# Localization-aware Channel Pruning for Object Detection

Zihao Xie[a], Li Zhu[a,*], Lin Zhao[a], Bo Tao[c], Liman Liu[d], Wenbing Tao[a,b]

[a]*National Key Laboratory of Science and Technology on Multi-spectral Information Processing, School of Artificial Intelligence and Automation, Huazhong University of Science and Technology, Wuhan 430074, China*
[b]*Shenzhen Huazhong University of Science and Technology Research Institute, Shenzhen, 518057, China*
[c]*State Key Laboratory of Digital Manufacturing Equipment and Technology, School of Mechanical Science and Engineering, Huazhong University of Science and Technology, Wuhan, Hubei 430074, PR China*
[d]*School of Biomedical Engineering, South-Central University for Nationalities, Wuhan 430074, China*

## Abstract

Channel pruning is one of the important methods for deep model compression. Most of existing pruning methods mainly focus on classification. Few of them conduct systematic research on object detection. However, object detection is different from classification, which requires not only semantic information but also localization information. In this paper, based on discrimination-aware channel pruning (DCP) which is state-of-the-art pruning method for classification, we propose a localization-aware auxiliary network to find out the channels with key information for classification and regression so that we can conduct channel pruning directly for object detection, which saves lots of time and computing resources. In order to capture the localization information, we first design the auxiliary network with a contextual RoIAlign layer which can obtain precise localization information of the default boxes by pixel alignment and enlarges the receptive fields of the default boxes when pruning shallow layers. Then, we construct a loss function for object detection task which tends to keep the channels that contain the key information for classification and regression. Extensive

---

[*]Corresponding author

*Email addresses:* `zihaoxie@hust.edu.cn` (Zihao Xie), `lizhu2016@hust.edu.cn` (Li Zhu), `wenbingtao@hust.edu.cn` (Wenbing Tao)

experiments demonstrate the effectiveness of our method. On MS COCO, we prune 70% parameters of the SSD based on ResNet-50 with modest accuracy drop, which outperforms the-state-of-art method.

*Keywords:* channel pruning, object detection, localization-aware

## 1. Introduction

Since AlexNet [1] won the ImageNet Challenge: ILSVRC 2012 [2], deep convolutional neural network (CNNs) have been widely applied to various computer vision tasks, from basic image classification tasks [3] to some more advanced applications, e.g., object detection [4, 5], semantic segmentation [6], video analysis [7] and many others. In these fields, CNNs have achieved state-of-the-art performance compared with traditional methods based on manually designed visual features.

However, deep models often have a huge number of parameters and its size is very large, which incurs not only huge memory requirement but also unbearable computation burden. As a result, a typical deep model is hard to be deployed on resource constrained devices, e.g., mobile phones or embedded gadgets. To make CNNs available on resource-constrained devices, there are lots of studies on model compression, which aims to reduce the model redundancy without significant degeneration in performance. Channel pruning [8, 9, 10] is one of the important methods. Different from simply making sparse connections [11, 12], channel pruning reduces the model size by directly removing redundant channels and can achieve fast inference without special software or hardware implementation.

In order to determine which channels to reserve, existing reconstruction-based methods [8, 9, 10] usually minimize the reconstruction error of feature maps between the original model and the pruned one. However, a well-reconstructed feature map may not be optimal for there is a gap between intermediate feature map and the performance of final output. Information redundancy channels could be mistakenly kept to minimize the reconstruction error of feature maps.

2

To find the channels with true discriminative power for the network, DCP [13] attend to conduct channel selection by introducing additional discrimination-aware losses that are actually correlated with the final performance. It constructs the discrimination-aware losses by a fully connected layer which works on the entire feature map. However, the discrimination-aware loss of DCP is designed for classification task. Since object detection network uses the classification network as backbone, a simple method to conduct DCP for object detection is to fine-tune the pruned model, which was trained on classification dataset, for the object detection task. But the information that the two tasks need is not exactly the same. The classification task needs strong semantic information while what the object detection task needs is not only semantic information but also localization information. Hence, the existing training scheme may not be optimal due to the mismatched goals of feature learning for classification and object detection task.

In this paper, we propose a method called localization-aware channel pruning (LCP), which conducts channel pruning directly for object detection. We propose a localization-aware auxiliary network for object detection task. First, we design the auxiliary network with a contextual RoIAlign layer which can obtain precise localization information of the default boxes by pixel alignment and enlarges the receptive fields of the default boxes when pruning shallow layers. Then, we construct a loss function for object detection task which tends to keep the channels that contain the key information for classification and regression. Our main contributions are summarized as follows. (1) We propose a localization-aware auxiliary network which can find out the channels with key information so that we can conduct channel pruning directly on object detecion dataset, which saves lots of time and computing resources. (2) We propose a contextual RoIAlign layer which enlarges the receptive fields of the default boxes in shallow layers. (3) Extensive experiments on benchmark datasets show that the proposed method is theoretically reasonable and practically effective. For example, our method can prune 70% parameters of SSD [4] based on ResNet-50 [3] with modest accuracy drop on VOC2007, which outperforms the-state-of-art

method.

## 2. Related Works

### 2.1. Network Quantization

Network quantization compresses the original network by reducing the number of bits required to represent each weight. Han et al. [11] propose a complete deep network compression pipeline: First trim the unimportant connections and retrain the sparsely connected network. Weight sharing is then used to quantize the weight of the connection, and then the quantized weight and codebook are Huffman encoded to further reduce the compression ratio. Courbariaux et al. [14] propose to accelerate the model by reducing the weight and accuracy of the output, because this will greatly reduce the memory size and access times of the network, and replace the arithmetic operator with a bit-wise operator. Li et al. [15] consider that multi-weights have better generalization capabilities than binarization and the distribution of weights is close to a combination of a normal distribution and a uniform distribution. Zhou et al. [16] propose a method which can convert the full-precision CNN into a low-precision network, making the weights 0 or 2 without loss or even higher precision (shifting can be performed on embedded devices such as FPGAs). For more recent works, Yu et al. [17], to reduce the communication complexity, propose a general scheme for quantizing both model parameters and gradients. Zhao et al. [18] attend to the statistical properties of sparse CNNs and present focused quantization, a novel quantization strategy based on power-of-two values, which exploits the weight distributions after fine-grained pruning.

### 2.2. Sparse or Low-rank Connections

Wen et al. [19] propose a learning method called Structured Sparsity Learning, which can learn a sparse structure to reduce computational cost, and the learned structural sparseness can be effectively accelerate for hardware. Guo et al. [20] propose a new network compression method, called dynamic network

surgery, is to reduce network complexity through dynamic connection pruning. Unlike previous methods of greedy pruning, this approach integrates join stitching throughout the process to avoid incorrect trimming and maintenance of the network. Jin et al. [21] proposes to reduce the computational complexity of the model by training a sparsely high network. By adding a $l_0$ paradigm about weights to the loss function of the network, the sparsity of weights can be reduced. For more recent works, Kim [22] propose novel accuracy metrics to represent the accuracy and complexity relationship for a given neural network and use these metrics in a non-iterative fashion to obtain the right rank configuration which satisfies the constraints on FLOPs and memory while maintaining sufficient accuracy. Liu et al. [23] propose a layerwise sparse coding (LSC) method to maximize the compression ratio by extremely reducing the amount of meta-data.

### 2.3. Channel Pruning

Finding unimportant weights in the network has a long history. LeCun [24] and Hassibi [25] consider using the Hessian, which contains second order derivative, performs better than using the magnitude of the weights. Computing the Hessian is expensive and thus is not widely used. Han [11] et al. proposed an iterative pruning method to remove the redundancy in deep models. Their main insight is that small-weight connectivity below a threshold should be discarded. In practice, this can be aided by applying $l_1$ or $l_2$ regularization to push connectivity values to become smaller. The major weakness of this strategy is the loss of universality and flexibility, thus seems to be less practical in real applications. Li et al. [26] measure the importance of channels by calculating the sum of absolute values of weights. Hu et al. [27] define average percentage of zeros (APoZ) to measure the activation of neurons. Neurons with higher values of APoZ are considered more redundant in the network. With a sparsity regularizer in the objective function [28, 29], training based methods are proposed to learn the compact models in the training phase. With the consideration of efficiency, reconstruction-methods [8, 9] transform the channel

5

selection problem into the optimization of reconstruction error and solve it by a greedy algorithm or LASSO regression. DCP [13] aimed at selecting the most discriminative channels for each layer by considering both the reconstruction error and the discrimination-aware loss. PP [30] jointly prunes and fine-tunes CNN model parameters, with an adaptive pruning rate. Liu et al. [31] propose a novel meta learning approach for automatic channel pruning of very deep neural networks. AutoPrune [32] prunes the network through optimizing a set of trainable auxiliary parameters instead of original weights.

### 2.4. Object Detection

Current state-of-the-art object detectors with deep learning can be mainly divided into two major categories: two-stage detectors and one-stage detectors. Two-stage detectors first generate region proposals which may potentially be objects and then make predictions for these proposals. Faster R-CNN [5] is a representative two-stage detector, which was able to make predictions at 5FPS on GPU and achieved state-of-the-art results on many public benchmark datasets, such as Pascal VOC 2007, 2012 and MSCOCO. Currently, there are huge number of detector variants based on Faster R-CNN for different usage [33, 34]. Mask R-CNN [35] extends Faster R-CNN to the field of instance segmentation. Based on Mask R-CNN, Huang et al. [36] proposed a mask-quality aware framework, named Mask Scoring R-CNN, which learned the quality of the predicted masks and calibrated the misalignment between mask quality and mask confidence score.

Different from two-stage detectors, one-stage detectors do not generate proposals and directly make predictions on the whole feature map. SSD [4] is a representative one-stage detector. However, the class imbalance between foreground and background is a severe problem in one-stage detector. RetinaNet [37] matigates the class imbalance problem by introducing focal loss which reduces loss of easy samples. The previous approaches required designing anchor boxes manually to train a detector. Currently, a series of anchor-free object detectors were developed, where the goal was to predict keypoints of the bounding

6

box, instead of trying to fit an object to an anchor. Law and Deng proposed a novel anchor-free framework CornerNet [38] which detected objects as a pair of corners. Later there were several other variants of anchor-free detectors [39, 40]

## 3. Proposed Method

Fig. 1 is the overall frame diagram. The blue network in Figure 1 is the original model which could be any detectors, such as SSD [4], Faster R-CNN [5] and so on. The orange network in Figure 1 is the model to be pruned which is initialized exactly the same as the original model. The middle is the auxiliary network. We use the auxiliary network to prune the model by constructing the localization-aware loss. The localization-aware loss consists of two parts, one part is the reconstruction error and the other is the loss of auxiliary network. Their details will be discussed later. After the loss is constructed, we could fine-tune the network and use the gradient of the localization-aware loss to decide which channel to preserve. After repeating this operation layer by layer, the pruning is finished.

The auxiliary network we propose mainly consists of two parts. First, a contextual RoIAlign layer is designed to extract the features of the boxes. Then, a loss is designed for object detection task which can reserve the important channels. The details of the proposed approach are elaborated below.

### 3.1. Contextual RoIAlign Layer

For object detection task, if we predict the bounding boxes directly on the entire feature maps, there will be a huge amount of parameters and unnecessary noises. So, it is important to extract the feature of region of interest (RoI) , which can be better used for classification and regression. To obtain precise localization information and find out the channels which are important for classification and regression, RoIAlign layer is a good choice which properly align the extracted features with the input. RoIAlign use bilinear interpolation to compute the exact values of the input features at four regularly sampled locations in each RoI bin, and aggregate the result (using max or average), see
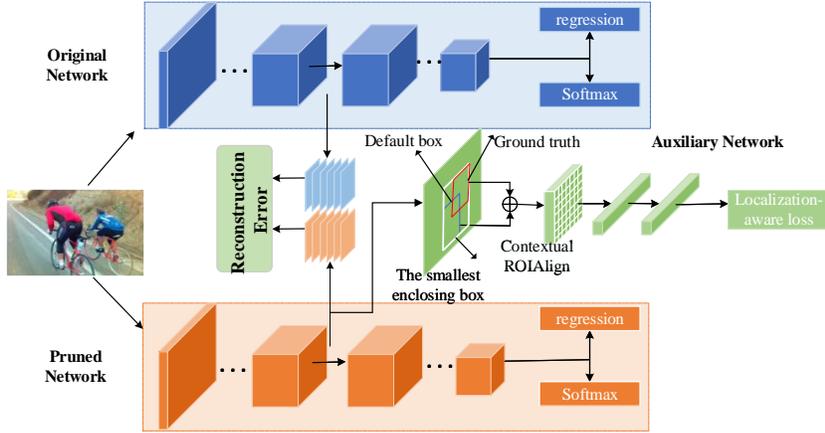
7

Figure 1: Illustration of localization-aware channel pruning. The auxiliary network is used to supervise layer-wise channel selection. The blue network the original model, the orange network is the model to be pruned, the green network is the auxiliary network. The reconstruction error and the auxiliary network are used to construct localization-aware loss.

Fig. 2 for details. However, the default boxes generated by the detector do not always completely cover the object area. From Fig. 3, we can see that, the defalut box is sometimes bigger than the ground truth and sometimes smaller than it. So, the receptive fields may be insufficient if we only extract the features of the default box especially when we prune the shallow layers. To solve this problem, we propose a contextual RoIAlign layer, which introduces larger context information. The orange part in Figure 4 is feature extraction network. The network first extract the feature map of the whole image, then obtain the features of the boxes by RoIAlign operation. To introduce larger context information, we further gather the information of the default box and its context by adding the feature of the default box and its enclosing convex object.

For better description of the algorithm, some notations are given first. For a training sample, $(x_{a1}, y_{a1}, x_{a2}, y_{a2})$ represents the coordinates of ground truth box $A$, $(x_{b1}, y_{b1}, x_{b2}, y_{b2})$ denotes the coordinates of the matched default box $B$. We further use $\mathcal{F}$ to denote the feature map and $\mathcal{F}_{\mathcal{S}}$ represents the features of
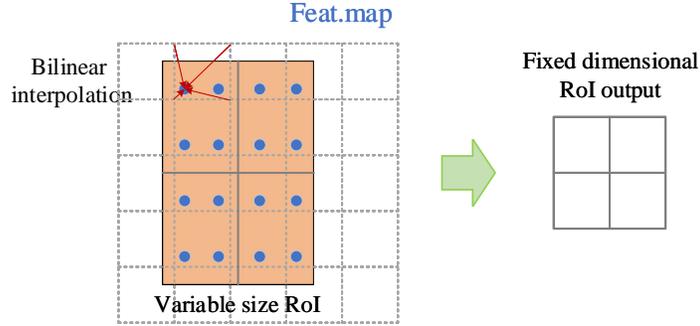
8

Figure 2: RoIAlign: The dashed grid represents a feature map, the solid lines an RoI (with 2× bins in this example), and the dots the 4 sampling points in each bin. RoIAlign computes the value of each sampling point by bilinear interpolation from the nearby grid points on the feature map.

area $S$, $RoIAlign$ represents the RoIAlign operation. First, we calculate the $IoU$ of box A and B:

$$IoU_{AB} = \frac{A \cap B}{A \cup B} \tag{1}$$

B is a positive sample only if $IoU_{AB}$ is larger than a preset threshold. We do not conduct contextual RoIAlign for B when B is negative sample. If B is a positive sample, then we calculate the smallest enclosing convex object C for A and B:

$$x_{c1} = min(x_{a1}, x_{b1}) \tag{2}$$

$$y_{c1} = min(y_{a1}, y_{b1}) \tag{3}$$

$$x_{c2} = max(x_{a2}, x_{b2}) \tag{4}$$

$$y_{c2} = max(y_{a2}, y_{b2}) \tag{5}$$

Figure 3: The features of default boxes do not always contain enough context information, especially when we prune shallow layers. The blue is default box, the red is ground truth.

where $(x_{c1}, y_{c1}, x_{c2}, y_{c2})$ are the coordinates of C. Finally, the output of contextual RoIAlign layer is defined as:

$$\mathcal{F}_{\mathcal{O}} = RoIAlign(\mathcal{F}_{\mathcal{B}}) + RoIAlign(\mathcal{F}_{\mathcal{C}}) \qquad (6)$$

Now we can get the precise features of default box B, the process can refer to Fig 4.

*3.2. Construction of the Loss of Auxiliary Network*

After we construct the contextual RoIAlign layer, we need to consturct a loss for object detection task so that we can use the gradient of the auxiliary network to conduct model pruning. The details are discussed below.

Considering that the object detection task needs both the classification and localization information. The loss of auxiliary consists of two parts. One is the loss for classification and the other is the loss for regression. For the classification part, We still use discrimination-aware loss which has been proved to be useful on classification task. For the regression part, we choose GIoU [41] loss function. It is reasonable to use GIoU as loss function for boxes regression. It considers not only overlapping areas but also non-overlapping areas, which better reflects the overlap of the boxes. The GIoU of two arbitrary shapes $A$ and $B$ is defined as:
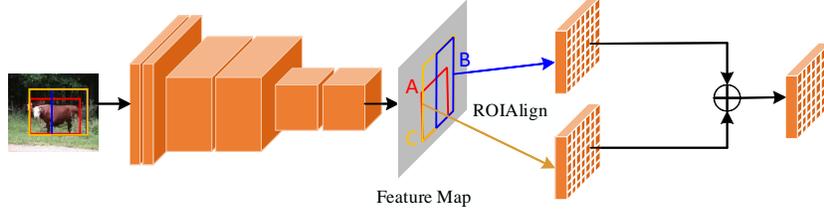
Figure 4: Contextual RoIAlign: The red, blue and orange boxes represent the ground truth, default box and its smallest enclosing box, respectively. The network first extract the feature map of the whole image, then obtain the features of the boxes by RoIAlign operation. To introduce larger context information, we further gather the information of the default box and its context by adding the feature of the default box and its enclosing convex object.

$$GIoU_{AB} = IoU_{AB} - \frac{C - U}{C} \tag{7}$$

where $U = A + B - IoU_{AB}$, $IoU_{AB}$ and $C$ are calculated by Eq. 1 - Eq. 5. Fig. 5 is a schematic diagram of GIoU. Then, we use $G_i$ to denote the GIoU of the $i$-th predicted box and the ground truth, $E_i$ to represent the cross entropy of the $i$-th predicted box. Then, in the pruning stage, $\mathcal{L}_{ac}$ represents the classification loss, $\mathcal{L}_{ar}$ represents the regression loss, $\mathcal{L}_a$ represents the localization-aware loss of the auxiliary network. Finally, the loss of positive samples in pruning stage is defined as:

$$\mathcal{L}_{ac} = \sum_i E_i \tag{8}$$

$$\mathcal{L}_{ar} = \sum_i m(1 - G_i) \tag{9}$$

$$\mathcal{L}_a = \mathcal{L}_{ac} + \mathcal{L}_{ar} \tag{10}$$

where $m$ is a constant coefficient.

### 3.3. Localization-aware Channel Pruning

After we construct the auxiliary network and the localization-aware loss, we can conduct channel pruning with them layer by layer. The pruning process
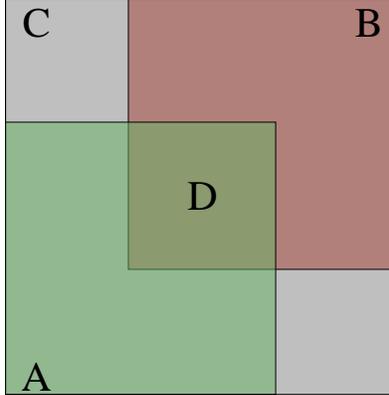
Figure 5: A, B are two arbitrary shapes, C is the smallest enclosing convex of A and B, D is the $[IoU]$ of A and B.

of the whole model is described in Algorithm 1. For better description of the channel selection algorithm, some notations are given first. Considering a $L$ layers of the CNN model and we are pruning the $l$-th layer, $X$ represents the output feature map of the $lth$ layer, $W$ denotes the convolution filter of the $(l+1)$-th layer of the pruned model and $*$ represents the convolution operation. We further use $F \in R^{N \times HY}$ to denote output feature maps of the $(l+1)$-th layer of the original model. Here, $N$, $H$, $Y$ represents the number of output channels, the height and the width of the feature maps respectively. Finally we use $\mathcal{L}_c$ and $\mathcal{L}_r$ to denote classification loss and regression loss of the pruned network.

To find out the channels which really contribute to the network, we should fine-tune the auxiliary network and pruned network first and the fine-tune loss is defined as the sum of the losses of them:

$$\mathcal{L}_f = \mathcal{L}_a + \mathcal{L}_c + \mathcal{L}_r \qquad (11)$$

In order to minimizing the reconstruction error of a layer, we introduce a reconstruction loss as DCP does which can be defined as the Euclidean distance

12

of feature maps between the original model and the pruned one:

$$\mathcal{L}_{re} = \frac{1}{2Q} \|F - X * W_{\mathcal{C}}\|_2^2 \tag{12}$$

where $Q = M \times H \times Y$, $\mathcal{C}$ represents the selected channels, $W_{\mathcal{C}}$ represents the submatrix indexed by $\mathcal{C}$.

---

**Algorithm 1** The proposed method

---

**Input:** number of layers $L$, weights of original model $\{W^l : 0 < l < L\}$, the training set $\{x_i, y_i\}$, the pruning rate $\eta$.

**Output:**          $\{W_{\mathcal{C}}^l : 0 < l < L\}$:          weights     of     the     pruned model.

1: Initialize $W_{\mathcal{C}}^l$ with $W^l$ for $\forall 1 \leq l \leq L$

2: **for** $l = 1, 2, \cdots, L$ **do**

3:     Construct the fine-tune loss $\mathcal{L}_f$ shown as in Eq. 11

4:     Fine-tune the auxiliary network and the pruned model by $\mathcal{L}_f$

5:     Construct the joint loss $\mathcal{L}$ shown as in Eq. 13

6:     Conduct channel selection for layer $l$ by Eq. 14

7:     Update $W_{\mathcal{C}}^l$ w.r.t. the selected channels by Eq. 15

8: **end for**

9: return the pruned model

---

Taking into account the reconstruction error, the localization-aware loss of the auxiliary network, the problem of channel pruning can be formulated to minimize the following joint loss function:

$$\min_{W_{\mathcal{C}}} \quad \mathcal{L}(W_{\mathcal{C}}) = \mathcal{L}_{re}(W_{\mathcal{C}}) + \alpha \mathcal{L}_a(W_{\mathcal{C}})$$

$$s.t. \quad \|\mathcal{C}\|_0 \leq \mathcal{K} \tag{13}$$

where $\alpha$ is a constant, $\mathcal{K}$ is the number of channels to be selected. Directly optimizing Eq. 13 is NP-hard. Following general greedy methods in DCP, we conduct channel pruning by considering the gradient of Eq. 13. Specifically, the
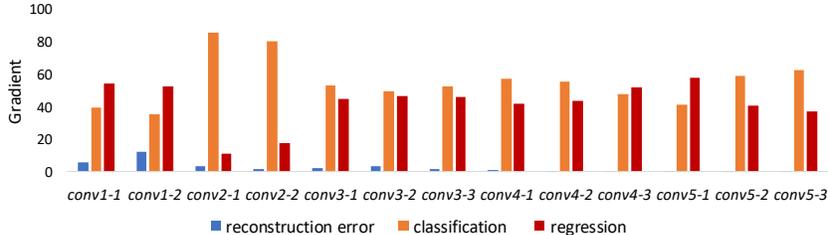
Figure 6: The percentage of the gradients generated by the three loss functions.

importance of the $k$-th channel is defined as:

$$\mathcal{S}_k = \sum_{i=1}^{H} \sum_{j=1}^{W} \|\frac{\partial \mathcal{L}}{\partial W_{k,i,j}}\|_2^2 \tag{14}$$

$\mathcal{S}_k$ is the square sum of gradient of the $k$-th channel. Then we reserve the channels with the $i$ largest importance and remove others. After this, the selected channels is further optimized by stochastic gradient (SGD). $W_{\mathcal{C}}$ is updated by:

$$W_{\mathcal{C}} = W_{\mathcal{C}} - \gamma \frac{\partial \mathcal{L}}{\partial W_{\mathcal{C}}} \tag{15}$$

where $\gamma$ represents the learning rate. After updating $W_{\mathcal{C}}$, the channel pruning of a single layer is finished.

## 4. Experiments

We evaluate LCP on the popular 2D object detector SSD [4]. Several state-of-the-art methods are adopted as the baselines, including ThiNet [9] and DCP [13]. In order to verify the effectiveness of our method, we use VGG and ResNet to extract feature respectively.

### 4.1. Dataset and Evaluation

The results of all baselines are reported on standard object detection benchmarks, i.e. the PASCAL VOC [42] . **PASCAL VOC2007 and 2012:** The

Pascal Visual Object Classes (VOC) benchmark is one of the most widely used datasets for classification, object detection and semantic segmentation. We use the union of VOC2007 and VOC2012 trainval as training set, which contains 16551 images and objects from 20 pre-defined categories annotated with bounding boxes. And we use the VOC2007 test as test set which contains 4592 images. In order to verify the effectiveness of our method, on PASCAL VOC, we first compare our method only with ThiNet based on VGG-16 because the authors of DCP do not release the VGG model. To this end, we compare our method with DCP and ThiNet based on ResNet-50. Then we conduct the ablation experiment of our method on PASCAL VOC. In order to more fully verify the effectiveness of our method, we also perform experiments on the MS COCO2017 dataset.

In this paper, we use $07metric$ for all experiments on PASCAL VOC. For experiments on MS COCO, the main performance measure used in this benchmark is shown by AP, which is averaging mAP across different value of IoU thresholds, i.e. $IoU = \{.5, .55, \cdots, .95\}$.

### 4.2. Implementation details

Our experiments are based on SSD and the input size of the SSD is $300 \times 300$. We use VGGNet and ResNet as the feature extraction network for experiments. For ThiNet, we implement it for object detection. And the three methods prune the same number of channels for each layer. Other common parameters are described in detail below.

For VGGNet [43], we use VGG-16 without Batch Normalization layer and prune the SSD from conv1-1 to conv5-3. The network is fine-tuned for 10 epochs every time a layer is pruned and the learning rate is started at 0.001 and divided by 10 at epoch 5. After the model is pruned, we fine-tune it for 60k iterations and the learning rate is started at 0.0005 and divided by 10 at iteration 30k and 45k, respectively.

For ResNet [3], we use the layers of ResNet-50 from conv1-x to conv4-x for feature extracting. The network is fine-tuned for 15 epochs every time a layer is

Table 1: The pruning results on PASCAL VOC2007. We conduct channel pruning from conv1-1 to conv5-3.

| Method | backbone | $\eta$ | flops↓ | params↓ | mAP |
|--------|----------|--------|--------|---------|-----|
| Original | VGG-16 | 0 | 0 | 0 | 77.4 |
| ThiNet | VGG-16 | 0.5 | 50% | 50% | 74.6 |
| LCP(our) | VGG-16 | 0.5 | 50% | 50% | **77.2** |
| ThiNet | VGG-16 | 0.75 | 75% | 75% | 72.7 |
| LCP(our) | VGG-16 | 0.75 | 75% | 75% | **75.2** |

pruned and the learning rate is started at 0.001 and divided by 10 at epoch 5 and 10, respectively. After the model is pruned, we fine-tune it for 120k iterations and the learning rate is started at 0.001 and divided by 10 at iteration 80k and 100k, respectively.

For the loss of auxiliary network, we set $m$ to 50.

### 4.3. Experiments on PASCAL VOC

On PASCAL VOC, we prune the VGG-16 from conv1-1 to conv5-3 with compression ratio 0.75, which is 4x faster. We report the results in Tab. 1. From the results, we can see that our method achieves the best performance under the same acceleration rate. The accuracy of reconstruction based method like ThiNet drops a lot. But for our LCP, there is not much degradation in the performance of object detection. It is proved that our method retain the channels which really contribute to the final performance. Then we conduct the experiment based ResNet-50. We report the results in Tab. 2. From the results, LCP achieves the best performance regardless of pruning by 75% or pruning by 50%, which proves that our method can reserve the channels which contain key information for classification and regression. In addition, the ThiNet outperforms the DCP when pruning ratio is 0.7, which indicates that pruning the model on classification dataset for object detection is not optimal.

Table 2: The pruning results on PASCAL VOC2007. We conduct channel pruning from conv2-x to conv4-x.

| Method | backbone | $\eta$ | flops↓ | params↓ | mAP |
|---|---|---|---|---|---|
| Original | ResNet-50 | 0 | 0 | 0 | 73.7 |
| DCP | ResNet-50 | 0.5 | 50% | 50% | 72.4 |
| ThiNet | ResNet-50 | 0.5 | 50% | 50% | 72.2 |
| LCP(our) | ResNet-50 | 0.5 | 50% | 50% | **73.3** |
| DCP | ResNet-50 | 0.7 | 70% | 70% | 70.2 |
| ThiNet | ResNet-50 | 0.7 | 70% | 70% | 70.8 |
| LCP(our) | ResNet-50 | 0.7 | 70% | 70% | **71.7** |

Table 3: The pruning results on MS COCO2017. The backbone is ResNet-50, We conduct channel pruning from conv2-x to conv4-x with compression ratio 0.7. Small, medium, large are the size of objects.

| Method | small | medium | large | $AP_{50}$ | $AP_{75}$ | mAP |
|---|---|---|---|---|---|---|
| Original | 4.2 | 22.5 | 39.0 | 37.3 | 22.7 | 21.9 |
| DCP | 2.8 | 17.2 | 33.0 | 31.8 | 17.8 | 17.8 |
| LCP | **4.1** | 20.4 | 38.2 | 35.5 | **21.6** | **20.9** |
| Relative improv.% | **46.4** | 18.6 | 15.8 | 11.6 | **21.3** | **17.4** |

## 4.4. Experiments on MS COCO

In this section, we prune the ResNet-50 by 70% on COCO2017. We report the results in Tab. 3 and Tab. 4. From the results, our method achieves a better performance than the DCP and ThiNet, which further illustrates the effectiveness of our approach. It is noted that compared with DCP, LCP has larger gain on small objects. In addition, the higher the IoU threshold, the greater improvement of our method. This indicates that our method retains more localization information and can obtain more accurate predictions.

## 4.5. Ablation Analysis

**Gradient Analysis.** In this section, we prune the VGG-16 from conv1-1 to conv5-3 with compression ratio 0.75 On PASCAL VOC. Then we count

Table 4: The pruning results on COCO. We conduct channel pruning from conv2-x to conv4-x.

| Method | backbone | $\eta$ | flops↓ | params↓ | mAP |
|--------|----------|--------|--------|---------|-----|
| Original | ResNet-50 | 0 | 0 | 0 | 21.9 |
| DCP | ResNet-50 | 0.5 | 50% | 50% | 21.2 |
| ThiNet | ResNet-50 | 0.5 | 50% | 50% | 22.6 |
| LCP(our) | ResNet-50 | 0.5 | 50% | 50% | **23.1** |
| DCP | ResNet-50 | 0.7 | 70% | 70% | 17.8 |
| ThiNet | ResNet-50 | 0.7 | 70% | 70% | 20.2 |
| LCP(our) | ResNet-50 | 0.7 | 70% | 70% | **20.9** |

Table 5: The pruning results on PASCAL VOC2007. We conduct channel pruning from conv2-x to conv4-x. CR means Contextual RoIAlign.

| Method | backbone | flops↓ | params↓ | mAP |
|--------|----------|--------|---------|-----|
| DCP | ResNet-50 | 70% | 70% | 70.2 |
| LCP+RoIAlign | ResNet-50 | 70% | 70% | **71.1** |
| LCP+CR | ResNet-50 | 70% | 70% | **71.7** |

the percentage of the gradients generated by the three losses during the pruning process. From Fig. 6, we see that the gradient of regression loss play a important role during the pruning process, which proves that the localization information is necessary. The gradient generated by reconstruction error only works in the shallow layers while the localization-aware loss contributes to the channel pruning process each layer.

**Component Analysis.** In this section, in order to verify the effectiveness of the two points we propose, we prune the SSD based on ResNet-50 by 70% with different combinations of our points. We report the results in Tab.5. From the results, we can get that each part of the method we propose contributes to the performance.

**Loss Analysis.** In order to explore the importance of the gradient of regression loss, we prune the SSD based on VGG-16 by 75% with different losses. We report the results in Tab. 6. From the results, we can know that the performance of our method drops a lot without the gradient of the regression loss during the pruning

Table 6: The pruning results on PASCAL VOC2007. We conduct channel pruning from conv1-1 to conv5-3 .

| Method | backbone | $\eta$ | flops↓ | params↓ | mAP |
|---|---|---|---|---|---|
| Original | VGG-16 | 0 | 0 | 0 | 77.4 |
| $\mathcal{L}_{re}+\mathcal{L}_{ac}$ | VGG-16 | 0.75 | 75% | 75% | 74.7 |
| $\mathcal{L}_{re}+\mathcal{L}_{ac}+\mathcal{L}_{ar}$ | VGG-16 | 0.75 | 75% | 75% | **75.2** |

stage, which shows that the regression branch contains important localization information.

### 4.6. Qualitative Results

To demonstrate the effectiveness of our proposed method in details, we use the detection analysis tool from [44]. Figure 7 shows that our model can detect various object categories with high quality (large blue area). The recall is higher than 90%, and is much higher with the weak (0.1 jaccard overlap) criteria. We can also observe that comparing with ThiNet, Our method has larger correct area, which indicates our superior performance.

### 4.7. Visualization of predictions

In this section, we prune the SSD based on VGG-16 by 75% and we compare the original model with the pruned models. From Fig. 8, we can find that the predictions of our method are closed to the predictions of the original model while the predictions of ThiNet are far away. It is proved that our method reserve more localization information for bounding box regression.

### 5. Conclusions

In this paper, we propose a localization-aware auxiliary network which allows us to conduct channel pruning directly for object detection. First, we design the auxiliary network with a contextual RoIAlign layer which can obtain precise localization information of the default boxes by pixel alignment and enlarges the
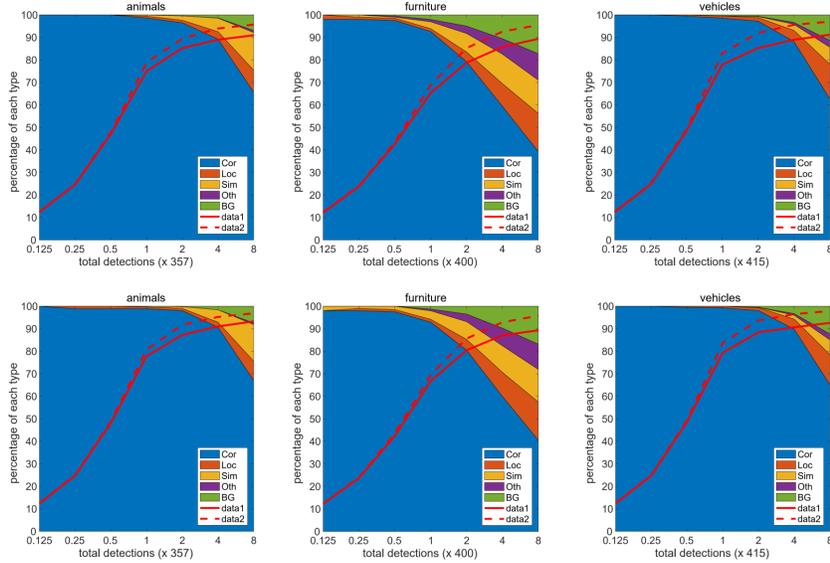
Figure 7: Visualization of the performance of ThiNet (top row) and our method (bottom row) on animals, furniture, and vehicles classes in the VOC 2007 test set. The figures show the cumulative fraction of detections that are correct (Cor) or false positives due to poor localization (Loc), confusion with similar categories (Sim), with others (Oth), or with background (BG). The solid red line reflects the change of recall with the strong criteria (0.5 jaccard overlap) as the number of detections increases. The dashed red line uses the weak criteria (0.1 jaccard overlap).

receptive fields of the default boxes when pruning shallow layers. Then, we construct a loss function for object detection task which tends to keep the channels that contain the key information for classification and regression. Visualization shows our method reserves layers with more localization information. Moreover, extensive experiments demonstrate the effectiveness of our method.

## 6. Acknowledge
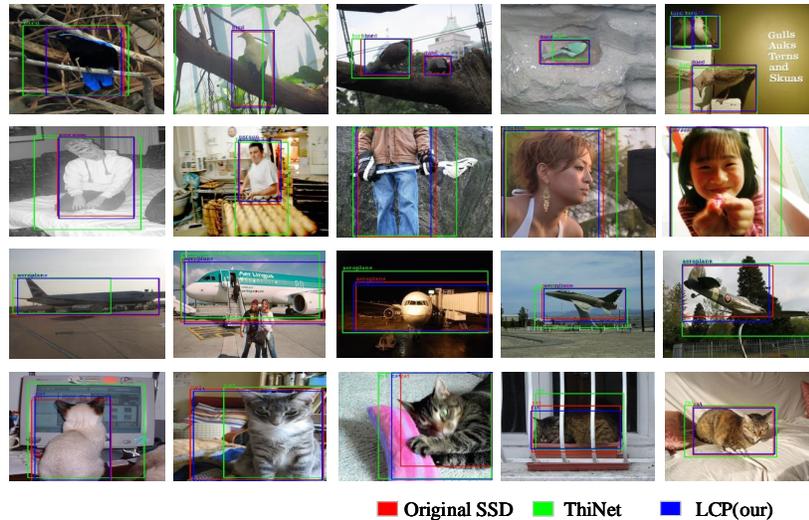
Figure 8: The predictions of original SSD, models pruned by Thinet and our LCP. We prune the VGG-16 by 75% on PASCAL VOC.

## References

[1] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: Advances in neural information processing systems, 2012, pp. 1097–1105.

[2] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., Imagenet large scale visual recognition challenge, International journal of computer vision 115 (3) (2015) 211–252.

[3] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.

[4] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, A. C. Berg, Ssd: Single shot multibox detector, in: European conference on computer vision, Springer, 2016, pp. 21–37.

[5] S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, in: Advances in neural information processing systems, 2015, pp. 91–99.

[6] H. Noh, S. Hong, B. Han, Learning deconvolution network for semantic segmentation, in: Proceedings of the IEEE international conference on computer vision, 2015, pp. 1520–1528.

[7] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, L. Van Gool, Temporal segment networks: Towards good practices for deep action recognition, in: European conference on computer vision, Springer, 2016, pp. 20–36.

[8] Y. He, X. Zhang, J. Sun, Channel pruning for accelerating very deep neural networks, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 1389–1397.

[9] J.-H. Luo, J. Wu, W. Lin, Thinet: A filter level pruning method for deep neural network compression, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 5058–5066.

[10] C. Jiang, G. Li, C. Qian, K. Tang, Efficient dnn neuron pruning by minimizing layer-wise nonlinear reconstruction error., in: IJCAI, Vol. 2018, 2018, pp. 2–2.

[11] S. Han, H. Mao, W. J. Dally, Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, arXiv preprint arXiv:1510.00149.

[12] S. Han, J. Pool, J. Tran, W. Dally, Learning both weights and connections for efficient neural network, in: Advances in neural information processing systems, 2015, pp. 1135–1143.

[13] Z. Zhuang, M. Tan, B. Zhuang, J. Liu, Y. Guo, Q. Wu, J. Huang, J. Zhu, Discrimination-aware channel pruning for deep neural networks, in: Advances in Neural Information Processing Systems, 2018, pp. 875–886.

[14] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, Y. Bengio, Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1, arXiv preprint arXiv:1602.02830.

[15] F. Li, B. Zhang, B. Liu, Ternary weight networks, arXiv preprint arXiv:1605.04711.

[16] A. Zhou, A. Yao, Y. Guo, L. Xu, Y. Chen, Incremental network quantization: Towards lossless cnns with low-precision weights, arXiv preprint arXiv:1702.03044.

[17] Y. Yu, J. Wu, L. Huang, Double quantization for communication-efficient distributed optimization, in: Advances in Neural Information Processing Systems, 2019, pp. 4440–4451.

[18] Y. Zhao, X. Gao, D. Bates, R. Mullins, C.-Z. Xu, Focused quantization for sparse cnns, in: Advances in Neural Information Processing Systems, 2019, pp. 5585–5594.

[19] W. Wen, C. Wu, Y. Wang, Y. Chen, H. Li, Learning structured sparsity in deep neural networks, in: Advances in neural information processing systems, 2016, pp. 2074–2082.

[20] Y. Guo, A. Yao, Y. Chen, Dynamic network surgery for efficient dnns, in: Advances In Neural Information Processing Systems, 2016, pp. 1379–1387.

[21] X. Jin, X. Yuan, J. Feng, S. Yan, Training skinny deep neural networks with iterative hard thresholding methods, arXiv preprint arXiv:1607.05423.

[22] H. Kim, M. U. K. Khan, C.-M. Kyung, Efficient neural network compression, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[23] X. Liu, W. Li, J. Huo, L. Y. Y. Gao, Layerwise sparse coding for pruned deep neural networks with extreme compression ratio.

[24] Y. LeCun, J. S. Denker, S. A. Solla, Optimal brain damage, in: Advances in neural information processing systems, 1990, pp. 598–605.

[25] B. Hassibi, D. G. Stork, Second order derivatives for network pruning: Optimal brain surgeon, in: Advances in neural information processing systems, 1993, pp. 164–171.

[26] H. Li, A. Kadav, I. Durdanovic, H. Samet, H. P. Graf, Pruning filters for efficient convnets, arXiv preprint arXiv:1608.08710.

[27] H. Hu, R. Peng, Y.-W. Tai, C.-K. Tang, Network trimming: A data-driven neuron pruning approach towards efficient deep architectures, arXiv preprint arXiv:1607.03250.

[28] J. M. Alvarez, M. Salzmann, Learning the number of neurons in deep networks, in: Advances in Neural Information Processing Systems, 2016, pp. 2270–2278.

[29] Z. Liu, J. Li, Z. Shen, G. Huang, S. Yan, C. Zhang, Learning efficient convolutional networks through network slimming, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 2736–2744.

[30] P. Singh, V. K. Verma, P. Rai, V. P. Namboodiri, Play and prune: Adaptive filter pruning for deep model compression, arXiv preprint arXiv:1905.04446.

[31] Z. Liu, H. Mu, X. Zhang, Z. Guo, X. Yang, K.-T. Cheng, J. Sun, Metapruning: Meta learning for automatic neural network channel pruning, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 3296–3305.

[32] X. Xiao, Z. Wang, S. Rajasekaran, Autoprune: Automatic network pruning by regularizing auxiliary parameters, in: Advances in Neural Information Processing Systems, 2019, pp. 13681–13691.

[33] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, S. Belongie, Feature pyramid networks for object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 2117–2125.

[34] Z. Cai, N. Vasconcelos, Cascade r-cnn: Delving into high quality object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 6154–6162.

[35] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask r-cnn, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2961–2969.

[36] Z. Huang, L. Huang, Y. Gong, C. Huang, X. Wang, Mask scoring r-cnn, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 6409–6418.

[37] T.-Y. Lin, P. Goyal, R. Girshick, K. He, P. Dollár, Focal loss for dense object detection, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2980–2988.

[38] H. Law, J. Deng, Cornernet: Detecting objects as paired keypoints, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 734–750.

[39] X. Zhou, D. Wang, P. Krähenbühl, Objects as points, arXiv preprint arXiv:1904.07850.

[40] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, Q. Tian, Centernet: Keypoint triplets for object detection, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 6569–6578.

[41] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, S. Savarese, Generalized intersection over union: A metric and a loss for bounding box regression, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 658–666.

[42] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes (voc) challenge, International journal of computer vision 88 (2) (2010) 303–338.

[43] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.

[44] D. Hoiem, Y. Chodpathumwan, Q. Dai, Diagnosing error in object detectors, in: European conference on computer vision, Springer, 2012, pp. 340–353.