

Neural Architecture Adaptation for Object Detection by Searching Channel Dimensions and Mapping Pre-trained Parameters

Harim Jung	Myeong-Seok Oh	Cheoljong Yang	Seong-Whan Lee
Dept. Artificial Intelligence Korea University Seoul, Republic of Korea Email: hr_jung@korea.ac.kr	Dept. Computer Engineering Korea University Seoul, Republic of Korea Email: ms_oh@korea.ac.kr	Vision AI Lab, AI Center NCSOFT Seongnam, Republic of Korea Email: cjang@ncsoft.com	Dept. Artificial Intelligence Korea University Seoul, Republic of Korea Email: sw.lee@korea.ac.kr

Abstract—Most object detection frameworks use backbone architectures originally designed for image classification, conventionally with pre-trained parameters on ImageNet. However, image classification and object detection are essentially different tasks and there is no guarantee that the optimal backbone for classification is also optimal for object detection. Recent neural architecture search (NAS) research has demonstrated that automatically designing a backbone specifically for object detection helps improve the overall accuracy. In this paper, we introduce a neural architecture adaptation method that can optimize the given backbone for detection purposes, while still allowing the use of pre-trained parameters. We propose to adapt both the micro- and macro-architecture by searching for specific operations and the number of layers, in addition to the output channel dimensions of each block. It is important to find the optimal channel depth, as it greatly affects the feature representation capability and computation cost. We conduct experiments with our searched backbone for object detection and demonstrate that our backbone outperforms both manually designed and searched state-of-the-art backbones on the COCO dataset.

I. INTRODUCTION

In recent years, there has been a growing interest in neural architecture search (NAS), which automates the process of designing neural network architectures. NAS is useful since the manual design of neural architectures heavily relies on human experts' limited experience and prior knowledge, which may lead to tedious trial and error. The automatically searched architectures have shown highly competitive performance compared to manually designed architectures in various tasks [1]–[6], including image classification [7]–[10], language modeling [11], [12], object detection [13]–[16], and semantic segmentation [17], [18]. Object detection is a task of localizing and classifying various scales of objects in a given image and NAS for object detection can be mainly divided into three categories: backbone [13], [14], neck [19], [20], and head search [21], [22]. The search can be conducted through various

strategies and most previous works have used reinforcement learning [19]–[21], evolutionary algorithms [14], [23], or gradient-based learning [8], [13], [24].

Most existing object detectors directly utilize the backbone designed for image classification and its pre-trained parameters on ImageNet [25]. However, image classification and object detection are essentially different tasks and there is no guarantee that the optimal backbone for image classification is also optimal for object detection. Therefore, it is necessary to specifically search for a backbone architecture that can extract features appropriate for detecting and classifying multiple objects of different scales. One of the barriers of applying NAS to object detection is the pre-training cost, as backbone pre-training is still a necessary but costly procedure in object detection, in order to achieve faster convergence and higher accuracy. The search process can incur extra costs from pre-training the supernet as in DetNAS [14], which encompasses all candidate architectures. Therefore, instead of searching for an entire network architecture from scratch as in [8], [14], [17], we adopt an efficient function-preserving neural architecture adaptation method to search the optimal backbone for object detection.

For the first time, we propose to adapt an existing backbone network in both the macro- and micro-level, while still permitting the use of ImageNet pre-trained parameters. We adjust the micro-architecture by modifying the specific operations and the macro-architecture by modifying the output channel dimension of each block, as well as the number of operations or layers of the network. Channel dimension can also be interpreted as the number of filters of convolutional operators within the block, which decides the type and amount of information to extract from previous feature maps. It is important to find the appropriate channel depths especially within an object detection framework since feature maps of multiple levels are utilized for localizing and classifying various scales of objects. In our work, we propose to search the optimal channel dimensions of the source backbone network, while considering the trade-off between the computational cost and accuracy of the given detection framework. We summarize

This work was supported by NCSOFT, in part by Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. 2019-0-00079, Artificial Intelligence Graduate School Program (Korea University)).

our main contributions as follows:

- We propose to adapt the source backbone within an object detection framework by searching for the optimal operations and output channel dimensions, while enabling the mapping of ImageNet pre-trained parameters in the searching and fine-tuning stages.
- We propose an efficient method for channel search which uses a shared block for all channel dimension candidates with a non-overlapping channel masking technique, enabling fast search and creating a decoupling effect among different channel dimension candidates.
- Our method outperforms other state-of-the-art hand-crafted and searched backbones in object detection accuracy.

II. RELATED WORK

A. Neural Architecture Search Methods

Neural Architecture Search (NAS) aims at automatically designing neural network architectures. The main issue of NAS is to decide what to search for and how to search for those components. Early works have used reinforcement learning (RL) [19]–[21], [26], [27] and evolutionary algorithm (EA) [14], [23], [28], [29] as the search strategy. However, these methods often require huge computational costs in the searching stage, usually thousands of GPU days. Later on, ENAS [11] proposed a more efficient approach of RL by allowing child models to share parameters, eliminating the need for training from scratch every time. In recent years, many efficient one-shot approaches [12], [23] were proposed. By introducing the idea of supernet, which is a representation of all possible architectures in the search space, they dramatically reduced the overall search cost. In particular, differentiable NAS [8], [12], [13], [24] enabled fast search by training with stochastic gradient descent. Unlike previous approaches that used evolutionary or reinforcement learning over a discrete and non-differentiable search space, differentiable NAS implemented a continuous relaxation of the search space. We also use a differentiable search strategy to efficiently carry out the proposed architecture adaptation.

B. Neural Architecture Search for Object Detection

Object detection is one of the fundamental tasks in computer vision and it aims to both localize and classify each object instance given an image. Two-stage detectors, such as R-CNN [30], Fast-RCNN [31], and Faster R-CNN [32], usually have a separate region proposal network (RPN) that selects anchor boxes likely to include object instances before classifying them, whereas one-stage detectors, such as SSD [33], YOLO [34], and RetinaNet [35], perform localization and classification at once. In general, object detectors consist of three components: a backbone that extracts features from the input image, a neck attached to the backbone that fuses the features, and a head that localizes and classifies object instances from the extracted features.

Previous NAS literature usually selected one component to search for out of the three parts. DetNAS [14] was the first

work to search for a backbone specifically for object detection. DetNAS first pre-trained the supernet with ImageNet [25] and fine-tuned it on COCO [36]. Then, the architecture search on the trained supernet was conducted using EA. MobileDets [26] suggested using RL to find the expansion or compression factor of channels. Neck search algorithms mostly aim to find a novel architecture for the feature pyramid network (FPN). For instance, NAS-FPN [19] attempted to find a FPN architecture in a novel scalable search space covering all cross-scale connections. Head search algorithms aim to find optimal sub-networks for both classification and localization. NAS-FCOS [22], in addition to searching for the FPN, aimed to search the prediction head of an anchor-free object detector named FCOS. Moreover, Hit-Detector [37] proposed a hierarchical trinity search to simultaneously search all three components including the backbone, neck, and head. They argued that it is suboptimal to search for each part separately, as they essentially communicate with each other throughout the overall framework.

It is important to consider the fact that backbones for object detection require ImageNet pre-training for faster convergence and higher accuracy in most cases. The cost of pre-training is non-negligible since ImageNet is a huge dataset. In order to utilize the pre-trained parameters, architecture adaptation or transformation methods for object detection were proposed. NATS [24] and Liu *et al.* [38] proposed to adjust the dilation rates of convolution operators at the channel level. FNA++ [13] proposed to adjust kernel sizes and expansion ratios of convolutional blocks. While FNA++ [13] focuses on modifying the micro-architecture (building blocks, operations, etc.), we propose to also modify the macro-architecture (number of layers, channel dimensions, etc.).

III. METHOD

A. Neural Architecture Search

1) *Search Space*: We expand the source network to form a supernet, which represents the overall search space as shown in Fig. 1. Since the source network is a hand-crafted architecture originally designed for the classification task, the later layers for classification are disregarded. Moreover, the first two blocks are identical to the source network. These pre-defined stem layers are followed by K searchable blocks. In this work, we experiment with MobileNetV2 [39] as the source network within the RetinaNet framework [35], but our search method can be generally applied to other one-stage or two-stage detectors, as well as other backbone networks. The overall search space is described in detail in Table I.

Operation candidates. We define N operation candidates, which include variations of the inverted residual convolution modules (MBConv) of MobileNetV2 [39] and skip connection. MBConv operation includes a sequence of layers, as shown in Fig. 3. More specifically, the search space includes MBConv operations with kernel sizes of $\{3, 5, 7\}$ and expansion factors of $\{3, 6\}$, consistent for all blocks. Kernel sizes control the distribution of effective receptive fields, which is especially an important factor for detecting various scales of objects,

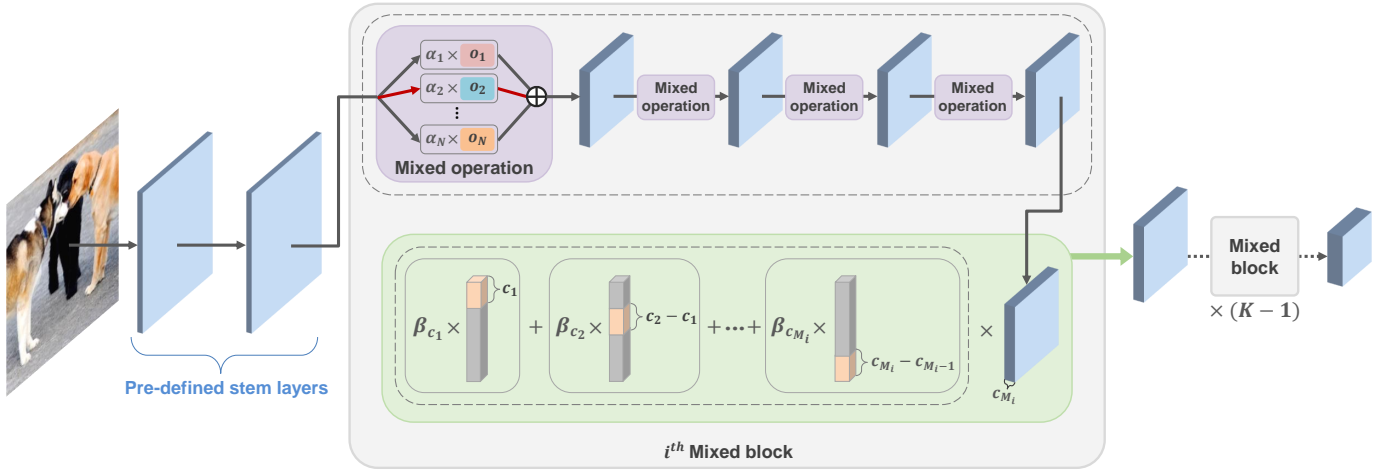


Fig. 1. Illustration of the overall supernet. The supernet expresses the whole search space with K searchable blocks, where each block has a total of M output channel dimension candidates and n searchable operations with N operation candidates. A mixed operation is defined as a weighted sum of all operation candidates and a mixed block is defined as an element-wise multiplication between a weighted sum of non-overlapping channel mask vectors and one shared feature map with the maximum channel dimension candidate c_{M_i} . After the supernet is sufficiently trained, the operation with the highest α and the channel dimension with the highest β is chosen to form the final searched architecture.

TABLE I
SEARCH SPACE BASED ON MOBILENETV2 [39]. "SB" DENOTES SEARCHABLE BLOCK, n DENOTES THE NUMBER OF MBConv OPERATIONS AND s INDICATES THE STRIDE OF THE FIRST OPERATION WITHIN EACH BLOCK. e STANDS FOR EXPANSION RATIO AND k STANDS FOR KERNEL SIZE. FOR OUTPUT CHANNEL DIMENSION c_{out} , THE TUPLES OF THREE ELEMENTS REPRESENT (MINIMUM, MAXIMUM, STEP).

Block type	n	s	k	e	c_{out}
2D Conv	1	2	3	-	32
MBConv ($k3_e1$)	1	1	3	1	16
TBS	4	2	{3, 5, 7}	{3, 6}	(16, 28, 2)
TBS	4	2	{3, 5, 7}	{3, 6}	(28, 48, 2)
TBS	4	2	{3, 5, 7}	{3, 6}	(48, 72, 2)
TBS	4	1	{3, 5, 7}	{3, 6}	(72, 128, 4)
TBS	4	2	{3, 5, 7}	{3, 6}	(128, 256, 4)
TBS	1	1	{3, 5, 7}	{3, 6}	(256, 400, 4)

while expansion factors control how much the channels are expanded and reduced within each operation. The blocks of MobileNetV2 [39] include at most 4 operations, so we allow up to 4 operations to be searched in each block for further flexibility. When skip connection is chosen, the corresponding operation is removed, adjusting the total number of layers. In order to avoid the case where all operations are skipped within each block, we exclude the skip connection candidate for the first MBConv operation.

Channel dimension candidates. We define M_i output channel dimension candidates for each block B_i and the specific set of candidates differ for each block. A block contains multiple MBConv operations with the same output channel dimension. The output channels may either increase, decrease or stay the same compared to the input dimension. The output channel dimensions of the first two blocks are fixed, following the source network. The channel dimensions

of the source network MobileNetV2 [39] for the 6 searchable blocks are [32, 16, 24, 32, 64, 96, 160, 320] and the specific channel candidates are described in Table I. Empirically, we found that keeping or increasing the channel dimensions for subsequent blocks throughout the network is most effective.

2) *Continuous Relaxation of Search Space:* In order to make the search space differentiable, it is necessary to relax the discrete supernet to be continuous. We achieve this using the concept of a mixed operation and a mixed block.

Mixed operation. Let \mathbb{O} be the set of operation candidates in the search space. In the supernet, we assign an architecture parameter α_o to each candidate operation $o \in \mathbb{O}$. Following DARTS [12], in order to make the discrete architecture search space into a continuous one, we relax the discrete choice of a specific operation to a softmax function over all operation candidates. By passing α through softmax, each operation is assigned a probability value. The mixed operation of the ℓ th operation is defined as a weighted sum of all operation candidates and can be expressed as:

$$\bar{o}_\ell(x_{\ell-1}) = \sum_{o \in \mathbb{O}} \frac{\exp(\alpha_o)}{\sum_{o' \in \mathbb{O}} \exp(\alpha_{o'})} o(x_{\ell-1}), \quad (1)$$

where $x_{\ell-1}$ denotes the input tensor, which is the output feature map of the previous operation $\bar{o}_{\ell-1}$. At the end of search, a discrete architecture can be obtained by replacing each mixed operation \bar{o}_ℓ with the most likely operation, that is, the operation with the highest α parameter.

Mixed block. In a similar fashion to how we defined mixed operation, we define a mixed block as a weighted sum of all candidate blocks with different output channel dimensions. In this work, a block is defined as a sequence of mixed operations with the same output channel dimensions. Let \mathbb{C} be the set of output channel dimension candidates of the i th mixed block.

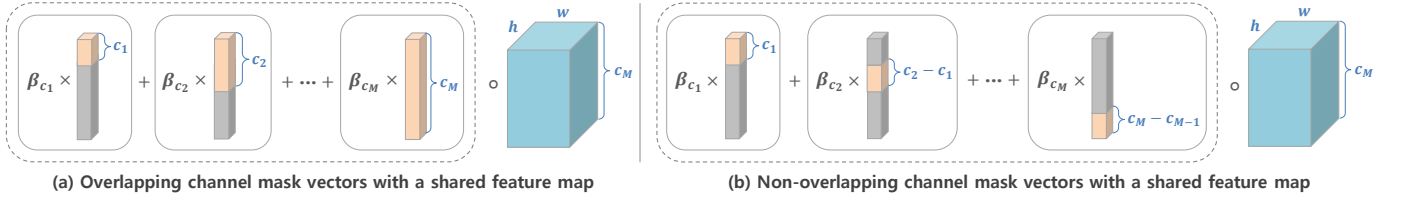


Fig. 2. Mixed block with overlapping channel masks and the proposed non-overlapping channel masks. There are M candidate blocks, each with different output channel dimensions up to c_M with beta parameters up to β_{c_M} . w and h denote the width and height of the shared feature map.

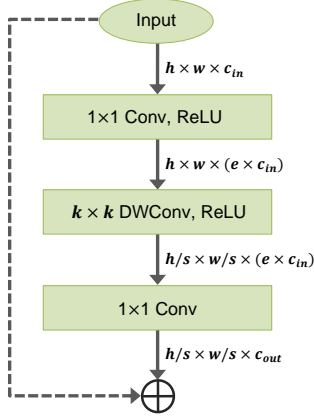


Fig. 3. MBConv operation. The search space includes operation candidates with various kernel sizes k , expansion factors e , and output channel dimensions c_{out} .

In the supernet, we assign an architecture parameter β_c to each candidate $c \in \mathbb{C}$. The i th mixed block can be expressed as:

$$\begin{aligned} & \bar{B}_i[\mathbb{C}, n, s](x_{i-1}) \\ &= B_i[c_{M_i}, n, s](x_{i-1}) \circ \sum_{c \in \mathbb{C}} \frac{\exp(\beta_c)}{\sum_{c' \in \mathbb{C}} \exp(\beta_{c'})} \mathbb{M}_c, \end{aligned} \quad (2)$$

where the inputs of B_i are n that denotes the number of mixed operations and s that refers to the stride of the first operation within the block. The strides of the following operations are set to 1. The specific set of candidates differ for each mixed block, having up to M_i candidates. \mathbb{M}_c denotes the zero mask column vector corresponding to each channel dimension candidate c .

A naive approach would be creating individual blocks for each output channel candidate and adding them through zero-padding, as the output feature maps originally cannot be added due to mismatching channel dimensions. However, this method incurs linearly increasing cost as more channel candidates are added. Therefore, we suggest a more efficient approach of using one shared block that has the largest channel dimension c_{M_i} and channel masking vectors, which only requires one forward pass and one feature map. One possible approach is to simply mask out the extra channels exceeding the current channel candidate, as shown in Fig. 2 (a). However, even when the weight of a higher channel dimension becomes very high, its value cannot easily overcome the accumulated weights of lower dimensions and as a result, the supernet was not properly trained. Therefore, we propose to use non-overlapping channel

mask vectors, as shown in Fig. 2 (b). This approach creates a decoupling effect among different channel candidates. When the weighted sum of channel masks is multiplied to the shared block through Hadamard product, the resulting mixed block embeds the various weights of different channel options. In other words, the channel dimensions with higher β weights would be clearly emphasized within the mixed block, while those with lower β weights would have less effect.

3) *Optimization*: We optimize the supernet by both training the operation and architecture parameters using two sets of training data of the same size, $trainA$ and $trainB$. We aim for a bi-level optimization with α and β as the upper-level variables and w as the lower-level variable:

$$\begin{aligned} & \min_{\alpha, \beta} \mathcal{L}_{trainB}(w^*(\alpha, \beta), \alpha, \beta) \\ & \text{s.t. } w^*(\alpha, \beta) = \underset{w}{\operatorname{argmin}} \mathcal{L}_{trainA}(w, \alpha, \beta). \end{aligned} \quad (3)$$

At first, we solely train the operation parameters w for some epochs until the model accuracy is not too low. After the operations are sufficiently trained, we start to alternatively train the operation parameters and architecture parameters, which approximates to a bi-level optimization.

Although searching for a model with the highest detection accuracy is important, it is also crucial to consider the computational cost, which decides the speed and efficiency of the detector framework. We use Multiply-Adds operations (MAdds) as the metric for measuring the computational cost. We define the loss function as a multi-objective optimization by adding a cost-based regularization term. This term adds some cost constraint so that an efficient model can be searched. The overall loss function is defined as:

$$\mathcal{L}(w, \alpha, \beta) = \mathcal{L}_{model}(w, \alpha, \beta) + \lambda C(\alpha, \beta), \quad (4)$$

where \mathcal{L}_{model} denotes the loss from the model including the classification loss and regression loss, while λ controls the magnitude of the cost regularization.

In order to estimate the model cost during the search stage, we initially create a lookup table recording the cost of each operation in the search space. In a similar fashion to how we define mixed operation and mixed block, we consider the probabilities embedded in the architecture parameters and

TABLE II
COMPARISON OF OBJECT DETECTION RESULTS OF RETINANET [35] ON COCO TEST-DEV [36].
NETWORK TYPE DENOTES WHETHER THE BACKBONE WAS HAND-CRAFTED OR SEARCHED THROUGH NAS.
FOR BASELINE MODELS, WE DIRECTLY CITE THE SEARCH COST BASED ON THE GPU USED IN THEIR ORIGINAL PAPERS.

Detector	Network Type	Backbone	Params	MAdds	mAP (%)	Search Cost (GPU-days)
RetinaNet	Hand-crafted	MobileNetV2 [39]	11.49M	133.05B	32.8	-
		ResNet-50 [40]	37.97M	202.84B	35.5	-
	NAS	DetNAS [14]	13.41M	133.26B	33.3	44.0*
		FNA [41]	11.73M	133.03B	33.9	6.0 [†]
		FNA++ [13]	11.90M	132.86B	34.7	5.3 [†]
		FBNet-C [42]	12.65M	134.17B	34.9	9.0 [‡]
		Proxyless (GPU) [43]	14.62M	135.81B	35.0	8.3*
		DenseNAS-C [8]	13.24M	133.91B	35.1	2.7*
		Ours	11.63M	134.92B	35.5	1.7[§]

*NVIDIA V100. [†]TITAN Xp. [‡]GPU not specified, [§]NVIDIA A100.

compute the MAdds of the whole architecture, expressed as:

$$C(\alpha, \beta) = \sum_{i=1}^K \bar{C}^i(\alpha, \beta), \quad (5)$$

$$\bar{C}^i(\alpha, \beta) = \sum_{c \in \mathcal{C}^i} \frac{\exp(\beta_c)}{\sum_{c' \in \mathcal{C}^i} \exp(\beta_{c'})} C_c^i(\alpha), \quad (6)$$

$$C_c^i(\alpha) = \sum_{\ell=1}^n \bar{C}_c^{i,\ell}(\alpha), \quad (7)$$

$$\bar{C}_c^{i,\ell}(\alpha) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_o)}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'})} C_{c,o}^{i,\ell}, \quad (8)$$

where the total cost C is calculated by summing the cost of K mixed blocks, which is a weighted sum of β weights and the cost of each candidate block C_c^i with channel dimension candidate c as in Eq. 6. In Eq. 7, the cost of the i th mixed block is computed by summing the cost of n mixed operations within the block, which is a weighted sum of α weights and the cost of each operation candidate o as in Eq. 8.

B. Mapping Pre-trained Parameters

After the supernet is trained, we derive the final architecture by selecting the operation with the highest α parameter and the output channel dimension of each block with the highest β parameter. The pre-trained parameters of the source network are mapped to the supernet and the searched architecture. Our method enables parameter mapping for different number of operations, channel dimensions, and kernel sizes.

1) *Number of operations*: Our method enables depth search by adjusting the number of operations in each block (stage), which ultimately decides the total number of layers of the searched network. Our search space includes the skip connection operator, which removes the corresponding operation in the searched network. The parameters of extra operations or layers are copied from the last layer in the original block.

2) *Channel dimension*: Our method enables width search by adjusting the number of output channels of each block and the expansion rate of MBCConv operations. When a higher

channel dimension is chosen, the parameters of corresponding dimensions are directly mapped and additional channels are assigned 0. When a lower channel dimension is chosen, the exceeding channels are disregarded.

3) *Kernel size*: Our method enables operation search by adjusting the kernel sizes of MBCConv operations. When the searched kernel size is larger than that of the source network, the original parameters are mapped in the central region and the surrounding region is assigned 0. When the searched kernel size is smaller, only the overlapping central region is mapped from the source network and the rest are disregarded. We add small random noises to the mapped parameters to allow backpropagation of those assigned 0.

IV. EXPERIMENTS

A. Implementation Details

Our model is implemented based on PyTorch [44]. For RetinaNet [35], the input image is resized to 800×1333 during inference and to 800×1088 for calculating Multiply-Adds operations (MAdds) following previous works. The search process runs for 14 epochs, where only operation parameters are trained for the first 8 epochs, and from epoch 9, architecture parameters are activated and trained alternatively with operation parameters. The search cost takes approximately 1.7 GPU-days on a single NVIDIA A100 GPU with a batch size of 8. Operation parameters are trained with the SGD optimizer using a learning rate of 0.02, momentum of 0.9, and weight decay of 0.0001. Architecture parameters are trained with the Adam optimizer using a learning rate of 0.0003 and weight decay of 0.001. Lastly, we fine-tune the searched network for 24 epochs, which takes approximately 3.6 GPU-days on NVIDIA A100 GPU with a batch size of 32. We use the SGD optimizer with a learning rate of 0.05, momentum of 0.9, and weight decay of 0.00005.

B. Quantitative Results

We report results using mean Average Precision (mAP), Average Precision (AP) with IoU threshold of 0.50 and 0.75,

TABLE III

COMPARISON BETWEEN USING INDIVIDUAL BLOCKS FOR EACH CHANNEL DIMENSION CANDIDATE AND USING ONE SHARED BLOCK FOR ALL CANDIDATES. THE TUPLES OF THREE ELEMENTS OF CHANNEL DIMENSION CANDIDATES REPRESENT THE LOWEST CHANNEL DIMENSION, HIGHEST CHANNEL DIMENSION, AND STEPS WITHIN THE RANGE. OBJECT DETECTION RESULTS OF RETINANET [35] ON COCO MINIVAL [36] ARE REPORTED.

Method	Channel dimension candidates	Max. number of channel candidates per block	Search cost	Params	MAdds	mAP (%)
Individual block	{24, 32, 64}, {32, 64, 96}, {64, 96, 160}, {96, 160, 320}, {160, 320, 480}, {320, 480, 640}	3	9.48 GPU-days	12.74M	133.84B	33.8
Shared block	(16, 28, 2), (28, 48, 2), (48, 72, 2), (72, 128, 4), (128, 256, 4), (256, 400, 4)	37	1.69 GPU-days	11.63M	134.92B	35.3

TABLE IV

COMPARISON BETWEEN OVERLAPPING AND NON-OVERLAPPING CHANNEL MASK TECHNIQUE FOR MIXED BLOCKS, AS WELL AS BETWEEN TRAINING FROM SCRATCH AND FINE-TUNING FROM PRE-TRAINED PARAMETERS. OBJECT DETECTION RESULTS OF RETINANET [35] ON COCO MINIVAL [36] ARE REPORTED.

Channel masking	Training	mAP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
Overlapping	Pre-trained	26.5	42.4	27.8	14.5	27.0	35.9
Non-overlapping	Scratch	33.0	52.8	35.5	18.0	36.5	43.4
	Pre-trained	35.3	54.2	37.3	19.4	38.4	46.1

AP of small, medium, and large object scales. In addition, we report the number of parameters and the number of MAdds to compare the size and computation cost of various backbones. Table II reports detection results of RetinaNet [35] using various backbones on the COCO test-dev set [36]. Our searched backbone was able to achieve the state-of-the-art accuracy on mAP, compared to other hand-crafted and searched models. Meanwhile, our model had the lowest number of parameters, enabling efficient training and inference. In addition, our searched model achieved comparable detection performance to ResNet-50 [40], despite having approximately 70% less parameters and 35% less MAdds. Compared to MobileNetV2 [39], our searched model improved the mAP by 2.5%, with only 1.3% more parameters and 1.5% more MAdds.

C. Ablation Studies

1) *Comparison between using individual blocks and one shared block for mixed blocks:* We compare two methods of defining the mixed block, between using a shared block with channel mask vectors and using individual candidate blocks with different channel dimensions through zero padding. In the shared block approach, all channel dimension candidates ultimately share the operation weights. The results reported in Table III show that the search cost of the shared block method is significantly lower than using individual candidate blocks. The low computational cost of the shared block method allowed us to add significantly more channel dimension candidates to search from, which we believe led to better detection accuracy. Considering the specific channel dimension candidates, we empirically found that gradually increasing the channel depth throughout the network yielded better results.

2) *Comparison between overlapping and non-overlapping channel masks for mixed blocks:* Table IV compares the results

of two channel masking methods, the method using overlapping mask vectors proposed in FBNetV2 [45], shown in Fig. 2 (a), and our proposed method using non-overlapping mask vectors, shown in Fig. 2 (b). Our proposed non-overlapping channel masking method achieved higher accuracy overall under the same setting. Empirically, the overlapping channel masking method did not train well, as the gradient values of higher channel dimensions for backpropagation progressively became almost identical. This may be due to the fact that the shallow channel dimensions of overlapping mask vectors are summed up to 1 or nearly 1. To avoid this effect, we suggested to use non-overlapping channel masking vectors to clearly emphasize the effect of optimal channel dimension candidates and to create a decoupling effect between each channel dimension candidate.

3) *Comparison between training with pre-trained parameters and from scratch:* Existing studies on object detection mostly utilize ImageNet pre-trained parameters of the given backbone. However, searching for a completely new architecture makes the use of existing pre-trained models impossible. This results in high computational cost, especially due to the fact that the large supernet must be pre-trained as well for better search results. Table IV demonstrates that fine-tuning based on pre-trained parameters yield better results than training from scratch. The results prove the importance of our architecture adaptation approach which allows for the mapping of pre-trained parameters.

V. CONCLUSION

In this paper, we proposed a function-preserving architecture adaptation method to search an optimal backbone specifically for object detection, beyond image classification. We achieved this by adjusting the number of layers, operations, and channel dimensions, in both the macro- and micro-level. In particular, searching for the optimal output channel dimensions of each block of the network increased its feature representation capability. One of the biggest advantages of our method is the low searching and training cost, achieved through the non-overlapping channel masking approach and the pre-trained parameter mapping scheme for the supernet and the searched backbone. We conducted experiments on the single-stage detection framework with our searched backbone and our model outperformed both manually designed and NAS-based state-of-the-art backbones on the COCO dataset.

REFERENCES

- [1] S.-W. Lee and S.-Y. Kim, "Integrated segmentation and recognition of handwritten numerals with cascade neural network," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 29, no. 2, pp. 285–290, 1999.
- [2] H.-D. Yang and S.-W. Lee, "Reconstruction of 3d human body pose from stereo image sequences based on top-down learning," *Pattern Recognition*, vol. 40, no. 11, pp. 3120–3131, 2007.
- [3] S.-W. Lee, J. H. Kim, and F. C. Groen, "Translation-, rotation- and scale-invariant recognition of hand-drawn symbols in schematic diagrams," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 4, no. 01, pp. 1–25, 1990.
- [4] S.-W. Lee and A. Verri, *Pattern Recognition with Support Vector Machines: First International Workshop, SVM 2002, Niagara Falls, Canada, August 10, 2002. Proceedings.* Springer, 2003, vol. 2388.
- [5] Y.-K. Lim, S.-H. Choi, and S.-W. Lee, "Text extraction in mpeg compressed video for content-based indexing," in *Proceedings 15th International Conference on Pattern Recognition. ICPR-2000*, vol. 4. IEEE, 2000, pp. 409–412.
- [6] D. Xi, I. T. Podolak, and S.-W. Lee, "Facial component extraction and face recognition with support vector machines," in *Proceedings of Fifth IEEE International Conference on Automatic Face Gesture Recognition*. IEEE, 2002, pp. 83–88.
- [7] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8697–8710.
- [8] J. Fang, Y. Sun, Q. Zhang, Y. Li, W. Liu, and X. Wang, "Densely connected search space for more flexible neural architecture search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 628–10 637.
- [9] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," in *Proceedings of International Conference on Learning Representations*, 2016.
- [10] A. Jordao, F. Akio, M. Lie, and W. R. Schwartz, "Stage-wise neural architecture search," in *Proceedings of International Conference on Pattern Recognition*, 2021, pp. 1985–1992.
- [11] H. Pham, M. Guan, B. Zoph, Q. Le, and J. Dean, "Efficient neural architecture search via parameters sharing," in *Proceedings of International Conference on Machine Learning*, 2018, pp. 4095–4104.
- [12] H. Liu, K. Simonyan, and Y. Yang, "DARTS: Differentiable architecture search," in *Proceedings of International Conference on Learning Representations*, 2018.
- [13] J. Fang, Y. Sun, Q. Zhang, K. Peng, Y. Li, W. Liu, and X. Wang, "FNA++: Fast network adaptation via parameter remapping and architecture search," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 9, pp. 2990–3004, 2021.
- [14] Y. Chen, T. Yang, X. Zhang, G. Meng, X. Xiao, and J. Sun, "DetNAS: Backbone search for object detection," in *Advances in Neural Information Processing Systems*, vol. 32, 2019, pp. 6642–6652.
- [15] H. Zhang, L. Wang, J. Sun, L. Sun, H. Kobashi, and N. Imamura, "NAS-EOD: an end-to-end neural architecture search method for efficient object detection," in *Proceedings of International Conference on Pattern Recognition*, 2021, pp. 1446–1451.
- [16] L. Yao, R. Pi, H. Xu, W. Zhang, Z. Li, and T. Zhang, "Joint-DetNAS: Upgrade your detector with nas, pruning and dynamic distillation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 175–10 184.
- [17] C. Liu, L.-C. Chen, F. Schroff, H. Adam, W. Hua, A. L. Yuille, and L. Fei-Fei, "Auto-DeepLab: Hierarchical neural architecture search for semantic image segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 82–92.
- [18] L.-C. Chen, M. D. Collins, Y. Zhu, G. Papandreou, B. Zoph, F. Schroff, H. Adam, and J. Shlens, "Searching for efficient multi-scale architectures for dense image prediction," in *Advances in Neural Information Processing Systems*, vol. 31, 2018, pp. 8713–8724.
- [19] G. Ghiasi, T.-Y. Lin, and Q. V. Le, "NAS-FPN: Learning scalable feature pyramid architecture for object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7036–7045.
- [20] H. Xu, L. Yao, W. Zhang, X. Liang, and Z. Li, "Auto-FPN: Automatic network architecture adaptation for object detection beyond classification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6649–6658.
- [21] B. Chen, G. Ghiasi, H. Liu, T.-Y. Lin, D. Kalenichenko, H. Adam, and Q. V. Le, "MnasFPN: Learning latency-aware pyramid architecture for object detection on mobile devices," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 607–13 616.
- [22] N. Wang, Y. Gao, H. Chen, P. Wang, Z. Tian, C. Shen, and Y. Zhang, "NAS-FCOS: Fast neural architecture search for object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 943–11 951.
- [23] R. Guo, C. Lin, C. Li, K. Tian, M. Sun, L. Sheng, and J. Yan, "Powering one-shot topological nas with stabilized share-parameter proxy," in *Proceedings of European Conference on Computer Vision*, 2020, pp. 625–641.
- [24] J. Peng, M. Sun, Z. Zhang, T. Tan, and J. Yan, "Efficient neural architecture transformation search in channel-level for object detection," in *Advances in Neural Information Processing Systems*, vol. 32, 2019, pp. 14 313–14 322.
- [25] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [26] Y. Xiong, H. Liu, S. Gupta, B. Akin, G. Bender, Y. Wang, P.-J. Kindermans, M. Tan, V. Singh, and B. Chen, "MobileDETS: Searching for object detection architectures for mobile accelerators," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3825–3834.
- [27] Z. Zhong, J. Yan, W. Wu, J. Shao, and C.-L. Liu, "Practical block-wise neural network architecture generation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2423–2432.
- [28] E. Real, S. Moore, A. Selle, S. Saxena, Y. L. Suematsu, J. Tan, Q. V. Le, and A. Kurakin, "Large-scale evolution of image classifiers," in *Proceedings of International Conference on Machine Learning*, 2017, pp. 2902–2911.
- [29] T. Elsken, J. H. Metzen, and F. Hutter, "Efficient multi-objective neural architecture search via lamarckian evolution," in *Proceedings of International Conference on Learning Representations*, 2019.
- [30] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 580–587.
- [31] R. Girshick, "Fast R-CNN," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448.
- [32] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems*, vol. 28, 2015, pp. 91–99.
- [33] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *Proceedings of European Conference on Computer Vision*, 2016, pp. 21–37.
- [34] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 779–788.
- [35] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2980–2988.
- [36] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *Proceedings of European Conference on Computer Vision*, 2014, pp. 740–755.
- [37] J. Guo, K. Han, Y. Wang, C. Zhang, Z. Yang, H. Wu, X. Chen, and C. Xu, "Hit-detector: Hierarchical trinity architecture search for object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 405–11 414.
- [38] J. Liu, C. Li, F. Liang, C. Lin, M. Sun, J. Yan, W. Ouyang, and D. Xu, "Inception convolution with efficient dilation search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 486–11 495.
- [39] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.

- [40] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [41] J. Fang, Y. Sun, K. Peng, Q. Zhang, Y. Li, W. Liu, and X. Wang, "Fast neural network adaptation via parameter remapping and architecture search," in *Proceedings of International Conference on Learning Representations*, 2020.
- [42] B. Wu, X. Dai, P. Zhang, Y. Wang, F. Sun, Y. Wu, Y. Tian, P. Vajda, Y. Jia, and K. Keutzer, "FBNet: Hardware-aware efficient convnet design via differentiable neural architecture search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 734–10 742.
- [43] H. Cai, L. Zhu, and S. Han, "ProxylessNAS: Direct neural architecture search on target task and hardware," in *Proceedings of International Conference on Learning Representations*, 2019.
- [44] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "PyTorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems*, vol. 32, 2019, pp. 8024–8035.
- [45] A. Wan, X. Dai, P. Zhang, Z. He, Y. Tian, S. Xie, B. Wu, M. Yu, T. Xu, K. Chen *et al.*, "FBNetV2: Differentiable neural architecture search for spatial and channel dimensions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 965–12 974.