

Inception Convolution with Efficient Dilation Search

Jie Liu ¹, Chuming Li ², Feng Liang ², Chen Lin ³, Ming Sun ²

Junjie Yan ², Wanli Ouyang ⁴, Dong Xu ⁴

¹Beihang University, ²SenseTime Research, ³University of Oxford, ⁴The University of Sydney

ljie@buaa.edu.cn chen.lin@eng.ox.ac.uk

{lichuming, liangfeng, sunming1, yanjunjie}@sensetime.com

{wanli.ouyang, dong.xu}@sydney.edu.au

Abstract

Dilation convolution is a critical mutant of standard convolution neural network to control effective receptive fields and handle large scale variance of objects without introducing additional computation. However, fitting the effective reception field to data with dilated convolution is less discussed in the literature. To fully explore its potentials, we proposed a new mutant of dilated convolution, namely inception (dilated) convolution where the convolutions have independent dilation among different axes, channels and layers. To explore a practical method for fitting the complex inception convolution to the data, a simple while effective dilation search algorithm(EDO) based on statistical optimization is developed. The search method operates in a zero-cost manner which is extremely fast to apply on large scale datasets. Empirical results reveal that our method obtains consistent performance gains in an extensive range of benchmarks. For instance, by simply replace the 3×3 standard convolutions in ResNet-50 backbone with inception convolution, we improve the mAP of Faster-RCNN on MS-COCO from 36.4% to 39.2%. Furthermore, using the same replacement in ResNet-101 backbone, we achieve a huge improvement over AP score from 60.2% to 68.5% on COCO val2017 for the bottom up human pose estimation.

1. Introduction

The receptive field is an important concept of convolution neural network and has been extensively studied. The authors [31] prove that the intensity in a receptive field is roughly a Gaussian distribution and only few pixels around the central part of the receptive field effectively contribute to the response of the output neuron. Furthermore, a more consciously defined effective receptive field (ERF) has been tested for different tasks in previous works [25, 34].

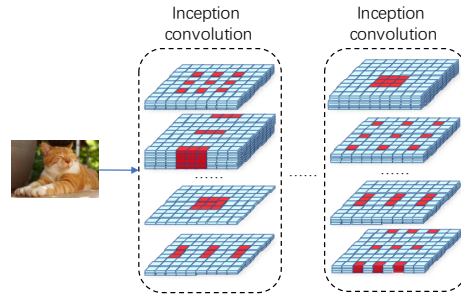


Figure 1. Illustration of inception convolution

The requirements of ERF vary in different tasks due to the size variance of the input images and the range of interesting object scales. For instance, in image classification, the input sizes tend to be small(e.g. 224×224), while in object detection, the input sizes are much bigger and the objects bear a large range of scales. Even for the same task with a fixed network, the optimal ERF for a certain convolution layer may be different from the standard convolution operations as discussed in [25, 33]. The varying requirements of ERF imply the necessity of a general and practical ERF optimization algorithm for different tasks.

As discussed in [31], dilation of the dilated convolution kernels is a highly effective hyper-parameter to adjust the distribution of ERFs among different tasks. The work in [25] proposed to assign different dilation values at different stages of a CNN and achieves consistent improvements. NATS [34] steps further and divide a convolution into different groups with each having independent dilation values.

However, they apply the skeleton network architecture search methods in relatively coarse search spaces, which neglect the fine-grained inner structure of dilated convolution. Therefore, in this work, we focus on exploring the search problem in the dilation-domain to efficiently adjust ERFs.

First of all, we would like to have a more flexible search space compared with [25]. Flexibility produces the power of fitting ERFs to different datasets. We propose a new mutant of dilated convolution, namely Inception Convolution.

tion, which contains as much as possible dilation patterns as shown in Figure 1. In the space of Inception Convolution, the dilation of each axis, each channel, and each convolution layer is independently defined. The inception convolution provides a dense range of possible ERFs. We further study the impact of inception convolution over search results in ablation study.

For optimization in our search space, we refer to the proliferating works in neural architecture search(NAS), which enables automatic optimization of the combination of neural network operators. DARTS and single path one-shot(SPOS) are two main families of NAS methods. DARTS trains a supernet where the discrete operation selection is relaxed into a continuous weighted sum of all candidate operations’ output. After training, in each block, the operation with the largest architecture weight is chosen. SPOS randomly selects an operation sequence(subnet) in each training step of the supernet and the same operation in different sequences share the same weights. After training, SPOS selects the best operation sequence via sampling and evaluation of multiple sequences inheriting the shared-weights.

However, both DARTS and SPOS are not suitable for our search space. In DARTS, all operations in a block are applied to the input during training to make the architecture weights aware of each operation’s importance, but the number of dilation patterns for a convolution layer (block) is large, i.e. 16 if the two axes each has 4 choices. It means DARTS needs 16 sequential calculations, thus has low GPU utility and big computational cost. SPOS samples operation sequences during training. However, in our search space, the number of dilation patterns even in a single convolution layer is huge, i.e. d_{max}^{2C} where C is the channel number and d_{max} is the maximum dilation. The huge number of dilation patterns pose extreme difficulty for designing a fair sampling strategy for SPOS.

In this paper, we propose a simple yet efficient dilation optimization algorithm(EDO). In EDO, each layer of the supernet is a standard convolution operation whose kernel covers all possible dilation patterns. After the pre-training of the supernet, we select the dilation pattern for each channel in each convolution layer via solving a statistical optimization problem. For each layer, the selection is solved with the pre-trained weights, via the minimization of the L_1 error between the expectation of the output of (1) the original convolution layer and (2) the cropped out dilation convolution with the selected dilation patterns.

EDO supports efficient channel-wise dilation optimization over our complete dilation-domain search space in a quite simple way. Compared with the search based method [15], the search cost of our methods is almost zero since the only cost is calculating the statistics of trained weights. Compared with differential methods [29, 3], it converts se-

quential calculation of different patterns into a parallel way, thus has lower computation cost and higher GPU utilization. Further, compared with SPOS, we don’t need to design complex mechanisms to guarantee the fairness of the sampling or the accurate ranking of the subnets.

Our contributions are summarized as:

- We proposed a new mutant of dilated convolution, namely Inception convolution, which can fit ERF to data efficiently.
- We propose a zero-cost statistic based architecture search algorithm (EDO) for Inception Convolution to fit the architecture to data at no cost.

Empirical results reveal that EDO achieves general improvements on an extensible range of tasks and models without any additional computation budget. On the ImageNet dataset, we outperform ResNet-50 by 1.1%. On the COCO dataset, our IC-ResNet50 and IC-ResNet101 can achieve 38.9% and 41.9% AP using Faster-RCNN-FPN, outperforming the baseline by 2.5% and 3.1% respectively. For bottom-up human pose estimation, we outperform ResNet-101 by 8.3% on COCO val2017. Furthermore, our inception convolution also yields a significant improvement in instance segmentation and crowd human detection.

2. Related Work

2.1. Receptive Field

We shall say we are inspired by the Inception module at some places (because it is used in the title) and state the difference. The Receptive Field is a crucial issue in Convolutional Neural Networks(CNNs). Traditional CNNs [20, 40, 18] stack multiple convolutional layers to enhance the receptive field. The Inception networks [42, 43, 41] introduce different size filters operating on the same level to aggregate different receptive fields. STN [19] presents a spatial transformer module that can actively transform a feature map by producing an appropriate transformation for each input sample. Deformable convolution [10] predicts sampling offsets with respect to the preceding feature maps to adjust the receptive fields automatically. Scale-Adaptive Convolutions [48] are proposed to predict local flexible-size dilations at each position to cover objects of various sizes. However, while effective in performance, these methods are unfriendly to hardware optimization in real-time applications.

Dilated(Atrous) convolution [46, 8] affects the receptive field by performing convolution at sparsely sampled locations, which have been widely used in semantic segmentation [7, 49], object detection [24, 23, 33]. PSConv [21] manually mixes up a spectrum of dilation rates in one convolution, which is proven to be sub-optimal in our exper-

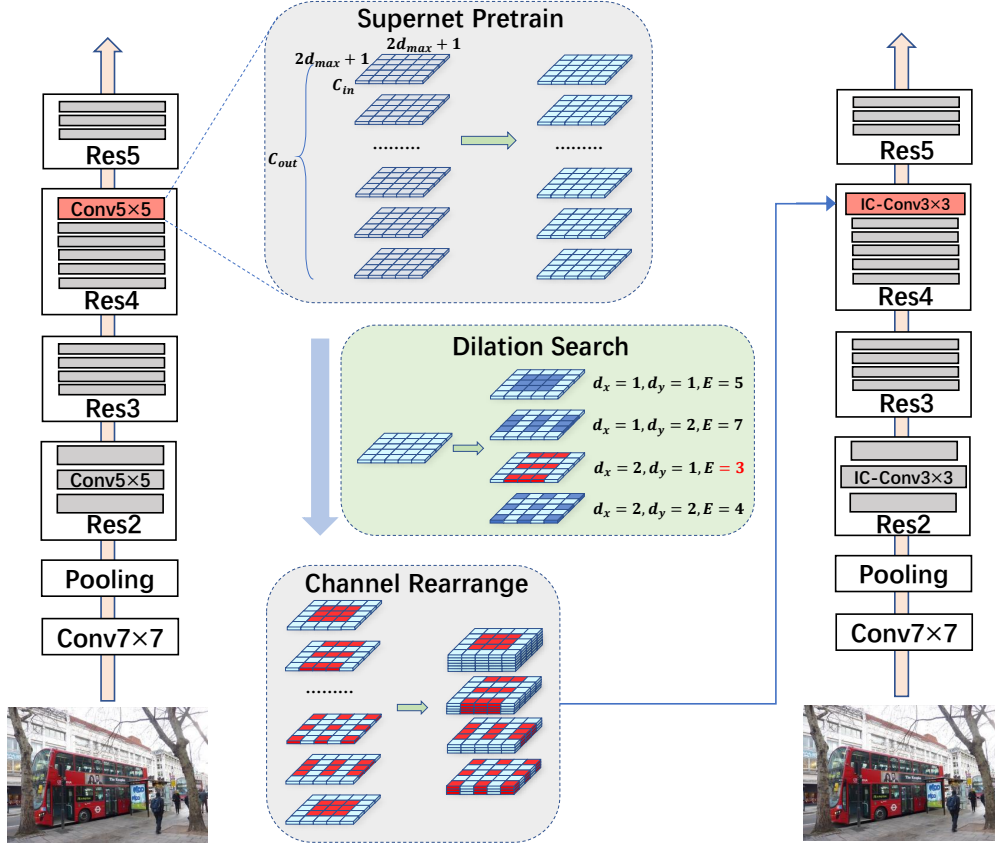


Figure 2. An overview of our EDO algorithm. Let’s take ResNet50 for example. Firstly, we get a pretrained R50 with $(2d_{max} + 1) \times (2d_{max} + 1)$ convolution as the bottleneck (In this figure, the kernel size of the supernet is 5x5 ($d_{max} = 2$)). Secondly, for each output filter in the convolution operation, we compute the L1 error of the expected output with all possible dilation patterns and select the dilation pattern of this filter with minimal expectation (dilation (2,1) with $E = 3$ in the above case). Finally, We rearrange filters so that the kernel of same dilation is arranged together, and produce our Inception Convolution.

iments. We aim at searching for efficient inception convolution, a mixed dilated convolution, which is friendly to hardware optimization.

2.2. Neural Architecture Search

Neural architecture search, the process of automating architecture design, has attracted a lot of attention. Early NAS approaches [50, 51] are computationally expensive due to the evaluating of each candidate. To reduce the searching costs, ENAS [35] introduces weight sharing strategy between the candidate architectures. Recently, One-shot NAS methods [1, 29, 3, 2, 22, 14, 30] build a directed acyclic graph G (a.k.a. supernet) to subsume all architectures in the search space to further reduce the cost. More relevant to our work, there are some NAS works searching dilation convolutions. NAT [34] utilizes DARTS [29] to search for dilated rate on group level in the CNN backbone. CRNAS [25] searches for different dilation rates for different building blocks by using a single path one-shot NAS [15]. However, the search spaces of the aforementioned methods are

limited. Moreover, simply combining the existing searching method, like DARTS [29] or SPOS [15] cannot handle the extremely large number of possible operations in a single convolution layer.

3. Method

3.1. Problem Formulation

To fully explore the flexibility of dilation in ERF fitting, we consider a complete dilation domain, namely Inception Convolution. An Inception Convolution has independent dilation for the two axes in each channel and is formally represented by:

$$\mathcal{d} = \{\mathbf{d}_x^i, \mathbf{d}_y^i | \mathbf{d}_x^i, \mathbf{d}_y^i \in 1, 2, \dots, d_{max}, i \in \{1, 2, \dots, C^{out}\}\}, \quad (1)$$

where \mathbf{d}_x^i and \mathbf{d}_y^i are the dilation in x and y axis of the filter at the i -th output channel, ranging from 1 to d_{max} , and C^{out} denotes the output channel numbers. The candidate structure number of a single Inception Convolution is $d_{max}^{2C^{out}}$. In this paper, we aim to develop an algorithm to ef-

ficiently fit the ERF among different tasks via the selection of set \mathbf{d} (The visualization of ERF can be seen in Supplementary).

3.2. Solution

Recently, NAS evolves as an effective way to generate high-performance architectures in specified search spaces. DARTS and SPOS are two mainstream families of NAS methods. However, as Inception Convolution contains d_{max}^2 dilation patterns and d_{max}^{2C} candidates, both DARTS and SPOS are not expandable enough for efficient search in our domain as discussed in Section 1.

Recall that DARTS alternately trains architecture weights and operation weights on two separate datasets, and use architecture weights to indicate the importance of the corresponding operations. Although the resulted architecture weights are likely to distribute uniform among operations, operations obtained in this way still consist of a good network. It reveals that the weights in a pre-trained supernet are informative to guide the operation selection. In this work, we follow this idea and formulate a statistical optimization problem to select dilations based on the corresponding weights in a pre-trained supernet. The resulted algorithm, EDO, is simple, effective, and efficient, without the need for a careful hyper-parameter tuning.

Supernet Given a network architecture and a task it is required to fit on, our supernet reserves its architecture while changes the kernel sizes to cover all candidate dilation patterns. Formally, for a convolution layer with kernel size $2k+1$, we replace it in the supernet with $2kd_{max}+1$, which is the maximum width and height for all candidate dilation patterns. The supernet is pre-trained on the given task.

Statistical Optimization For each convolution layer with weights $\mathbf{W} \in R^{C_{out} \times C_{in} \times (2kd_{max}+1) \times (2kd_{max}+1)}$, we define $\mathbf{W}^i \in R^{C_{in} \times (2kd_{max}+1) \times (2kd_{max}+1)}$ as the weights of the i -th convolution filter expanded in our supernet, and $\mathbf{W}_{\mathbf{d}_x^i, \mathbf{d}_y^i}^i \in R^{C_{in} \times (2k+1) \times (2k+1)}$ denotes the dilation convolution filter cropped from \mathbf{W}^i with the positions determined by \mathbf{d}_x^i and \mathbf{d}_y^i . We formulate the dilation selection as an optimization problem where the L_1 error between the expectation of the output of the pre-trained expanded weights \mathbf{W} and the cropped dilation weights $\mathbf{W}_{\mathbf{d}}$ is minimized, formally expressed as:

$$\min_{\mathbf{d}} \|E[\mathbf{W}\mathbf{X}] - E[\mathbf{W}_{\mathbf{d}}\mathbf{X}]\|_1, \quad (2)$$

$$s.t. \quad \mathbf{d}_x^i, \mathbf{d}_y^i \in \{0, 1, \dots, d_{max}\}. \quad (3)$$

Here $\mathbf{X} \in R^{B \times C_{in} \times H \times W}$ is the input of this convolution layer with batch size B , height H and width W . As \mathbf{W} and $\mathbf{W}_{\mathbf{d}}$ are independent of \mathbf{X} , the target of optimization is

further expressed as:

$$\|E[\mathbf{W}\mathbf{X}] - E[\mathbf{W}_{\mathbf{d}}\mathbf{X}]\|_1 = \|\mathbf{W}E[\mathbf{X}] - \mathbf{W}_{\mathbf{d}}E[\mathbf{X}]\|_1, \quad (4)$$

$$= \|(\mathbf{W} - \mathbf{W}_{\mathbf{d}})E[\mathbf{X}]\|_1. \quad (5)$$

Accurate solution to the above problem involves arduous computation to average \mathbf{X} over the whole training dataset. Considering batch normalization is currently a common element of CNNs, we introduce an assumption that \mathbf{X} has passed a batch normalization operation and the γ and β for each channel of \mathbf{X} has the same value due to γ and β of different input channels are identical in different runs of training. Along with the shift invariance among H and W and the permutation invariance among B , we derive that \mathbf{X} has the same distribution among all positions, and the target is simplified more:

$$\|(\mathbf{W} - \mathbf{W}_{\mathbf{d}})E[\mathbf{X}]\|_1 = E\|(\mathbf{W} - \mathbf{W}_{\mathbf{d}})\mathbf{1}\|_1, \quad (6)$$

$$= E \sum_1^{C_{out}} \|(\mathbf{W}^i - \mathbf{W}_{\mathbf{d}_x^i, \mathbf{d}_y^i}^i)\mathbf{1}\|_1. \quad (7)$$

Where E is the expectation of all positions in \mathbf{X} and $\mathbf{1}$ is an all-ones matrix with the same shape as \mathbf{X} . With the above deductions, the dilation selections \mathbf{d} of a certain convolution can be easily obtained, via independently traversing all dilation patterns $(\mathbf{d}_x^i, \mathbf{d}_y^i)$ for each filter \mathbf{W}^i along with little cost to calculate Eq. 7.

3.3. Discussion

The Relationship with DARTS An intuitive application of DARTS on our inception convolution is the way introduced in [34]. In this way d_{max}^2 operations are calculated sequentially, with the total cost $C_{in} \times C_{out} \times (2k+1)^2 \times d_{max}^2$. However, in our EDO algorithm, the d_{max}^2 operations are calculated in parallel, with the total cost $C_{in} \times C_{out} \times (2kd_{max}+1)^2$. For most CNNs, where k is usually 1, $C_{in} \times C_{out} \times (2kd_{max}+1)^2$ is only 56% of $C_{in} \times C_{out} \times (2k+1)^2 \times d_{max}^2$ with d_{max} set 4. Therefore, EDO is more computationally efficient than DARTS.

Additionally, as proved in [47], DARTS degenerates into random sampling in some cases because the principal eigenvalue of the hessian matrix of the architecture parameters α appears large. However, we direct define the statistical optimization problem over the pre-trained network weights rather than introducing the architecture parameters α which are not robust.

The Relationship with NATS and CRNAS NATS and CRNAS also propose search space powered with flexible dilations, while they are less complete than our inception convolution. CRNAS search dilation independently for each

stage, thus it is based on SPOS. NATS divides a convolution into groups and search among a few dilation patterns(usually 5 patterns) for each group with DARTS. Inception convolution is channel-wise and contains all dilation patterns(at least 16) under the max dilation d_{max} . The pipeline of our proposed method is shown in Figure 2.

4. Experiment

4.1. Image Recognition

4.1.1 Dataset and implementation details

For image recognition, we evaluate our method on the ImageNet dataset [37] with 1.28M training images and 50k validation images. We first train our supernet with the largest kernel size (9) following the standard training procedure in [18]. More specifically, we use stochastic gradient descent (SGD) as an optimizer with 0.9 momentum and 0.0001 weight decay. The supernet is trained for 100 epochs with the batch size 1024 without any tricks. We adopt the cosine learning rate scheduler with the initial learning rate of 0.4. Then we perform EDO to obtain the best inception convolutions, as described in Section. 3. The resulted IC-Net is retrained using the same procedure as the supernet training.

4.1.2 Main results

We search inception convolution on various types of networks, from MobileNet-V2 [38] to ResNeXt [45]. As shown in Table 1, the searched inception convolution consistently boosts the ImageNet performance. IC-ResNet18 and IC-ResNet50 outperform the baselines by 1.07% and 1.11% respectively. Inception convolution is also compatible with the networks consisting of depth-wise convolution or group convolution. For instance, inception convolution obtains 1.21% and 0.62% gain respectively on MobileNet-V2 and ResNeXt-101.

Table 1. ImageNet top-1/top-5 accuracy comparisons on the validation set. All the models are measured using single center crop evaluation. **The baseline results are re-implemented by ourselves in the same code with IC-Net.**

Architecture	Conv Type	Top-1/5 Acc.(%)
MobileNet-V2 [38]	standard	70.71 / 89.81
	IC-Conv	71.92 / 90.54
ResNet-18 [18]	standard	70.67 / 89.74
	IC-Conv	71.74 / 90.91
ResNet-50 [18]	standard	76.19 / 92.93
	IC-Conv	77.30 / 93.58
ResNeXt-101 (32x4d) [45]	standard	78.71 / 94.20
	IC-Conv	79.33 / 94.74

4.2. Object Detection

4.2.1 Dataset and implementation details

In the following experiments, unless otherwise stated, we will only replace the convolutions with our inception convolutions in the backbone. For object detection, we use the MS-COCO [28] for the experiment. The dataset is widely believed challenging in particular due to the huge variation of object scales and a large number of objects per image. The supernet with the largest (9) kernel size is used as the pretrain model to generate the inception convolutions. For detector training, we use stochastic gradient descent (SGD) as an optimizer with 0.9 momentum and 0.0001 weight decay. The model is trained for 13 epochs, known as 1× schedule [13]. We use multi-GPU training over 8 1080TI GPUs with a total batch size of 16. The initial learning rate is 0.00125 per image and is divided by 10 at 8 and 11 epochs. Warm-up is adopted for both baselines and our searched models.

4.2.2 Main results

Our searched inception convolution shows great potential on various types of detectors, from Faster RCNN [36] to Cascade RCNN [4]. The same type of detectors is trained using exactly the same 1× training procedure [13]. We replace all the 3×3 convolutions in the pretrained backbone network by IC-Conv, while the convolutions in the FPN neck are kept as standard convolutions. As shown in Table 3, our searched inception convolution boosts the COCO performance on an extensive range of backbones. For Faster RCNN [36, 26], our IC-ResNet50, IC-ResNet101 and IC-ResNeXt101-32x4d outperforms the baseline by a large 2.5%, 3.1% and 1.6% margin respectively. For a more powerful Cascade RCNN [4], our inception convolution is also compatible. For instance, we obtain 45.7% AP using IC-ResNeXt101-32x4d which is 1.3% higher than the baseline. Our inception convolution is especially effective for large objects(4.1%, 4.1% and 2.4% improvement for AP_l), mainly benefiting from the large receptive field provided by effective dilation.

4.2.3 Transferability Verification

Different detector We transfer our searched inception convolution backbone to various types of object detectors, including one stage detector RetinaNet [27], anchor box free detector FCOS [44], NAS powered NAS-FPN [12], transformer based detector DETR [6]. The experimental results on COCO *minival* are shown in Table 2. Equipped with our IC-ResNet50, we can obtain a average 1.8% COCO AP gain without additional FLOPS cost. In particular, our inception convolution is compatible with the new transformer based detector DETR. Our IC-R50 DETR

Table 2. Transfer inception convolution to different detectors. AP(%) is reported on COCO *minival*. Our inception convolution (denoted by ‘IC-x’) achieves consistent improvements. For DETR, we use the official released 150 epochs training scripts rather than 500 epochs in the original paper, due to the limited computation resource.

Detector	Backbone	Conv Type	AP
Faster-RCNN-C4 [36]	ResNet-50	standard	35.0
		IC-Conv	38.5 _(+3.5)
RetinaNet [27]	ResNet-50	standard	36.0
		IC-Conv	37.9 _(+1.9)
DETR [6]	ResNet-50	standard	39.7
		IC-Conv	40.7 _(+1.0)
FCOS [44]	ResNet-50	standard	37.2
		IC-Conv	38.8 _(+1.6)
Faster-RCNN-NASFPN [12]	ResNet-50	standard	40.2
		IC-Conv	41.1 _(+0.9)

can obtain 40.7% COCO AP which is 1.0% higher than the vanilla R50 DETR.

Different dataset We also transfer our inception convolution backbone to another object detection dataset VOC [11] by using RFCN framework [9]. As shown in Table 4, our IC-ResNet50 and IC-ResNet101 achieves AP₅₀ improvement by 1.94% and 1.59% comparing with the already high baseline. Training details can be seen in Supplementary.

4.3. Instance Segmentation

4.3.1 Dataset and implementation details

Segmentation is another task that is highly sensitive to the ERF [16]. Therefore, we further search for inception convolution in instance segmentation task by using the Mask RCNN [17] framework. We use the MS-COCO [28] as a dataset for instance segmentation experiments. We use the whole COCO *trainval135* as training set and validate on COCO *minival*. The supernet with the largest (9) kernel size is also used as the pretrain model to generate the inception convolutions. For Mask RCNN training, we use the same training configurations as Section. 4.2.1.

4.3.2 Main results

We search inception convolution on Mask RCNN [17] and transfer it to Cascade Mask RCNN [5]. As shown in Table 5, the searched inception convolution consistently boosts the instance segmentation performance. For Mask RCNN, our IC-ResNet50, IC-ResNet101 and IC-ResNeXt101-32x4d outperform the baselines by 2.8%, 2.8% and 2.0% respectively. In particular, for a stronger Cascade Mask RCNN with ResNeXt101-32x4d, our IC-ResNeXt101-32x4d can further improve by a 1.5% margin.

4.4. Crowd Human Detection

Crowd human detection is also a challenging computer vision task that requires fine ERF. In crowd scenarios, because different people have large variations of poses and scales, a more strong ERF is called for correctly perceive neighborhood features. Furthermore, the detection bounding box of the human body is usually a vertical rectangle while the ERF of a standard convolution is square. It implies that an ERF of vertical rectangle shape is more suitable for human detection.

4.4.1 Dataset and implementation details

We use CrowdHuman [39], which contains 15K images for training, as a dataset to evaluate our inception convolution. RFCN with R50 is adopted as a framework and the supernet kernel size is also 9. For RFCN training, we use SGD as an optimizer with 0.9 momentum and 0.0001 weight decay. The model is trained for 20 epochs and the learning rate is divided by 10 at 10 and 15 epochs.

4.4.2 Main results

As shown in Table 6, inception convolution also outperforms its baseline on CrowdHuman by 0.8 % MR^{-2} , which is the most important metric on CrowdHuman, implying its generality over different detection tasks. We conduct further visualizations and plot the proportion of square filters($d_x^i = d_y^i$), vertical filters($d_x^i > d_y^i$) and horizontal filters($d_x^i < d_y^i$) separately in Figure 3. It is observed that, at the former layers, vertical and horizontal filters have comparable amounts. While in the latter layers, where the resolution is low and a kernel can cover a human shape, **vertical filters have a larger amount than horizontal filters**. It reveals how EDO optimizes the dilation for such a specific task.

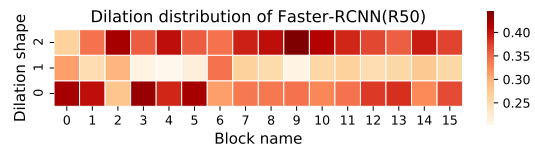


Figure 3. The dilation’s distribution of our inception convolution on CrowdHuman dataset with Faster-RCNN-C4 (R50). For the vertical coordinates, zero represents horizontal filters (e.g. (1,4),(2,3)), 1 represents squares filters (e.g., (1,1),(3,3)) and 2 represents vertical filters (e.g., (4,1),(3,2)). Vertical filters are more than horizontal ones at the high layers where a kernel can cover the human body shape. It implies EDO fits the shape of ERF to certain object categories.

4.5. Human Pose Estimation

Given an input image, human pose estimation aims to detect the locations of keypoints or parts (e.g., elbow, wrist,

Table 3. Detection performance AP (%) on COCO *minival* for different backbones using our inception convolution (denoted by ‘IC-x’)

Detector	Backbone	Conv Type	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
Faster-RCNN [36]	R50	standard	36.4	58.6	39.2	21.7	40.2	46.4
		IC-Conv	38.9(+2.5)	61.6	41.8	22.9	42.3	50.5
	R101	standard	38.8	60.9	42.1	22.6	42.9	50.5
		IC-Conv	41.9(+3.1)	64.2	45.3	25.5	45.8	54.6
	X101-32x4d	standard	40.5	63.1	44.4	24.9	44.8	52.0
		IC-Conv	42.1(+1.6)	64.7	45.7	25.5	46.1	54.4
Cascade-RCNN [4]	R50	standard	40.5	59.2	44.0	22.5	43.9	53.6
		IC-Conv	42.4(+1.9)	62.0	46.0	25.2	45.9	56.0
	R101	standard	42.6	60.9	46.2	23.8	46.2	56.9
		IC-Conv	45.0(+2.4)	64.8	48.7	26.9	49.0	59.6
	X101-32x4d	standard	44.4	63.6	48.5	25.8	48.5	58.1
		IC-Conv	45.7(+1.3)	65.5	49.7	27.1	49.8	59.8

Table 4. RFCN detection performance on VOC test2007. Our inception convolution models are denoted by ‘IC-’.

	R50	IC-R50	R101	IC-R101
AP ₅₀	79.66	81.60	81.36	82.95

etc). To fully validate the superiority of our inception convolution, we chose to solve a more difficult multi-person pose estimation problems with the bottom-up approach where all the human parts in an image should be detected and the keypoints of the same person should be associated.

4.5.1 Dataset and implementation details

We train our model on COCO train2017 dataset, including 57K images and 150K person instances. We evaluate our approach on the val2017 set containing 5000 images. The supernet with the largest (13) kernel size is used as the pre-train model to generate the inception convolutions. Following the setting in MMPose¹, we adopt associative embedding [32] as the method for bottom-up human pose estimation and use the Adam optimizer. The base learning rate is set as $1e-3$, and is dropped to $1e-4$ and $1e-5$ at the 200th and 260th epochs, respectively. The training process is terminated within 300 epochs.

4.5.2 Main results

We report the results of our inception convolution and standard convolution in Table 7. Our approach is about **10%** higher than the baseline with the same input size, which is a **huge** improvement. Furthermore, an AP of more than 60 enables the resnet backbone, which is widely used and optimized for detection and classification tasks, to be reapplied to human pose estimation in industrial production. An unified backbone favors model deployment.

¹<https://github.com/open-mmlab/mmpose>

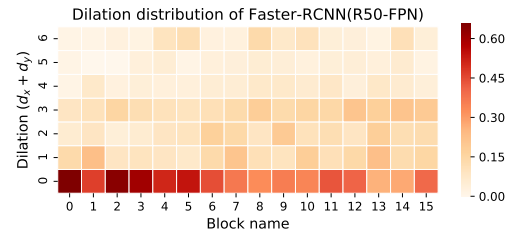


Figure 4. The dilation’s distribution of our inception convolution on COCO dataset with Faster-RCNN (R50-FPN). For visualization, we set the vertical axis as the sum of dilation in the x and y directions, and the horizontal axis represents the number of blocks.

4.6. Comparisons with the State-of-the-art Methods

According to the results available in other papers, we choose Faster-RCNN and Mask-RCNN for comparison. Table 8 gives a quantitative comparison of our method with POD [33], NATS [34] and PSConv [21]. Remarkably, our method achieves better performance than all the other methods on different backbones and detectors. Our proposed method has a larger channel-wise search space and a more efficient search algorithm compared with other approaches which we believe is the reason for the improved detection accuracy.

4.7. Ablation Study

4.7.1 Dilation Visualization

In this section, we visualize the distribution of dilations searched for Inception convolution using our EDO algorithm. As shown in Figure 4 and 5, $(d_x, d_y) = (1, 1)$ is the main part, which means that standard convolution is the skeleton part of our inception convolution and extracts the main information from feature, while dilations of other sizes are used to complement various scale information.

Table 5. Mask RCNN/ Cascade Mask RCNN detection and instance segmentation performance on COCO *minival* for different backbones using our inception convolution (denoted by ‘IC-x’). Box and mask are the AP(%) of the bounding box and segmentation results respectively.

Detector	Backbone	Conv Type	box AP						mask AP					
			AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
Mask-RCNN [17]	R50	standard	37.2	59.0	40.1	22.3	41.2	47.7	33.8	55.7	35.8	17.9	37.6	45.9
		IC-Conv	40.0 _(+2.8)	62.1	43.1	23.5	43.7	52.1	35.9 _(+1.9)	58.4	37.9	18.9	39.5	49.5
	R101	standard	39.8	61.8	43.5	23.2	43.9	51.9	35.6	58.1	38.1	18.4	39.5	49.0
		IC-Conv	42.6 _(+2.8)	64.6	46.4	25.6	46.6	55.4	37.9 _(+2.3)	61.2	40.3	20.5	41.7	52.0
	X101-32x4d	standard	41.4	63.6	45.2	24.6	45.9	53.2	37.1	60.1	39.6	19.4	41.3	50.9
		IC-Conv	43.4 _(+2.0)	65.7	47.4	27.2	47.2	56.3	38.4 _(+1.3)	62.1	40.4	21.3	42.0	53.0
Cascade-RCNN [5]	R50	standard	41.2	59.7	44.8	23.4	44.6	54.7	35	56.5	37.4	18.0	38.3	48.5
		IC-Conv	43.4 _(+2.2)	62.5	47.0	25.3	47.1	57.1	36.8 _(+1.8)	59.2	39.2	19.4	40.1	50.8
	R101	standard	43.1	61.9	46.8	24.8	47.0	56.7	37.0	59.0	39.6	19.0	40.7	50.8
		IC-Conv	45.7 _(+2.6)	65.2	49.8	26.8	49.6	61.0	38.7 _(+1.7)	61.9	41.3	20.5	42.2	53.9
	X101-32x4d	standard	44.9	64.1	48.9	26.1	48.3	59.4	37.9	60.6	40.6	20.0	41.2	52.6
		IC-Conv	46.4 _(+1.5)	66.0	50.5	27.1	50.3	61.0	39.1 _(+1.2)	62.6	41.6	20.6	42.7	53.7

Table 6. Faster RCNN-C4 detection performance on *CrowdHuman* for ResNet50 backbone using our inception convolution (denoted by ‘IC-x’). In Crowd human detection, MR^{-2} is the most critical indicator, and the lower the value, the better the performance.

Conv Type	Recall	AP	MR^{-2}
standard	79.29	75.61	59.60
IC-Conv	79.32	75.68	58.82

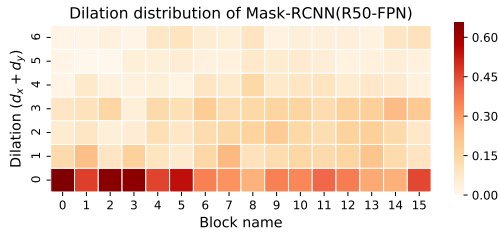


Figure 5. The dilation’s distribution of our inception convolution on COCO dataset with Mask-RCNN (R50-FPN).

4.7.2 Detailed Validation

In this section, we validate the effectiveness of both inception convolution and EDO by Faster-RCNN [36] using 1x scheduler. The supernet uses kernel size 9. As summarized in Table 9, Res50 indicates the vanilla ResNet. Res50_k9 denotes replacing the 3x3 convolution in ResNet50 with 9x9 standard convolution. Res50_IC_1d is the one dimension version of inception convolution, i.e. the dilation rate keeps the same across two directions. Res50_IC_bw is the block version of inception convolution i.e. the dilating rate keeps the same across all channels in a convolution layer. As shown in Table 9, even with the simple 1d inception convolution, we can obtain a higher COCO AP than ResNet50_k9 with 9 times less computation cost. It reveals that the receptive field aggregation in inception convolution is quite efficient and effective. Res50_IC_uni indicates the channel number of each dilation group keeps the same.

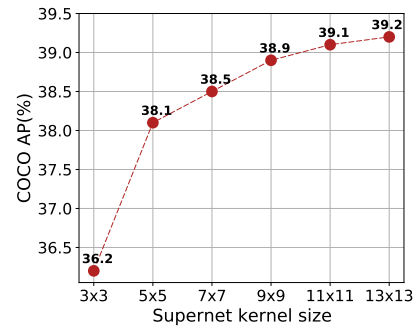


Figure 6. Detection performance on COCO when using different supernet kernel size.

Res50_IC_1e and Res50_IC_1l denote two kinds of targets of statistical optimization to search the dilation, i.e. L_1 error of expectations and expectations of L_1 error. Comparing with Res50_1e and Res50_IC_uni reveals our statistics method is quite effective.

4.7.3 Supernet initial kernel

The kernel size of the supernet is a crucial part of our method. Larger kernel size introduces more dilation combination but more computation complexity in the retrain process of the supernet. In this section, we do an ablation study about the supernet kernel size. We conduct all the experiments on COCO. Remarkably, as shown in Figure 6, when the kernel size grows from 3 to 13, the COCO AP grows consistently. This indicates that our method has great potential. By simply increasing the kernel size of the supernet from 9 to 13, all experiments may obtain higher gains.

4.7.4 Searching head/neck for detector

We further extend our method to searching the head and neck part of the detector. Specifically, we search all the 3x3 convolution in FPN [26] and RPN [36]. The ablation study can be found in Table 10.

Table 7. Comparisons with standard convolution on the COCO validation set. Top: w/o multi-scale test. Bottom: w/ multi-scale test.

Method	Backbone	Input Size	Conv Type	AP	AP ₅₀	AP ₇₅	AP _M	AP _L	AR
Associative Embedding [32]	w/o multi-scale test								
	R50	640x640	standard	51.0	78.3	52.3	50.0	52.0	59.3
			IC-Conv	62.2 _(+11.2)	84.6	67.8	55.2	72.2	67.5
	R101	640x640	standard	55.5	80.8	58.2	52.2	60	62.7
			IC-Conv	63.3 _(+7.8)	85.8	69.4	55.9	74.0	68.3
	w/ multi-scale test								
	R50	640x640	standard	55.8	81.0	58.1	56.3	55.3	63.8
			IC-Conv	65.8 _(+10.0)	85.9	71.3	60.3	73.9	70.7
	R101	640x640	standard	60.2	83.8	63.3	58.7	62.4	67.2
			IC-Conv	68.5 _(+8.3)	87.9	74.2	63.6	75.7	73

Table 8. Comparing our inception convolution with other dilation search methods. Values in baseline column are provided in the original paper of the corresponding methodology. The baseline of all methods is basically the same, which shows that our comparison is fair.

Method	Backbone	Conv type	baseline	box AP
Faster-RCNN [36]	R50	POD	36.2	37.9
		NATS	36.4	38.4
		PSCConv	36.4	38.4
		CRNAS	36.4	38.3
		IC-Conv	36.4	38.9
	R101	POD	38.6	40.1
		NATS	38.6	40.4
		PSCConv	38.5	40.9
		CRNAS	38.6	40.2
		IC-Conv	38.8	41.9
	X101-32x4d	NATS	40.5	41.6
		CRNAS	40.6	41.5
		PSCConv	40.1	41.3
		IC-Conv	40.5	42.1
Mask-RCNN [17]	R50	NATS	37.5	39.3
		PSCConv	37.3	39.4
		CRNAS	37.6	39.1
		IC-Conv	37.2	40.0
	R101	PSCConv	39.4	41.6
		CRNAS	39.7	41.5
		IC-Conv	39.8	42.6
	X101-32x4d	PSCConv	41.1	42.4
		IC-Conv	41.4	43.4

5. Conclusion

In this paper, we propose inception convolution, which is obtained by searching channel-wise dilation through our proposed EDO (Efficient Dilation Optimization) algorithm. IC-Conv can effectively allocate the receptive field in a convolution operation and aggregates the information obtained by the receptive field of multiple scales. What's more, IC-Conv along with EDO generalizes on an extensible range of tasks and can be plugged into arbitrary CNN architectures.

Table 9. Detection performance AP(%) on COCO *minival* using different pattern.

Method	COCO AP
Res50	36.4
Res50_k9	38.1
Res50_IC_bw	37.8
Res50_IC_1d	38.3
Res50_IC_uni	37.8
Res50_IC_el	38.3
Res50_IC_le	38.9

Table 10. Searching different parts of detector.

backbone	FPN [26]	RPN [36]	AP
✓			38.9
✓	✓		39.0
✓	✓	✓	39.2

References

- [1] Gabriel Bender, Pieter-Jan Kindermans, Barret Zoph, Vijay Vasudevan, and Quoc Le. Understanding and simplifying one-shot architecture search. In *International Conference on Machine Learning*, pages 550–559, 2018. **3**
- [2] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network and specialize it for efficient deployment. *arXiv preprint arXiv:1908.09791*, 2019. **3**
- [3] Han Cai, Ligeng Zhu, and Song Han. Proxylessnas: Direct neural architecture search on target task and hardware. *arXiv preprint arXiv:1812.00332*, 2018. **2, 3**
- [4] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018. **5, 7**
- [5] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: high quality object detection and instance segmentation. *IEEE*

- Transactions on Pattern Analysis and Machine Intelligence*, 2019. 6, 8
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *arXiv preprint arXiv:2005.12872*, 2020. 5, 6
 - [7] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 2
 - [8] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. 2
 - [9] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016. 6
 - [10] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017. 2
 - [11] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015. 6
 - [12] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Nas-fpn: Learning scalable feature pyramid architecture for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7036–7045, 2019. 5, 6
 - [13] Ross Girshick, Ilija Radosavovic, Georgia Gkioxari, Piotr Dollár, and Kaiming He. Detectron. <https://github.com/facebookresearch/detectron>, 2018. 5
 - [14] Ronghao Guo, Chen Lin, Chuming Li, Keyu Tian, Ming Sun, Lu Sheng, and Junjie Yan. Powering one-shot topological nas with stabilized share-parameter proxy. *arXiv preprint arXiv:2005.10511*, 2020. 3
 - [15] Zichao Guo, Xiangyu Zhang, Haoyuan Mu, Wen Heng, Zechun Liu, Yichen Wei, and Jian Sun. Single path one-shot neural architecture search with uniform sampling. In *European Conference on Computer Vision*, pages 544–560. Springer, 2020. 2, 3
 - [16] Ryuhei Hamaguchi, Aito Fujita, Keisuke Nemoto, Tomoyuki Imaizumi, and Shuhei Hikosaka. Effective use of dilated convolutions for segmenting small object instances in remote sensing imagery. In *2018 IEEE winter conference on applications of computer vision (WACV)*, pages 1442–1450. IEEE, 2018. 6
 - [17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 6, 8, 9
 - [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2, 5
 - [19] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in neural information processing systems*, pages 2017–2025, 2015. 2
 - [20] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. 2
 - [21] Duo Li, Anbang Yao, and Qifeng Chen. Psconv: Squeezing feature pyramid into one compact poly-scale convolutional layer. *arXiv preprint arXiv:2007.06191*, 2020. 2, 7
 - [22] Xiang Li, Chen Lin, Chuming Li, Ming Sun, Wei Wu, Junjie Yan, and Wanli Ouyang. Improving one-shot nas by suppressing the posterior fading. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13836–13845, 2020. 3
 - [23] Yanghao Li, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Scale-aware trident networks for object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 6054–6063, 2019. 2
 - [24] Zeming Li, Chao Peng, Gang Yu, Xiangyu Zhang, Yangdong Deng, and Jian Sun. Detnet: A backbone network for object detection. *arXiv preprint arXiv:1804.06215*, 2018. 2
 - [25] Feng Liang, Chen Lin, Ronghao Guo, Ming Sun, Wei Wu, Junjie Yan, and Wanli Ouyang. Computation reallocation for object detection. *arXiv preprint arXiv:1912.11234*, 2019. 1, 3
 - [26] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 5, 8, 9
 - [27] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 5, 6
 - [28] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 5, 6
 - [29] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018. 2, 3
 - [30] Jiaheng Liu, Shunfeng Zhou, Yichao Wu, Ken Chen, Wanli Ouyang, and Dong Xu. Block proposal neural architecture search. *IEEE Transactions on Image Processing*, 30:15–25, 2020. 3
 - [31] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. In *Advances in neural information processing systems*, pages 4898–4906, 2016. 1
 - [32] Alejandro Newell, Zhiao Huang, and Jia Deng. Associative embedding: End-to-end learning for joint detection and grouping. In *Advances in neural information processing systems*, pages 2277–2287, 2017. 7, 9
 - [33] Junran Peng, Ming Sun, Zhaoxiang Zhang, Tieniu Tan, and Junjie Yan. Pod: practical object detection with scale-sensitive network. In *Proceedings of the IEEE International*

- Conference on Computer Vision*, pages 9607–9616, 2019. 1, 2, 7
- [34] Junran Peng, Ming Sun, ZHAO-XIANG ZHANG, Tieniu Tan, and Junjie Yan. Efficient neural architecture transformation search in channel-level for object detection. In *Advances in Neural Information Processing Systems*, pages 14313–14322, 2019. 1, 3, 4, 7
 - [35] Hieu Pham, Melody Y Guan, Barret Zoph, Quoc V Le, and Jeff Dean. Efficient neural architecture search via parameter sharing. *arXiv preprint arXiv:1802.03268*, 2018. 3
 - [36] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 5, 6, 7, 8, 9
 - [37] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 5
 - [38] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 5
 - [39] Shuai Shao, Zijian Zhao, Boxun Li, Tete Xiao, Gang Yu, Xiangyu Zhang, and Jian Sun. Crowdhuman: A benchmark for detecting human in a crowd. *arXiv preprint arXiv:1805.00123*, 2018. 6
 - [40] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2
 - [41] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alex Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. *arXiv preprint arXiv:1602.07261*, 2016. 2
 - [42] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015. 2
 - [43] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 2
 - [44] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 9627–9636, 2019. 5, 6
 - [45] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 5
 - [46] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. 2
 - [47] Arber Zela, Thomas Elsken, Tonmoy Saikia, Yassine Marakchi, Thomas Brox, and Frank Hutter. Understanding and robustifying differentiable architecture search. *arXiv preprint arXiv:1909.09656*, 2019. 4
 - [48] Rui Zhang, Sheng Tang, Yongdong Zhang, Jintao Li, and Shuicheng Yan. Scale-adaptive convolutions for scene parsing. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2031–2039, 2017. 2
 - [49] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017. 2
 - [50] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016. 3
 - [51] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018. 3

A. Comparison with NAS-based methods

SPOS and DARTS are widely used as efficient and fast NAS algorithms, but they do not work well in our search space due to our huge dilation search space. As can be seen from the Table 11, our EDO algorithm has obvious advantages over SPOS and DARTS, with faster search speed and higher accuracy.

Table 11. Comparison between our EDO algorithm and the baseline method, SPOS and DARTS. The unit of column "Search Cost" in the table is "GPU-days". "20 + 530" means the cost for training the supernet is 20 GPU-days and the cost for selecting the subnet is 530 GPU-days.

Method	Search Cost	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
baseline	-	36.4	58.6	39.2	21.7	40.2	46.4
SPOS[15]	20 + 530	38.3	60.9	41.2	22.5	41.8	50.0
DARTS[29]	100	38.0	60.6	40.8	22.1	41.6	49.8
EDO(our)	12	38.9	61.6	41.8	22.9	42.3	50.5

A.1. Single path one-shot (SPOS)

We implement an efficient SPOS searching algorithm based on mask tensors.

Supernet Training For each convolution layer with the weight $W \in R^{C_{out} \times C_{in} \times (2kd_{max}+1) \times (2kd_{max}+1)}$, we insert a mask tensor $m \in R^{C_{out} \times (2kd_{max}+1) \times (2kd_{max}+1)}$, whose value is 0 or 1, to realize the sampling of different dilation patterns during training. As shown in Figure 7, for each channel in the mask, we fill it with one of d_{max}^2 dilation patterns at random with the value 0 or 1 at each training step. We use $W \times m$ in the forward process to ensure that only the parts of W corresponding to the current dilation are updated.

The training details and results We verify the superiority of our EDO over SPOS and DARTS by using FasterRCNN (R50) on the COCO dataset and the configurations are the same as Section 4.2.1 in the main paper. After training the supernet, we randomly sample 500 subnets for evaluating and then we pick the top10 subnets with the highest mAP values for retraining. As shown in Figure 8, the difference of mAP between the worst subnet and the best subnet is about 0.8, which is a very small value. The results indicate that the SPOS supernet is not distinguishable enough.

A.2. DARTS

Following the regime mentioned in [32], we introduce DARTS to search our inception convolution at the channel level. As shown in Table 11, it seems that DARTS degenerates into random sampling in our dilation search space because it only achieves a small improvement of about 0.2% mAP when compared with uniform sampling which is 37.8%.

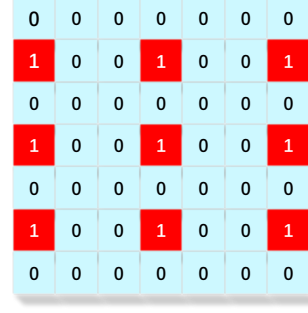


Figure 7. The schematic diagram of a mask tensor and the corresponding dilation pattern.

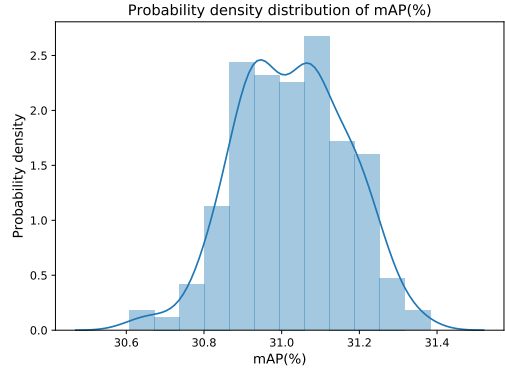


Figure 8. The probability density distribution of mAP(%).

B. Visualization of ERFs

Following the instructions in [30] of the main paper, we visualize the ERF of the neuron at the center of the last convolutional layer. More specifically, the input image is set to all 1, and only the gradient of the value right in the center of the output feature map across channels is calculated against the input. We draw the absolute value of the gradient. As in [32] of the main paper, to only focus on the strength of connections, ReLUs are removed during visualization. As shown in the Figure 9, the sizes of ERFs in our inception convolution network are much larger than ERFs of the vanilla structures. Furthermore, the shapes of our region are more complex than those of the standard networks, which indicates that these type of ERFs could better handle the object detection task and leads to high performance.

C. Training details of RFCN on Pascal Voc dataset

We use VOC trainval2007+trainval2012 as our training dataset and VOC test2007 as our validation dataset. The kernel size of the supernet is 9. We use SGD as an optimizer with 0.9 momentum and 0.0001 weight decay. The model is trained for 20 epochs and the learning rate is divided by 10 at the 16th and 18th epochs. Warm-up is adopted for both baselines and our inception models.

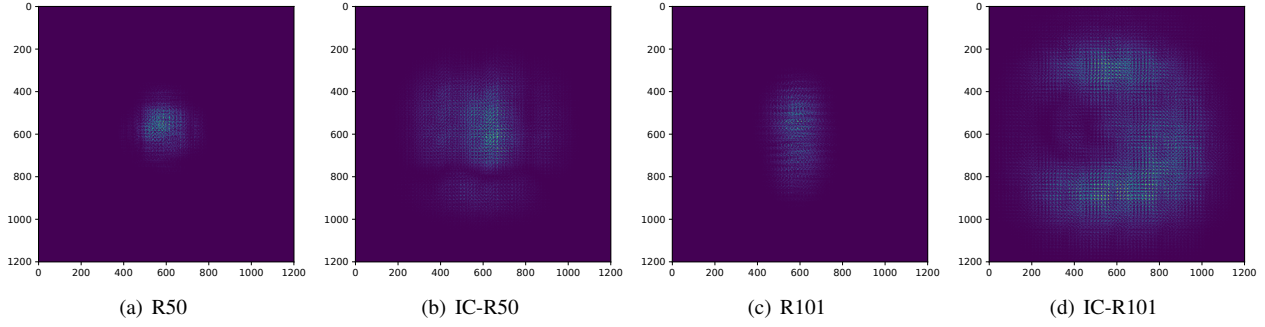


Figure 9. Visualization of ERFs in the inception architectures and the vanilla architectures with Faster-RCNN on the COCO dataset.

D. Searched dilation distribution

The searched dilation of EDO, SPOS and DARTs are provided, which are *frcnn-edo-r50.json*, *frcnn-spos-r50.json*, *frcnn-darts-r50.json*. The dilation distribution searched by EDO is relatively obvious, in which ordinary 3×3 normal convolution accounts for the majority, while other dilation patterns account for the minority. However, the distributions of dilation searched by SPOS and DARTS are very similar to a uniform distribution, indicating meaningless searching.