dense	17.4	14.8M	300 × 300	81.7

Table 4: **PASCAL VOC 2007 test detection results.** SSD300* is updated version by the authors after the paper publication. SSD300S[†] indicates training SSD300* from scratch with ResNet-101 or VGGNet, which serves as our baseline. Note that the speed of Faster R-CNN with ResNet-101 (2.4 fps) is tested on K40, while others are tested on Titan X. The result of SSD300S with ResNet-101 (63.8% mAP, without the pre-trained model) is produced with the default setting of SSD, which may not be optimal.

Dense Prediction Structure. We analyze the dense prediction structure from three aspects: speed, accuracy and parameters. As shown in Table 4, DSOD with dense front-end structure runs slightly lower than the plain structure (17.4 fps vs. 20.6 fps) on a Titan X GPU, due to the overhead from additional down-sampling blocks. However, the dense structure improves mAP from 77.3% to 77.7%, while reduces the parameters from 18.2M to 14.8M. Table 3 gives more details (rows 9 and 10). We also tried to replace the prediction layers in SSD with the proposed dense prediction layers. The accuracy on VOC 2007 test set can be improved from 75.8% (original SSD) to 76.1% (with pre-trained models), and 69.6% to 70.4% (w/o pre-trained models), when using the VGG-16 models as backbone. This verifies the effectiveness of the dense prediction layer.

What if pre-training on ImageNet? It is interesting to see the performance of DSOD with backbone network pre-trained on ImageNet. We trained one lite backbone network DS/64-12-16-1 on ImageNet, which obtains 66.8% top-1 accuracy and 87.8% top-5 accuracy on the validation-set (slightly worse than VGG-16). After fine-tuning the whole detection framework on “07+12” trainval set, we achieve 70.3% mAP on the VOC 2007 test-set. The corresponding training-from-scratch solution achieves 70.7%

accuracy, which is even slightly better. Future work will investigate this point more thoroughly.

4.1.3 Runtime Analysis

The inference speed is shown in the 6th column of Table 4. With 300×300 input, our full DSOD can process an image in 48.6ms (20.6 fps) on a single Titan X GPU with the plain prediction structure, and 57.5ms (17.4 fps) with the dense prediction structure. As a comparison, R-FCN runs at 90ms (11 fps) for ResNet-50 and 110ms (9 fps) for ResNet-101. The SSD300* runs at 82.6ms (12.1 fps) for ResNet-101 and 21.7ms (46 fps) for VGGNet. In addition, our model uses about only 1/2 parameters to SSD300 with VGGNet, 1/4 to SSD300 with ResNet-101, 1/4 to R-FCN with ResNet-101 and 1/10 to Faster R-CNN with VGGNet. A lite-version of DSOD (10.4M parameters, w/o any speed optimization) can run 25.8 fps with only 1% mAP drops.

4.2. Results on PASCAL VOC2007

Models are trained based on the union of VOC 2007 trainval and VOC 2012 trainval (“07+12”) following [21]. We use a batch size of 128. Note that this batch-size is beyond the capacity of GPU memories (even for an 8 GPU server, each with 12GB memory). We use a trick

Method	data	backbone network	pre-train	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
ION [1]	07+12+S	VGGNet	✓	76.4	87.5	84.7	76.8	63.8	58.3	82.6	79.0	90.9	57.8	82.0	64.7	88.9	86.5	84.7	82.3	51.4	78.2	69.2	85.2	73.5
Faster RCNN [27]	07++12	ResNet-101	✓	73.8	86.5	81.6	77.2	58.0	51.0	78.6	76.6	93.2	48.6	80.4	59.0	92.1	85.3	84.8	80.7	48.1	77.3	66.5	84.7	65.6
R-FCN _{multi-sc} [19]	07++12	ResNet-101	✓	77.6	86.9	83.4	81.5	63.8	62.4	81.6	81.1	93.1	58.0	83.8	60.8	92.7	86.0	84.6	84.4	59.0	80.8	68.6	86.1	72.9
YOLOv2 [26]	07++12	Darknet-19	✓	73.4	86.3	82.0	74.8	59.2	51.8	79.8	76.5	90.6	52.1	78.2	58.5	89.3	82.5	83.4	81.3	49.1	77.2	62.4	83.8	68.7
SSD300* [21]	07++12	VGGNet	✓	75.8	88.1	82.9	74.4	61.9	47.6	82.7	78.8	91.5	58.1	80.0	64.1	89.4	85.7	85.5	82.6	50.2	79.8	73.6	86.6	72.1
DSOD300	07++12	DS/64-192-48-1	✗	76.3	89.4	85.3	72.9	62.7	49.5	83.6	80.6	92.1	60.8	77.9	65.6	88.9	85.5	86.8	84.6	51.1	77.7	72.3	86.0	72.2
DSOD300	07++12+COCO	DS/64-192-48-1	✗	79.3	90.5	87.4	77.5	67.4	57.7	84.7	83.6	92.6	64.8	81.3	66.4	90.1	87.8	88.1	87.3	57.9	80.3	75.6	88.1	76.7

Table 5: **PASCAL VOC 2012 test detection results.** 07+12: 07 trainval + 12 trainval, 07+12+S: 07+12 plus segmentation labels, 07++12: 07 trainval + 07 test + 12 trainval. Result links are DSOD300 (07+12): <http://host.robots.ox.ac.uk:8080/anonymous/PIOBKI.html>; DSOD300 (07+12+COCO): <http://host.robots.ox.ac.uk:8080/anonymous/IOUHO.html>.

Method	data	network	pre-train	Avg. Precision, IoU:			Avg. Precision, Area:			Avg. Recall, #Dets:			Avg. Recall, Area:		
				0.5:0.95	0.5	0.75	S	M	L	1	10	100	S	M	L
Faster RCNN [27]	trainval	VGGNet	✓	21.9	42.7	-	-	-	-	-	-	-	-	-	-
ION [1]	train	VGGNet	✓	23.6	43.2	23.6	6.4	24.1	38.3	23.2	32.7	33.5	10.1	37.7	53.6
R-FCN [19]	trainval	ResNet-101	✓	29.2	51.5	-	10.3	32.4	43.3	-	-	-	-	-	-
R-FCN _{multi-sc} [19]	trainval	ResNet-101	✓	29.9	51.9	-	10.8	32.8	45.0	-	-	-	-	-	-
SSD300 (Huang et al.) [11]	< trainval35k	MobileNet	✓	18.8	-	-	-	-	-	-	-	-	-	-	-
SSD300 (Huang et al.) [11]	< trainval35k	Inception-v2	✓	21.6	-	-	-	-	-	-	-	-	-	-	-
YOLOv2 [26]	trainval35k	Darknet-19	✓	21.6	44.0	19.2	5.0	22.4	35.5	20.7	31.6	33.3	9.8	36.5	54.4
SSD300* [21]	trainval35k	VGGNet	✓	25.1	43.1	25.8	6.6	25.9	41.4	23.7	35.1	37.2	11.2	40.4	58.4
DSOD300	trainval	DS/64-192-48-1	✗	29.3	47.3	30.6	9.4	31.5	47.0	27.3	40.7	43.0	16.7	47.1	65.0

Table 6: **MS COCO test-dev 2015 detection results.**

to overcome the GPU memory constraints by accumulating gradients over two training iterations, which has been implemented on Caffe platform [14]. The initial learning rate is set to 0.1, and then divided by 10 after every 20k iterations. The training finished when reaching 100k iterations. Following [21], we use a weight decay of 0.0005 and a momentum of 0.9. All conv-layers are initialized with the “xavier” method [6].

Table 4 shows our results on VOC2007 test set. SSD300* is the updated SSD results which use the new data augmentation technique. Our DSOD300 with plain connection achieves 77.3%, which is slightly better than SSD300* (77.2%). DSOD300 with dense prediction structure improves the result to 77.7%. After adding COCO as training data, the performance is further improved to 81.7%.

4.3. Results on PASCAL VOC2012

For the VOC 2012 dataset, we use VOC 2012 trainval and VOC 2007 trainval + test for training, and test on VOC 2012 test set. The initial learning rate is set to 0.1 for the first 30k iterations, then divided by 10 after every 20k iterations. The total training iterations are 110k. Other settings are the same as those used in our VOC 2007 experiments. Our results of DSOD300 are shown in Table 5. DSOD300 achieves 76.3% mAP, which is consistently better than SSD300* (75.8%).

4.4. Results on MS COCO

Finally we evaluate our DSOD on the MS COCO dataset [20]. MS COCO contains 80k images for training, 40k for validation and 20k for testing (test-dev set). Following [27, 19], we use the trainval set (train set + validation set) for training. The batch size is also set as 128. The initial learning rate is set to 0.1 for the first 80k iterations, then divided by 10 after every 60k iterations. The

total number of training iterations is 320k.

Results are summarized in Table 6. Our DSOD300 achieves 29.3%/47.3% on the test-dev set, which outperforms the baseline SSD300* with a large margin. Our result is comparable to the single-scale R-FCN, and is close to the R-FCN_{multi-sc} which uses ResNet-101 as the pre-trained model. Interestingly, we observe that our result with 0.5 IoU is lower than R-FCN, but our [0.5:0.95] result is better or comparable. This indicates that our predicted locations are more accurate than R-FCN under the larger overlap settings. It is reasonable that our small object detection precision is slightly lower than R-FCN since our input image size (300×300) is much smaller than R-FCN’s (~600×1000). Even with this disadvantage, our large object detection precision is still much better than R-FCN. This further demonstrates the effectiveness of our approach. Figure 2 shows some qualitative detection examples on COCO with our DSOD300 model.

5. Discussion

Better Model Structure vs. More Training Data. An emerging idea in the computer vision community is that object detection or other vision tasks might be solved with deeper and larger neural networks backed with massive training data like ImageNet [3]. Thus more and more large-scale datasets have been collected and released recently, such as the Open Images dataset [16], which is 7.5x larger in the number of images and 6x larger of categories than that of ImageNet. We definitely agree that, under modest assumptions that given boundless training data and unlimited computational power, deep neural networks should perform extremely well. However, our proposed approach and experimental results imply an alternative view to handle this problem: a better model structure might enable similar or better performance compared with complex models trained

References

- [1] S. Bell, C. Lawrence Zitnick, et al. Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks. In *CVPR*, 2016. 7
- [2] L.-C. Chen, G. Papandreou, I. Kokkinos, et al. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015. 1
- [3] J. Deng, W. Dong, R. Socher, L.-J. Li, et al. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 1, 7
- [4] R. Girshick. Fast r-cnn. In *ICCV*, 2015. 1, 2, 3
- [5] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 1, 2, 3
- [6] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010. 7
- [7] S. Gupta, J. Hoffman, and J. Malik. Cross modal distillation for supervision transfer. In *CVPR*, 2016. 2
- [8] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Hypercolumns for object segmentation and fine-grained localization. In *CVPR*, 2015. 1
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 3, 4
- [10] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten. Densely connected convolutional networks. In *CVPR*, 2017. 1, 2, 3, 4
- [11] J. Huang, V. Rathod, C. Sun, et al. Speed/accuracy trade-offs for modern convolutional object detectors. In *CVPR*, 2017. 3, 7
- [12] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 3
- [13] S. Jégou, M. Drozdal, D. Vazquez, et al. The one hundred layers tiramisù: Fully convolutional densenets for semantic segmentation. *arXiv preprint arXiv:1611.09326*, 2016. 3
- [14] Y. Jia, E. Shelhamer, J. Donahue, et al. Caffe: Convolutional architecture for fast feature embedding. In *ACM MM*, 2014. 5, 7
- [15] K.-H. Kim, S. Hong, B. Roh, et al. Pvanet: Deep but lightweight neural networks for real-time object detection. *arXiv preprint arXiv:1608.08021*, 2016. 3
- [16] I. Krasin, T. Duerig, N. Alldrin, A. Veit, et al. Openimages: A public dataset for large-scale multi-label and multi-class image classification. <https://github.com/openimages>, 2016. 7
- [17] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1, 3
- [18] C.-Y. Lee, S. Xie, P. W. Gallagher, et al. Deeply-supervised nets. In *AISTATS*, 2015. 2, 4
- [19] Y. Li, K. He, J. Sun, et al. R-fcn: Object detection via region-based fully convolutional networks. In *NIPS*, 2016. 1, 2, 6, 7
- [20] T.-Y. Lin, M. Maire, S. Belongie, et al. Microsoft coco: Common objects in context. In *ECCV*, 2014. 7
- [21] W. Liu, D. Anguelov, D. Erhan, et al. Ssd: Single shot multi-box detector. In *ECCV*, 2016. 1, 2, 3, 4, 6, 7
- [22] W. Liu, A. Rabinovich, and A. C. Berg. Parsenet: Looking wider to see better. *arXiv preprint arXiv:1506.04579*, 2015. 5
- [23] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 1
- [24] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *CVPR*, 2014. 1
- [25] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 1, 2, 3
- [26] J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. In *CVPR*, 2017. 3, 6, 7
- [27] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 1, 2, 3, 6, 7
- [28] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 1, 3
- [29] N. Srivastava, G. E. Hinton, A. Krizhevsky, et al. Dropout: a simple way to prevent neural networks from overfitting. *JMLR*, 2014. 3
- [30] Y. Sun, D. Liang, X. Wang, and X. Tang. Deepid3: Face recognition with very deep neural networks. *arXiv preprint arXiv:1502.00873*, 2015. 4
- [31] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *ICLR workshop*, 2016. 3, 4
- [32] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, et al. Going deeper with convolutions. In *CVPR*, 2015. 1, 3, 4
- [33] C. Szegedy, V. Vanhoucke, S. Ioffe, et al. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 4
- [34] J. R. Uijlings, K. E. Van De Sande, T. Gevers, et al. Selective search for object recognition. *IJCV*, 2013. 2
- [35] S. Xie and Z. Tu. Holistically-nested edge detection. In *ICCV*, 2015. 2, 4
- [36] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016. 1