
Block-wise Separable Convolutions: An Alternative Way to Factorize Standard Convolutions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Convolutional neural networks (CNNs) have demonstrated great capability of solving
2 various computer vision tasks with nice prediction performance. Nevertheless,
3 the higher accuracy often comes with an increasing number of model parameters
4 and large computational cost. This raises challenges in deploying them in resource-
5 limited devices. In this paper, we introduce block-wise separable convolutions
6 (BlkSConv) to replace the standard convolutions in order to compress deep CNN
7 models. First, BlkSConv expresses the standard convolutional kernel as an ordered
8 set of block vectors each of which is a linear combination of fixed basis block
9 vectors. Then it eliminates most basis block vectors and their corresponding coef-
10 ficients to obtain an approximated convolutional kernel. Moreover, the proposed
11 BlkSConv operation can be efficiently realized via a combination of pointwise and
12 group-wise convolutions. Thus the constructed networks have smaller model size
13 and fewer multiply-adds operations while keeping comparable prediction accu-
14 racy. However, it is unknown how to search a qualified hyperparameter setting
15 of the block depth and number of basis block vectors. To address this problem,
16 we develop a hyperparameter search framework based on principal component
17 analysis (PCA) to help determine these two hyperparameters such that the cor-
18 responding network achieves nice prediction performance while simultaneously
19 satisfying the constraints of model size and model efficiency. Experimental results
20 demonstrate the prediction performance of constructed BlkSConv-based CNNs
21 where several convolutional layers are replaced by BlkSConv layers suggested
22 by the proposed PCA-based hyperparameter search algorithm. Our results show
23 that BlkSConv-based CNNs achieve competitive performance compared with the
24 standard convolutional models for the datasets including ImageNet, CIFAR-10/100,
25 Stanford Dogs, and Oxford Flowers.

26 1 Introduction

27 In the past decade, Deep Learning (DL) has been the basis of many successes in artificial intelligence,
28 including a variety of applications in computer vision, reinforcement learning, and natural language
29 processing. One of the most popular deep neural networks is Convolutional Neural Network (CNN).
30 With the help of various techniques such as residual connections and batch normalization, it is easy to
31 train deep CNNs with many layers on powerful GPUs. While large-scale CNN models have achieved
32 great successes, they require huge computational complexity and massive storage. For example,
33 VGG16 (27) has 138 million parameters and requires 154700 million multiply-add operations
34 (MAdds) to classify an image. It is a great challenge to deploy them in real-time applications,
35 especially on devices with limited resources such as mobile phones and embedded systems. Thus,
36 the prediction models are required be compact and fast while keeping acceptable accuracy. The main
37 approach to be compact is the model compression which aims at establishing a tradeoff between

38 model efficiency and accuracy. In the area of model compression, methods to construct efficient and
39 compact CNNs are mainly divided into two approaches: one approach is to compress trained CNNs
40 and the other approach is to design new compact CNNs and train them from scratch. Many works
41 based on the first approach suggested several techniques such as quantization (33), model pruning
42 (6; 24), Huffman coding (6), and low rank factorization (12).

43 Studies in the second approach mainly explored many ways for factorizing convolutions. For
44 instance, Szegedy et al. (30) improved GoogLeNet (29) through factorizing convolutions with larger
45 spatial filters by a two-layer convolutional architecture with smaller spatial filters. At present, most
46 factorizing methods are usually performed via a combination of depthwise convolution, pointwise
47 convolution, and groupwise convolution. For example, in (25), the depth-wise separable convolutions
48 (DSCs) were proposed where the standard convolution is decomposed into a depth-wise convolution
49 and a pointwise convolution. The ShuffleNets (36; 18) utilizes pointwise group convolution with
50 channel shuffle to decompose the standard convolution. Moreover, many lightweight models based
51 on DSCs or groupwise convolutions such as MobileNets (8; 23; 7) and ShuffleNets (36; 18) were
52 proposed to greatly reduce computation cost while maintaining accuracy.

53 In this paper, we follow the research path of the second approach and propose block-wise separable
54 convolutions (*BlkSConv*) to replace standard convolutions. *BlkSConv* approximates a standard
55 convolution as follows. A standard $k \times k \times M$ convolutional kernel can be represented as an ordered
56 set of block vectors of size $k \times k \times t$. Since each block vector can be written as a linear combination
57 of k^2t basis vectors of size $k \times k \times t$, this standard convolutional kernel can be viewed as an ordered
58 set of block vectors each of which is a linear combination of k^2t basis block vectors. Then *BlkSConv*
59 eliminates most basis block vectors and their corresponding coefficients to obtain an approximated
60 convolutional kernel. As shown on the left of Figure 1, the extreme version of *BlkSConv* is called the
61 basic *BlkSConv* where only one basis block vector is used. When carefully setting the depth of the
62 block vector, that is the parameter t , an approximated convolution of fewer parameters can be obtained
63 and the corresponding compact CNN has acceptable prediction performance compared to the standard
64 convolutions. To increase the prediction accuracy of the basic *BlkSConv*, an enhanced version is
65 proposed by increasing the number of basis block vectors, that is the parameter s , as shown on the
66 right of Figure 1. However, adding too many basis block vectors will significantly increase the model
67 size and computational cost. Thus there is a tradeoff between model efficiency/size and accuracy. To
68 realize the full potential of the enhanced *BlkSConv* in trading-off model efficiency/size and accuracy,
69 we propose a framework based on the principal component analysis to search for the hyperparameters
70 t and s of each *BlkSConv* layer for the given standard convolutional network. The proposed search
71 framework suggests a possible setting of parameters t and s such that the constructed model based on
72 these selected hyperparameters may achieve high prediction accuracy while simultaneously satisfying
73 the constraints of model size and model efficiency in terms of MAdds.

74 To summarize, our main contributions are as follows. First, we develop a new convolutional layer
75 called *BlkSConv* to approximate the standard convolutional layer. To approximate a standard convo-
76 lutional kernel, *BlkSConv* divides the kernel into blocks and approximates each block by a linear
77 combination of several fixed basis block vectors. The constructed networks have small model size
78 and cost fewer multiply-adds operations while maintaining acceptable prediction accuracy. Then, we
79 also develop a search framework to determine the block depth and the number of basis block vectors
80 such that the corresponding networks with selected hyperparameters achieve comparable prediction
81 performance while simultaneously satisfying the constraints of model size and model efficiency.
82 We also present experimental results to demonstrate the performance of selected *BlkSConv*-based
83 CNNs based on our proposed hyperparameter search algorithm. Our results show that selected
84 *BlkSConv*-based CNNs achieve competitive performance compared with the standard convolutional
85 models for the datasets including ImageNet, CIFAR-10/100, Stanford Dogs, and Oxford Flowers.

86 2 Related Work

87 Many efforts have been devoted to improve the efficiency of CNNs which could be roughly divided
88 into three categories. First, model pruning is a popular method to improve efficiency of CNNs. In
89 (6; 37), their methods remove redundancy in the trained CNN model by pruning connection. In
90 (6; 21; 20; 35), the calculation amount of the trained model is compressed via quantization. In
91 (17; 11; 16; 9; 28), model filters that have small contributions are removed and the corresponding
92 trained model is fine-tuned to preserve the performance.

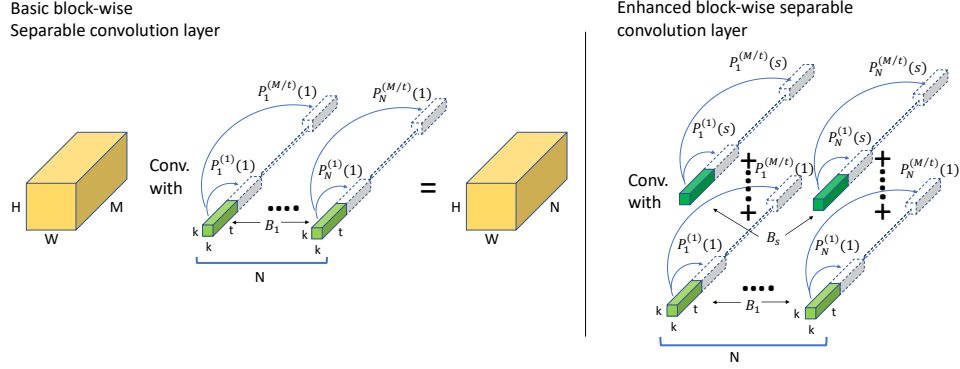


Figure 1: The proposed block-wise separable convolution and its enhanced version.

93 Second, many techniques are developed to factorize the standard convolutions. In (30), convolutions
 94 with larger spatial filters are factorized into two-layer convolutional architectures with smaller
 95 spatial filters. Through different combinations of depthwise convolution, pointwise convolution, and
 96 groupwise convolution, many well-known factorizing frameworks were developed. In (25), the depth-
 97 wise separable convolutions (DSCs) were proposed where the standard convolution is decomposed
 98 into a depth-wise convolution and a pointwise convolution. The ShuffleNets (36; 18) uses pointwise
 99 group convolution with channel shuffle to decompose the standard convolution. Moreover, many
 100 lightweight models based on DSCs or groupwise convolutions such as MobileNets (8; 23; 7) and
 101 ShuffleNets (36; 18) were proposed to greatly reduce computation cost while maintaining accuracy.

102 Recently, neural architecture search-based methods (34; 32; 38; 39; 31) have been proposed to
 103 automatically construct network architectures. These methods search over a set of network hyperpa-
 104 rameters including different types of convolutional layers and kernel sizes, to find a network structure
 105 which satisfies optimization constraints such as inference speed. Major search frameworks include
 106 genetic-based methods (34) and reinforcement learning based methods (38). These techniques were
 107 used in state-of-the-art CNN architectures such as MnasNet (31) and MobileNetV3 (7).

108 Convolution weights of trained CNNs are also analyzed in (1; 3; 26; 4). Following their analysis, sev-
 109 eral approaches toward reducing redundant weights were proposed. In (2; 12; 13), the convolutional
 110 kernels are approximated via low-rank factorization. In (4), the kernels are analyzed via principal
 111 component analysis.

112 3 Block-wise Separable Convolutions (BlkSConv)

113 For any natural number n , let $[n]$ denote the set $\{1, 2, \dots, n\}$. In a standard CNN, each convolutional
 114 layer converts an input tensor I of size $M \times X \times Y$ into an output tensor O of size $N \times X \times Y$
 115 by applying the filter kernels F_1, F_2, \dots, F_N , each of size $M \times \ell \times \ell$ with odd ℓ such that, for any
 116 $x, y, j \in [X] \times [Y] \times [N]$,

$$O(x, y, j) = \sum_{s_1=-(\ell-1)/2}^{(\ell-1)/2} \sum_{s_2=-(\ell-1)/2}^{(\ell-1)/2} \sum_{s_3=1}^M I(x + s_1, y + s_2, s_3) \cdot F_j(s_1, s_2, s_3). \quad (1)$$

117 During training, the weights of each kernel F_j are optimized via backpropagation. The total number
 118 of weight parameters to be optimized in each kernel F_j is $\ell^2 \cdot M$. In the subsequent work, we propose
 119 a framework to reduce the number of parameters of the standard convolutions while preserving its
 120 prediction performance. Then, in order to implement our new framework, we adopt a combination of
 121 pointwise and group-wise convolutions to efficiently realize the reduced convolutions. Combining
 122 these ideas, we introduce block-wise separable convolutions, denoted by *BlkSConv*. However,
 123 to generate a BlkSConv-based models, many hyperparameters should be determined for keeping
 124 prediction performance, model size, and model efficiency. Thus, we also propose an efficient
 125 hyperparameter search algorithm to select hyperparameters satisfying the given model constraints.

126 **3.1 Expressing a standard convolution via a linear combination of block vectors**

127 In this section, we propose block-wise separable convolutions. First, each convolutional kernel F_j
 128 of size $M \times \ell \times \ell$ can be expressed as a concatenation of M/t blocks $Q_j^{(1)}, Q_j^{(2)}, \dots, Q_j^{(M/t)}$ each
 129 of size $\ell \times \ell \times t$ where $Q_j^{(k)}(x, y, z) = F_j(x, y, z + (k - 1)t)$ for any $x, y, z \in [X] \times [Y] \times [t]$.
 130 We call t the block depth. Let $\{B_1, B_2, \dots, B_{t\ell^2}\}$ be a set of basis block vectors. Each $Q_j^{(k)}$ can
 131 be expressed uniquely as a linear combination of $B_1, B_2, \dots, B_{t\ell^2}$, that is, there exist $t\ell^2$ values
 132 $P_j^{(k)}(i) \in \mathbb{R}$ such that $Q_j^{(k)} = \sum_{i=1}^{t\ell^2} P_j^{(k)}(i) \cdot B_i$. In practice, $t\ell^2$ may be large. In order to reduce
 133 the model size, we require the number of basis block vectors is fewer than or equal to a fixed number
 134 s with $s < t\ell^2$. Now each $Q_j^{(k)}$ is replaced by the following linear combination of B_1, \dots, B_s , that
 135 is $\widehat{Q}_j^{(k)} = \sum_{i=1}^s P_j^{(k)}(i) \cdot B_i$. The corresponding convolutional kernel \widehat{F}_j is the concatenation of
 136 M/t blocks $\widehat{Q}_j^{(1)}, \dots, \widehat{Q}_j^{(M/t)}$. Therefore, the corresponding output tensor is

$$\widehat{O}(x, y, j) = \sum_{s_1, s_2 = -(\ell-1)/2}^{(\ell-1)/2} \sum_{s_3=1}^M I(x + s_1, y + s_2, s_3) \cdot \widehat{F}_j(s_1, s_2, s_3). \quad (2)$$

137 By Equation 2, the number of weight parameters in BlkSConv is $s \cdot (t \cdot \ell^2 + \frac{M}{t})$. To significantly
 138 reduce model size, we set $s = 1$. Figure 1 left illustrates the operation of BlkSConv when $s = 1$.
 139 In order to achieve the minimal model size, t can be set as \sqrt{M}/ℓ and the number of parameters
 140 becomes $2\ell\sqrt{M}$ while the parameter number of the standard and 1×1 pointwise convolutions are
 141 $M\ell^2$ and M , respectively. Thus, the constructed BlkSConv-based CNNs have smaller model size
 142 than existing lightweight CNN models. Take the ResNet34 (10) as an example where, in the last
 143 stage of the ResNet-34, the convolutional kernel size is 3×3 and the channel size is 512, that is $\ell = 3$
 144 and $M = 512$. In this case, the ratio between the parameter size of the BlkSConv-based convolutions
 145 and the standard convolutions is approximately 0.0295.

146 However, the prediction performance of the BlkSConv-based CNN with the smallest model size is
 147 usually worse than the standard CNNs. To increase accuracy, the number of basis block vectors
 148 should be increased, that is $s > 1$. Figure 1 right illustrates the operation of BlkSConv when $s > 1$.
 149 In this case, the number of parameters becomes $2s\ell\sqrt{M}$. Let us take convolutions in the last stage
 150 of ResNet-34 as examples. Let us set $t = 4$ in the BlkSConv. Now the ratio between the parameter
 151 size of the BlkSConv-based convolutions and the standard convolutions is approximately 0.0356.
 152 Thus we can add at least 5 basis block vectors to increase prediction accuracy. In this case, the ratio
 153 between the parameter size of the BlkSConv-based convolutions with 5 basis block vectors and the
 154 standard convolutions is approximately 0.178. In the experimental section, we demonstrate that the
 155 BlkSConv-based convolutions with few basis block vectors have prediction performance as well
 156 as the standard convolutions on ImageNet or even outperform the standard convolutions on several
 157 datasets when the backbone CNNs are ResNets.

158 The next problem is the computational efficiency of BlkSConv. If we compute the kernel \widehat{F}_j first and
 159 perform a regular convolution according to the kernel \widehat{F}_j , then it is obvious that the computational
 160 cost is larger than the cost for just performing a standard convolution. We will address this problem
 161 in the subsequent section.

162 **3.2 Implementation of BlkSConv via a combination of pointwise and group-wise convolutions**

163 In this section, we propose an efficient implementation method to realize BlkSConv. The flowchart
 164 of the proposed implementation is illustrated in Figure 2. To derive an efficient implementation for
 165 BlkSConv operation, we rewrite Equation 2 as follows.

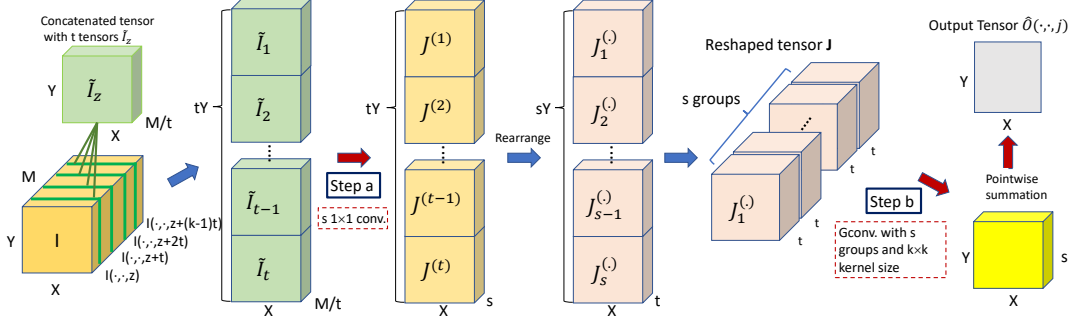


Figure 2: Flowchart of the block-wise separable convolution where Gconv. means the group-wise convolution.

$$\widehat{O}(x, y, j) = \sum_{s_1, s_2} \sum_{z=1}^t \sum_{k=1}^{M/t} I(x + s_1, y + s_2, z + (k-1)t) \cdot \widehat{Q}_j^{(k)}(s_1, s_2, z) \quad (3)$$

$$= \sum_{s_1, s_2} \sum_{z=1}^t \sum_{k=1}^{M/t} I(x + s_1, y + s_2, z + (k-1)t) \cdot \sum_{i=1}^s P_j^{(k)}(i) \cdot B_i(s_1, s_2, z) \quad (4)$$

$$= \sum_{i=1}^s \sum_{s_1, s_2} \sum_{z=1}^t B_i(s_1, s_2, z) \underbrace{\sum_{k=1}^{M/t} P_j^{(k)}(i) \cdot \overbrace{I(x + s_1, y + s_2, z + (k-1)t)}^{\widetilde{I}_z(x + s_1, y + s_2, k)}}_{J^{(z)}(x, y, i): \text{a point-wise convolution of } \widetilde{I}_z} \quad (5)$$

166 Let $\widetilde{I}_z(x, y, k)$ be a tensor of size $X \times Y \times M/t$ defined by $\widetilde{I}_z(x, y, k) \triangleq I(x, y, z + (k-1)t)$. We
 167 define $J^{(z)}(x, y, i) \triangleq \sum_{k=1}^{M/t} P_j^{(k)}(i) \cdot \widetilde{I}_z(x + s_1, y + s_2, k)$ which is a point-wise convolution of \widetilde{I}_z .
 168 Next, we define $J_i(x, y, z) \triangleq J^{(z)}(x, y, i)$ and let J be the reshaped tensor which is the concatenation
 169 of J_1, \dots, J_s , that is $J(x, y, z + (i-1)t) = J_i(x, y, z)$. Now Equation 5 can be rewritten as

$$\widehat{O}(x, y, j) = \sum_{i=1}^s \sum_{s_1, s_2} \sum_{z=1}^t B_i(s_1, s_2, z) \cdot J(x + s_1, y + s_2, z + (i-1)t). \quad (6)$$

170 Finally, Equation 6 is just a group-wise convolution of the tensor J with s groups.

171 Let us compute the computational cost (MAdds) of the implementation for BlkSConv. By Equation 5
 172 (Step a in Figure 2), the computational cost of s point-wise convolutions on the concatenation of
 173 $\widetilde{I}_1, \dots, \widetilde{I}_t$ is $sXYM$. In addition, by Equation 6 (Step b in Figure 2), the computational cost of the
 174 group-wise convolution on the tensor J is $sXYt\ell^2$. Finally, the computational cost of the point-wise
 175 summation in the last step is sXY . The total MAdds of a BlkSConv operation is $sXY(M + t\ell^2 + 1)$
 176 while the MAdds of a standard convolution is $XYM\ell^2$. Again, let us take convolutions in the last
 177 stage of ResNet-34 as examples. We set $s = 5$ and $t = 4$ as the hyperparameters of the BlkSConv-
 178 based convolution. Now the ratio between the MAdds of a BlkSConv-based convolution and a
 179 standard convolution is approximately 0.595. Thus the proposed BlkSConv operation is much more
 180 efficient than the standard convolution in practical cases. We remark that the proposed implementation
 181 requires much GPU memory due to using many group-wise and pointwise convolutions.

182 3.3 Hyperparameter search via principal component analysis

183 Designing a BlkSConv-based CNN involves hyperparameters including the block depth and the num-
 184 ber of basis block vectors in each convolutional layer that affect the performance of the corresponding
 185 CNN model. To realize an efficient BlkSConv-based CNN, we conduct a hyperparameter search algo-
 186 rithm based on principal component analysis of trained CNNs. Given a trained CNN, the algorithm
 187 generates the block depth and the number of basis block vectors for each standard convolutional layer

188 of the trained CNN in the following way. First, for each individual $\ell \times \ell \times M$ kernel K of the trained
189 CNN where we assume that $M = 2^\alpha$ for some $\alpha \in \mathbb{N}$, the kernel K is partitioned into M/t block
190 vectors $B_1, B_2, \dots, B_{M/t}$ each of size $\ell \times \ell \times t$ with $t \in \{1, 2, \dots, 2^\beta\}$ for some integer $\beta < \alpha$.
191 Next, we perform principal component analysis (PCA) on the set $\{B_1, B_2, \dots, B_{M/t}\}$. Then, for
192 a fixed integer γ and, for each $q \in \{1, 2, \dots, \gamma\}$, the algorithm computes the variance $V_{t,q}$ of the
193 kernel K which is explained by the first q principal components PC1, PC2, ..., PC q . In addition, let
194 $\text{CC}_{t,q}$ and $\text{MS}_{t,q}$ denote the MAdds and the model size of the BlkSConv under the setting that the
195 block depth is t and the number of basis block vectors is q , respectively. Note that the MAdds and the
196 model size of the standard convolution is exactly $\text{CC}_{M,1}$ and $\text{MS}_{M,1}$, respectively. After computing
197 all $V_{t,q}$, $\text{CC}_{t,q}$, and $\text{MS}_{t,q}$, the algorithm generates the feasible set

$$H_{\alpha_v, \alpha_c, \alpha_s} = \{(t, q) : V_{t,q} \geq \alpha_v, \text{CC}_{t,q} \leq \alpha_c \cdot \text{CC}_{M,1}, \text{and } \text{MS}_{t,q} \leq \alpha_s \cdot \text{MS}_{M,1}\} \quad (7)$$

198 for fixed positive constants $\alpha_v, \alpha_c, \alpha_s \in (0, 1)$. Finally, the algorithm chooses the hyperparameter
199 (t, q) from $H_{\alpha_v, \alpha_c, \alpha_s}$ according to the computational cost or the model size.

200 On one hand, note that the goal of BlkSConv is to maintain the prediction performance of the trained
201 standard CNN. In general, the prediction accuracy is proportional to the model size of the constructed
202 CNN. Therefore, in this sense, we choose the hyperparameters (\hat{t}, \hat{q}) from $H_{\alpha_v, \alpha_c, \alpha_s}$ such that the
203 constructed BlkSConv has the largest parameter size, that is

$$(\hat{t}, \hat{q}) = \arg \max_{(t,q) \in H_{\alpha_v, \alpha_c, \alpha_s}} \text{MS}_{t,q}. \quad (8)$$

204 One can expect that the generated BlkSConv-based CNN has nice prediction performance compared
205 to the original CNN with standard convolutions.

206 On the other hand, one of the advantage of BlkSConv operations is that BlkSConv can greatly reduce
207 the model size of the original standard CNN. Thus, in this sense, we can select the hyperparameters
208 (\tilde{t}, \tilde{q}) from $H_{\alpha_v, \alpha_c, \alpha_s}$ such that the constructed BlkSConv has the smallest parameter size, that is

$$(\tilde{t}, \tilde{q}) = \arg \min_{(t,q) \in H_{\alpha_v, \alpha_c, \alpha_s}} \text{MS}_{t,q}. \quad (9)$$

209 However, the prediction performance may degrade when the parameter size of the BlkSConv-based
210 model decreases. We will demonstrate in the experimental section that the BlkSConv-based CNNs
211 generated according to Equation 8 also have acceptable prediction accuracy compared to the standard
212 CNNs.

213 In summary, both Equation 8 and Equation 9 provide ways to determine hyperparameters from the
214 feasible set $H_{\alpha_v, \alpha_c, \alpha_s}$ such that corresponding BlkSConv-based CNNs have smaller model size and
215 fewer multiply-adds operations than the original CNN with standard convolutions.

216 Finally, let us consider the extreme case that two constants β and γ are set by $\beta = 0$ and $\gamma = 1$. Let
217 us further set the search parameter $\alpha_v = 0$. Under this restricted search condition, the cardinality of
218 the feasible set $H_{0, \alpha_c, \alpha_s}$ is always 1. Thus the outputs of Equation 8 and Equation 9 are the same.
219 In fact, the resulting BlkSConv-based CNN is exactly the same as the CNN where the standard
220 convolutions are replaced by the blueprint separable convolutions previously developed in (4).

221 4 Experiments

222 We evaluate BlkSConv and the proposed hyperparameter architecture search algorithm combining
223 with ResNet-10, ResNet-18, and ResNet-26 (5) on ImageNet (22), Stanford Dogs, (14), and Oxford
224 102 Flowers (19). The proposed methods are also evaluated combining with ResNet-20 and ResNet-56
225 on CIFAR 10/100 (15).

226 4.1 Hyperparameter Search Details

227 We apply the PCA-based hyperparameter search algorithm (HSA) developed in Section 3.3 on several
228 variants of ResNet models. In the first part, we consider the large-scale classification scenarios.
229 Several standard ResNets are trained on ImageNet first and their architectures are shown in Table 1.
230 The HSA for searching BlkSConv architectures is only applied to conv3_x, conv4_x, conv5_x
231 layers of these standard ResNets. Next, the search hyperparameters $\alpha_v, \alpha_c, \alpha_s$ are set as 0.5 or

Table 1: ResNet architectures used in the first part of the experiment on ImageNet, Stanford Dogs, and Oxford 102 Flowers. The PCA-based HSA is applied to conv3_x, conv4_x, conv5_x layers and the corresponding convolutional kernel is replaced by the BlkSConv module found by HSA.

ResNet-10 (L=1), ResNet-18 (L=2), ResNet-26 (L=3)					
Layers Names	Output Size	ResNet	Applying HSA	e.g.(ResNet-10)	
conv1	$112 \times 112 \times 64$	$7 \times 7, 64, \text{stride } 2$	No		
max pool	$56 \times 56 \times 64$	$3 \times 3, \text{stride } 2$			
conv2_x	$56 \times 56 \times 64$	$\begin{bmatrix} 3 \times 3, & 64 \\ 3 \times 3, & 64 \end{bmatrix} \times L$	No		
conv3_x	$28 \times 28 \times 128$	$\begin{bmatrix} 3 \times 3, & 128 \\ 3 \times 3, & 128 \end{bmatrix} \times L$	Yes	conv-s5t2	
conv4_x	$14 \times 14 \times 256$	$\begin{bmatrix} 3 \times 3, & 256 \\ 3 \times 3, & 256 \end{bmatrix} \times L$	Yes	conv-s5t2	
conv5_x	$7 \times 7 \times 512$	$\begin{bmatrix} 3 \times 3, & 512 \\ 3 \times 3, & 512 \end{bmatrix} \times L$	Yes	conv-s1t1	
average pool	$1 \times 1 \times 512$	7×7			
fully connected	1000	$512 \times 1000 \text{ fc}$			

Table 2: Performance results for BlkSConv-based ResNet-18 and ResNet-26 on ImageNet.

$(\alpha_v, \alpha_c, \alpha_s, SS)$	ResNet-18 on ImageNet			ResNet-26 on ImageNet		
	Accuracy	P_ratio	MA_ratio	Accuracy	P_ratio	MA_ratio
(0.50, 0.50, 0.50, max)	69.922	0.4065	0.4614	72.038	0.4241	0.4618
(0.50, 0.75, 0.75, max)	69.782	0.6014	0.6597	72.326	0.6308	0.6722
(0.50, 0.50, 0.50, min)	67.572	0.1264	0.2700	69.970	0.1250	0.2668
(0.50, 0.75, 0.75, min)	67.540	0.1246	0.2986	69.922	0.1243	0.2910
Standard (replaced layers)	70.728	10.8M	1213.8M	72.604	17.03M	1907.4M

232 0.75. It is possible that the feasible set $H_{\alpha_v, \alpha_c, \alpha_s}$ is empty. In this case, the corresponding standard
 233 convolutional layer is not replaced and is denoted by conv as shown in Table 1. Moreover, the
 234 proposed HSA has two selection strategies: one is based on the largest parameter size, denoted by
 235 $SS = \max$, and the other is based on the smallest parameter size, denoted by $SS = \min$ as shown in
 236 Table 2. The selected BlkSConv, $sitj$, which means i basis block vectors and j depth of the blocks.

237 The feasible set $H_{\alpha_v, \alpha_c, \alpha_s}$ is likely to be empty when the parameter α_v is large. In the case that α_v
 238 is large, it often requires many principal components to accumulate enough explained variance and thus
 239 this causes large numbers of parameters or MAdds. Therefore, the feasible set $H_{\alpha_v, \alpha_c, \alpha_s}$ is probably
 240 empty when we further require small α_c and α_s . On the other hand, the parameter α_v cannot be too
 241 small because the prediction performance of the network is highly proportional to the amount of the
 242 accumulated variance as discussed in Section 3.3 where we will demonstrate it in the ablation study
 243 of this section. For the above reason, we only present the results for ResNet-18 and ResNet-26 on
 244 ImageNet under the setting that $\alpha_v = 0.5$ which are shown in Table 2.

245 In the second part, we consider the small-scale classification on CIFAR10/100. We use the standard
 246 ResNet-20 and ResNet-56 as the experimental models where the architectures are slightly modified
 247 to suit the small-scale images. The proposed HSA is only applied to conv4_x layers of these two
 248 standard ResNets. More BlkSConv-based architecture search results can be found in the appendix.

249 4.2 Performance on large-scale classification: ImageNet

250 To evaluate the performance of BlkSConv-based models in large-scale recognition, we conduct
 251 experiments on ImageNet(22). Each model takes 3 days to be trained on a single GPU (Nvidia
 252 Tesla V100). ImageNet contains nearly 1.3M training images and 50,000 testing images. For

Table 3: Comparison among the BlkSConv-based and Standard ResNet on ImageNet and CIFAR.

Dataset	Models	Accuracy	Parameters	MAdds
ImageNet	ResNet-10 standard	63.386	4.64M	520M
	BlkSConv-ResNet-18 (0.5, 0.5, 0.5, max)	69.922	4.39M	560M
	BlkSConv-ResNet-26 (0.5, 0.5, 0.5, min)	69.970	2.13M	509M
CIFAR 100	ResNet-20 standard	67.994	202752	12.97M
	BlkSConv-ResNet-20 (0.5, 0.5, 0.5, max)	67.078	72704	5.04M
	BlkSConv-ResNet-56 (0.5, 0.5, 0.5, min)	69.994	149440	10.61M

Table 4: Performance results for BlkSConv-based ResNet-56 on CIFAR10/100.

$(\alpha_v, \alpha_c, \alpha_s, SS)$	ResNet-56 on CIFAR 10			ResNet-56 on CIFAR 100		
	Accuracy	P_ratio	MA_ratio	Accuracy	P_ratio	MA_ratio
(0.5, 0.5, 0.5, max)	93.372	0.3734	0.3829	70.636	0.3734	0.3829
(0.5, 0.75, 0.75, max)	93.338	0.6196	0.6292	70.668	0.6196	0.6292
(0.5, 0.5, 0.5, min)	93.324	0.2335	0.2462	69.994	0.2316	0.2570
(0.5, 0.75, 0.75, min)	93.324	0.2335	0.2462	69.994	0.2316	0.2570
Standard (replaced layers)	93.218	645120	41.28M	70.998	645120	41.28M

the experimental setup, ResNet-10, ResNet-18, and ResNet-26 are trained on ImageNet under the following setting. The number of epochs is 100 and the batch size is 256. SGD is used as the optimizer and the initial learning rate, the momentum, and the weight decay are set to 0.1, 0.9, and 10^{-4} , respectively. The learning rate is scheduled to decay by a factor of 0.1 at epochs 30, 60, and 90. We augment the data via random resized crop to 224px and random horizontal flip. The performance results are shown in Table 2, More experimental results can be found in the appendix.

On one hand, let us focus the cases that $\alpha_v = 0.5$ and $SS = \max$ in Table 2. The prediction accuracies of the selected BlkSConv-based models and the standard model are close within 1%. It confirms our expectation that BlkSConv-based models have smaller parameter sizes and fewer MAdds than standard models while preserving prediction performance if the proposed HSA adopts a selection strategy based on the maximum parameter size.

On the other hand, let us consider the case that $\alpha_v = 0.5$ and $SS = \min$ in Table 2. The parameters and MAdds of the BlkSConv-based models are only 12.6% and 29.8% of the standard model while the gap of their prediction accuracies is about 3%. We adopt an interesting way based on restricting the parameter size and MAdds to interpret the advantage of the generated BlkSConv-based models where the selection strategy SS is set as min. We also compare the standard ResNet-10, the BlkSConv-based ResNet-18, and the BlkSConv-based ResNet-26 in Table 3 where the parameter sizes or MAdds of three given models are similar. The BlkSConv-based ResNet-26 with parameter (0.5, 0.5, 0.5, min) and the BlkSConv-based ResNet-18 with parameter (0.5, 0.5, 0.5, max) greatly outperform the standard ResNet-10 where both the BlkSConv-based models lead to an accuracy gain of at least 6.5%. In addition, the BlkSConv-based ResNet-26 with parameter (0.5, 0.5, 0.5, min) only has half the parameter size of the BlkSConv-based ResNet-18 with parameter (0.5, 0.5, 0.5, max).

4.3 Performance on small-scale classification: CIFAR 10/100

The performance results are shown in Table 4. The BlkSConv-based ResNet-56 models have much smaller model sizes and fewer MAdds than the standard model while all BlkSConv-based variants outperform the standard model on CIFAR 10 and have comparable accuracies on CIFAR 100. In the bottom of Table 3, the BlkSConv-based ResNet-20 with (0.5, 0.5, 0.5, max) and the standard ResNet-20 both have a comparable accuracy while the BlkSConv-based ResNet-20 model is compressed 64% of the parameter size and MAdds is decreased 61% compared to the standard ResNet-20 model. Furthermore, the BlkSConv-based ResNet-56 with (0.5, 0.5, 0.5, min) and the standard ResNet-20 model both have similar parameter sizes and MAdds while the BlkSConv-based ResNet-56 model has an accuracy gain of 2%. More experimental results can be found in the appendix.

Table 5: Performance comparison among BlkSConv-based and the standard ResNet-18 models.

$(\alpha_v, \alpha_c, \alpha_s, SS)$	Stanford Dogs			Oxford 102 Flowers		
	Accuracy	P_ratio	MA_ratio	Accuracy	P_ratio	MA_ratio
(0.5, 0.5, 0.5, max)	53.005	0.5327	0.5394	65.546	0.5327	0.5394
(0.5, 0.75, 0.75, max)	53.359	0.7277	0.7377	64.567	0.7277	0.7377
(0.5, 0.5, 0.5, min)	53.159	0.3273	0.4006	65.289	0.4835	0.5179
(0.5, 0.75, 0.75, min)	53.615	0.3171	0.4611	63.217	0.4743	0.5872
Standard (replaced layers)	52.436	10.8M	1213.8M	62.238	10.8M	1213.8M

Table 6: Results on Stanford Dogs for different α_v with $SS = \min$.

ResNet-18 on Stanford Dogs ($\alpha_v, \alpha_c = 0.75, \alpha_s = 0.75, SS = \min$)							
α_v	Accuracy	P_ratio	MA_ratio	α_v	Accuracy	P_ratio	MA_ratio
0.0	50.918	0.0397	0.1781	0.3	52.839	0.1074	0.2377
0.1	50.499	0.0449	0.1777	0.4	52.035	0.2751	0.4684
0.2	51.040	0.0767	0.2458	0.5	53.615	0.3171	0.4611
				Standard	52.436	10.8M	1213.8M

285 **4.4 Performance on fine-grained classification**

286 We conduct experiments for fine-grained recognition on two datasets Stanford Dogs and Oxford
 287 102 Flowers. For the experimental setup, the standard ResNet-18 and its BlkSConv-variants are all
 288 trained from scratch by augmenting data through random crops, horizontal flips, and random gamma
 289 transform. We use SGD as the optimizer and the initial learning rate, the moment, and the weight
 290 decay are set to 0.1, 0.9, and 10^{-4} , respectively. The number of epochs is 200, and the learning rate
 291 is scheduled to decay at epochs 100, 150, and 200 by a factor of 0.1. The proposed BlkSConv-based
 292 ResNet-18 models significantly outperform the standard ResNet-18 model both on Stanford Dogs
 293 and Oxford 102 Flowers as shown in shown Table 5.

294 **4.5 Ablation Study: Necessity to have large explained variance**

295 Here, we demonstrate how the variance hyperparameter α_v affects the prediction accuracy of
 296 BlkSConv-based CNNs. We use ResNet-18 as the experimental model. After training the stan-
 297 dard ResNet-18 on Stanford Dogs, the next goal is to find several BlkSConv-variants of ResNet-
 298 18 all of which have different explained variances such that their accuracies can be compared.
 299 Note that $H_{a,0.75,0.75} \subseteq H_{b,0.75,0.75}$ for any a, b with $a \geq b$. Based on this observation, the
 300 model in $H_{b,0.75,0.75}$ which has the smallest parameter size is likely to have a small explained
 301 variance as well. Therefore, the selection strategy of the proposed HSA is set by $SS = \min$
 302 in order to select several BlkSConv-based ResNet-18 models with different explained variances.
 303 Now we apply the proposed HSA to the trained ResNet-18 under six search hyperparameters
 304 $\{(\alpha_v, 0.75, 0.75, \min) : \alpha_v = 0.0, 0.1, 0.2, \dots, 0.5\}$. The comparison result is shown in Table 6. It
 305 can be seen that the accuracy of the BlkSConv-based model is greater than that of the standard model
 306 only when the variance hyperparameter α_v is large enough, that is $\alpha_v \geq 0.5$.

307 **5 Conclusion**

308 In this paper, we introduce the block-wise separable convolutions (BlkSConv) to replace standard
 309 convolutions. An efficient implementation of the BlkSConv operation via a combination of pointwise
 310 and group-wise convolutions is also given. Moreover, we also propose an efficient hyperparameter
 311 search algorithm based on principal component analysis in order to select an optimal BlkSConv-based
 312 convolutional network under certain constraints on model size and model efficiency. Finally, the
 313 experimental results demonstrate the advantage of the BlkSConv-based CNN models selected by the
 314 proposed hyperparameter search algorithm.

315 **References**

- 316 [1] M. Denil, B. Shakibi, L. Dinh, M. A. Ranzato, and N. De Freitas, "Predicting parameters in deep
317 learning," in *Advances in neural information processing systems (NIPS)*, pages 2148–2156, 2013.
- 318 [2] E. Denton, W. Zaremba, J. Bruna, Y. Le-Cun, and R. Fergus, "Exploiting linear structure within
319 convolutional networks for efficient evaluation," in *Advances in neural information processing
320 systems (NIPS)*, pages 1269–1277, 2014.
- 321 [3] J. Guo, Y. Li, W. Lin, Y. Chen, and J. Li, "Network decoupling: From regular to depthwise
322 separable convolutions," arXiv preprint arXiv:1808.05517, 2018
- 323 [4] D. Haase and M. Amthor, "Rethinking Depthwise Separable Convolutions: How Intra-Kernel
324 Correlations Lead to Improved MobileNets," in *Proceedings of the IEEE Conference on Computer
325 Vision and Pattern Recognition (CVPR)*, pages 14600–14609, 2020.
- 326 [5] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. CoRR,
327 abs/1512.03385, 2015.
- 328 [6] S. Han, H. Mao, and W. J. Dally. "Deep compression: Compressing deep neural networks
329 with pruning, trained quantization and huffman coding," in *Proceedings of the International
330 Conference on Learning Representations*, 2015.
- 331 [7] A. G. Howard, M. Sandler, G. Chu, L.-C. Chen, M. Tan, W. Wang, et al., "Searching for
332 mobilenetv3", in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019.
- 333 [8] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, et al., "Mobilenets:
334 Efficient convolutional neural networks for mobile vision applications", 2017.
- 335 [9] Y. He, G. Kang, X. Dong, Y. Fu, and Y. Yang, "Soft filter pruning for accelerating deep Convolu-
336 tional Neural Networks," 2018.
- 337 [10] K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," in
338 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
339 770–778, 2016.
- 340 [11] Y. He, X. Zhang, and J. Sun, "Channel pruning for accelerating very deep Neural Networks," in
341 *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 1389–1397,
342 2017.
- 343 [12] M. Jaderberg, A. Vedaldi, and A. Zisserman, "Speeding up convolutional neural networks with
344 low rank expansions," in *Proceedings of the British Machine Vision Conference. (BMVA)*, 2014.
- 345 [13] J. Jin, A. Dundar, and E. Culurciello, "Flattened convolutional neural networks for feedforward
346 acceleration," arXiv preprint arXiv:1412.5474, 2014
- 347 [14] A. Khosla, N. Jayadevaprakash, B. Yao, and F.-F. Li. "Novel dataset for fine-grained image
348 categorization," in *first Workshop on Fine-Grained Visual Categorization, IEEE Conference on
349 Computer Vision and Pattern Recognition*, Colorado Springs, CO, June 2011.
- 350 [15] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. Technical
351 report, Citeseer, 2009.
- 352 [16] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H.P. Graf, "Pruning filters for efficient ConvNets",
353 2016.
- 354 [17] J.H. Luo, J. Wu, and W. Lin, "Thinet: A filter level pruning method for deep Neural Network
355 compression," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*,
356 pages 5058–5066, 2017.
- 357 [18] N. Ma, X. Zhang, H.-T. Zheng and J. Sun, "Shufflenet v2: Practical guidelines for efficient cnn
358 architecture design", in *Proceedings of the European Conference on Computer Vision (ECCV)*,
359 pages 116–131, 2018.

- 360 [19] M.-E. Nilsback and A. Zisserman. "Automated flower classification over a large number of
361 classes", In *2008 Sixth Indian Conference on Computer Vision, Graphics and Image Processing*,
362 pages 722–729, IEEE, 2008.
- 363 [20] E. Park, S. Yoo, and P. Vajda, "Value-aware quantization for training and inference of Neural
364 Networks", in *Proceedings of the European Conference on Computer Vision (ECCV)*, pages
365 580–595, 2018.
- 366 [21] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-Net: ImageNet classification
367 using binary Convolutional Neural Networks", in *Proceedings of the European Conference on
368 Computer Vision (ECCV)*, pages 525–542, 2016.
- 369 [22] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A.
370 Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International
371 journal of computer vision*, 115(3):211–252, 2015.
- 372 [23] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov and L. Chen, "MobileNetV2: Inverted Residuals
373 and Linear Bottlenecks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern
374 Recognition (CVPR)*, pages 4510–4520, 2018.
- 375 [24] A. See, M.-T. Luong, and C. D. Manning, "Compression of neural machine translation models
376 via pruning," in *Proceedings of the 20th SIGNLL Conference on Computational Natural Language
377 Learning*, pages 291–301, 2016.
- 378 [25] L. Sifre and S. Mallat. Rigid-motion scattering for image classification. PhD thesis, 2014.
- 379 [26] W. Shang, K. Sohn, D. Almeida, and H. Lee, "Understanding and improving convolutional neu-
380 ral networks via concatenated rectified linear units," in *Proceedings of international conference
381 on machine learning (ICML)*, pages 2217–2225, 2016.
- 382 [27] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image
383 recognition," in *Proceedings of the International Conference on Learning Representations*, 2015.
- 384 [28] P. Singh, V.K. Verma, P. Rai, and V. Nambodiri, "Leveraging filter correlations for deep model
385 compression," 2018.
- 386 [29] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and
387 A. Rabinovich, "Going deeper with convolutions" in *Proceedings of the IEEE Conference on
388 Computer Vision and Pattern Recognition (CVPR)*, pages 1–9, 2015.
- 389 [30] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens and Z. Wojna, "Rethinking the Inception Archi-
390 tecture for Computer Vision," in *Proceedings of the IEEE Conference on Computer Vision and
391 Pattern Recognition (CVPR)*, pages 2818–2826, 2016.
- 392 [31] M. Tan, B. Chen, R. Pang, V. Vasudevan, and Q. V. Le, "Mnasnet: Platform-aware neural
393 architecture search for mobile," arXiv preprint arXiv:1807.11626, 2018.
- 394 [32] B. Wu, X. Dai, P. Zhang, Y. Wang, F. Sun, Y. Wu, Y. Tian, P. Vajda, Y. Jia, and K. Keutzer,
395 "Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search,"
396 arXiv preprint arXiv:1812.03443, 2018.
- 397 [33] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng, "Quantized convolutional neural networks
398 for mobile devices," in *Proceedings of the IEEE Conference on Computer Vision and Pattern
399 Recognition*, pages 4820–4828, 2016.
- 400 [34] L. Xie and A. Yuille, "Genetic cnn," in *Proceedings of the IEEE International Conference on
401 Computer Vision (ICCV)*, pages 1379–1388, 2017.
- 402 [35] D. Zhang, J. Yang, D. Ye, and G. Hua, "LQ-Nets: Learned quantization for highly accurate
403 and compact deep Neural Networks," in *Proceedings of the European Conference on Computer
404 Vision (ECCV)*, pages 365–382, 2018.
- 405 [36] X. Zhang, X. Zhou, M. Lin and J. Sun, "ShuffleNet: An Extremely Efficient Convolutional
406 Neural Network for Mobile Devices," in *Proceedings of the IEEE Conference on Computer
407 Vision and Pattern Recognition (CVPR)*, pages 6848–6856, 2018.

- 408 [37] L. Zhu, R. Deng, M. Maire, Z. Deng, G. Mori, and P. Tan, "Sparsely aggregated Convolutional
409 Networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pages
410 186–201, 2018.
- 411 [38] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," in *Proceedings*
412 *of the International Conference on Learning Representations*, 2017.
- 413 [39] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable
414 image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern*
415 *Recognition (CVPR)*, pages 8697–8710, 2018.

416 **Checklist**

- 417 1. For all authors...
- 418 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
419 contributions and scope? [Yes] Section 1
- 420 (b) Did you describe the limitations of your work? [Yes] Section 3.2
- 421 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 422 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
423 them? [Yes]
- 424 2. If you are including theoretical results...
- 425 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 426 (b) Did you include complete proofs of all theoretical results? [N/A]
- 427 3. If you ran experiments...
- 428 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
429 mental results (either in the supplemental material or as a URL)? [Yes] Supplemental
430 material
- 431 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
432 were chosen)? [Yes] Section 4.1, Section 4.2, Section 4.4, Appendix
- 433 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
434 ments multiple times)? [No]
- 435 (d) Did you include the total amount of compute and the type of resources used (e.g., type
436 of GPUs, internal cluster, or cloud provider)? [Yes] Section 4.2, Appendix
- 437 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 438 (a) If your work uses existing assets, did you cite the creators? [Yes] Section 4
- 439 (b) Did you mention the license of the assets? [Yes] Appendix
- 440 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 441 (d) Did you discuss whether and how consent was obtained from people whose data you're
442 using/curating? [Yes] Appendix
- 443 (e) Did you discuss whether the data you are using/curating contains personally identifiable
444 information or offensive content? [Yes] Appendix
- 445 5. If you used crowdsourcing or conducted research with human subjects...
- 446 (a) Did you include the full text of instructions given to participants and screenshots, if
447 applicable? [N/A]
- 448 (b) Did you describe any potential participant risks, with links to Institutional Review
449 Board (IRB) approvals, if applicable? [N/A]
- 450 (c) Did you include the estimated hourly wage paid to participants and the total amount
451 spent on participant compensation? [N/A]