

MUNet: Macro Unit-based Convolutional Neural network for Mobile Devices

Dae Ha Kim, Seung Hyun Lee and Byung Cheol Song

Inha University

Incheon, South Korea

kdhht5022@gmail.com, lsh910703@gmail.com, bcsong@inha.ac.kr

Abstract

Deep neural networks perform better than traditional machine learning methods on various classification problems by producing good quality feature maps through successive convolution operation(s). However, when implementing a deep neural network in an embedded system or SoC for mobile devices, its large parameter size can be a significant burden on the internal memory design. In this paper, we propose a new deep neural network that reduces computation and the number of model parameters but maintains reasonable performance. The configuration of the proposed network is as follows: First, we propose a macro unit (MU) to reduce heavy computations and to learn various feature maps. Second, we employ asymmetric convolution of the well-known Inception network to further efficiently manipulate feature maps within the MU. Third, all the feature maps produced from MU(s) of each layer are concatenated and then the grouped feature map is distributed to all the MUs of the next layer for transferring richer information. Experimental results show that the proposed network achieves about 10% higher performance than DenseNet-BC in case of extremely small parameter size for CIFAR-100. The proposed network also has very few learning parameters and smaller floating point operations per second (FLOPS) than the other networks optimized for mobile devices such as MobileNet.

1. Introduction

Prior to the deep learning era, a lot of traditional machine learning approaches have been developed for solving various classification problems, and neural networks (NN) were one of them. Recently, deep NN techniques [1-3] that can be considered as a catalyst for this deep learning era have been devised. The end-to-end learning optimizes all processing steps simultaneously, and it leads to better performance and smaller systems. Because of outstanding performance beyond traditional machine learning techniques, deep learning is expanding its scope to a variety of applications, including computer vision, medical, and

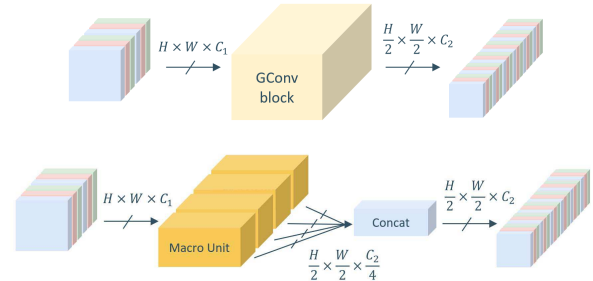


Figure 1: The main concept of the proposed method. *Up*: Learning scheme of group convolution block. *Down*: Learning scheme of proposed convolution block. Here H, W represent the horizontal and vertical sizes of the feature map, respectively and C_n represents the number of channels.

architectural applications. AlexNet [4] is known as the most famous convolutional NN (CNN) among early deep learning methods. Next, competition to develop deeper and better CNNs than AlexNet has become fierce, and so the well-known visual geometry group (VGG) network was born out of such competition [5]. However, since VGG, many researchers have encountered limitations in increasing network depth. He et al. caused sensation by presenting ResNet [7], which enables a deeper structure of the network using residual connections.

ResNeXt [8] and Xception [9] have also been proposed, which effectively learns the feature map by focusing on the convolution filter configuration rather than the feature map connection. Recently DenseNet [10], which is a network with differentiating connection methods between feature maps, based on the existing ResNet became famous.

On the other hand, extensive size of learning parameters and huge computational complexity make it difficult to implement existing CNNs with embedded software or SoC for mobile devices, drones, and social robots. To solve this problem, CNNs focusing on sparsity connection aimed at hardware-friendly structure have been developed [13, 14].

However, conventional CNNs, including [13] and [14], still have significant computational complexity and parameter size, making them cumbersome to implement. The fundamental problem of CNNs in terms of computational complexity and memory size is the use of so-called group convolution (GConv) [4], which requires a

large number of convolution filters. When constructing CNN with such group convolutions, it is advantageous to obtain various feature maps in spite of simple convolution structure. On the other hand, in the optimization process on the high dimension space, GConv may have some degree of redundancy. This causes unnecessary increase of parameter size and computational complexity [26, 35].

To solve the above-mentioned problem, MobileNet [16] and ShuffleNet [17] have adopted factorized convolution such as depth-wise separable convolution [6]. Nonetheless, the tradeoff between performance and parameter size/computation is still an unresolved issue in the research on CNN. The reason is that even if you use factorized convolution, many filters per convolution operation are still required. This can eventually lead to a redundancy in learning high dimension feature map [31].

This paper proposes an efficient CNN structure that shows satisfactory performance even with only a small number of convolution filters. As shown in Fig. 1, learning a feature map by dividing a single GConv block into multiple macro units (MUs) is a key concept of the proposed network. Each MU consists of much smaller number of convolution filters than GConv block, and furthermore it adopts separable convolution of Inception-V4 [15], which has considerable advantages in terms of parameter size and computation amount. Meanwhile, we propose a method to simultaneously concatenate feature maps generated from multiple MUs of each layer and distribute them to the next layer, as shown in Fig. 2, so that a rich feature map can be generated even with MUs composed of a small number of convolution filters. Therefore, the proposed network can learn sufficiently rich feature maps with relatively a small number of convolution filters.

We verified the classification performance of the proposed MU-based network (MUNet) on highly competitive benchmarking datasets, i.e., CIFAR, SVHN, and Tiny-ImageNet. For example, for CIFAR-100 dataset, MUNet shows a performance improvement of up to 10% in a small parameter region compared to the DenseNet-BC. Also, in the reasonably similar error rate range, MUNet exhibits more than 15 times less parameter size and fewer FLOPS than MobileNet.

2. Related Works

Novel convolution filter. Inception-V4 is an inception module that combines the results from multiple filters of different sizes. To reduce the computational complexity, the inception module adopts the asymmetric convolution instead of the two-dimensional (2D) convolution. Also, ResNeXt and Xception networks have shown that feature maps can be efficiently learned by implementing point-wise convolution and depth-wise convolution based on (1,1) convolution filter and (3,3) convolution filter.

Connection between feature maps. As mentioned earlier, ResNet is a way to increase the number of parameters and pursue a deeper network than AlexNet and VGG. Adding the by-pass connection to the existing VGG network is the key idea, which solved the gradient vanishing problem and helped the smooth transmission of feature map information. Thereafter, modified versions of ResNet appeared in which various methods were applied to ResNet [18, 19].

Beyond the ResNet family of algorithms, DenseNet has introduced two new concepts. First, DenseNet utilized residual connections of concatenate style rather than those of summation style, thus enabling rich learning of feature maps. Second, DenseNet has made connections between all the layers of the block module so that feature map and gradient information can be seamlessly conveyed across the entire network. In addition, FractalNet [11] is remarkable in that it builds up a deep network without residual connections and boosts overall performance.

Efficient Model Designs. MobileNet is based on depth-wise separable convolution which factorize a standard convolution into a depth-wise convolution and a 1x1 convolution called a point-wise convolution. Note that conventional convolution performs filtering and combining simultaneously. The depth-wise separable convolution of MobileNet splits this into two layers, i.e., one layer for filtering and the other layer for combining. So such a factorization has the effect of drastically reducing computation and model size. Therefore, MobileNet enables very fast inference time while showing reasonable performance.

ShuffleNet has a more advanced convolution structure than MobileNet. ShuffleNet is also based on a depth-wise separable convolution, but realizes a small network with few channels through the channel shuffle technique. So ShuffleNet has similar performance to MobileNet, but its parameter size and computation becomes smaller than that of MobileNet.

On the other hand, NASNet model [20], which has been attracting attention recently, is not a network designed to have a fundamentally small parameter size, so we do not use NASNet as a benchmarking network of MUNet in this paper.

3. Approach

As shown in Fig. 2, MU is a crucial component of the proposed network. Feature maps learned in independent MUs of each layer are concatenated and distributed to MUs of the next layer. After repeating the above steps in several layers as shown in the figure, the final feature map is produced and then global pooling is applied. Finally, the vector information is obtained through the fully connected (FC) layer of the same level as the class of the

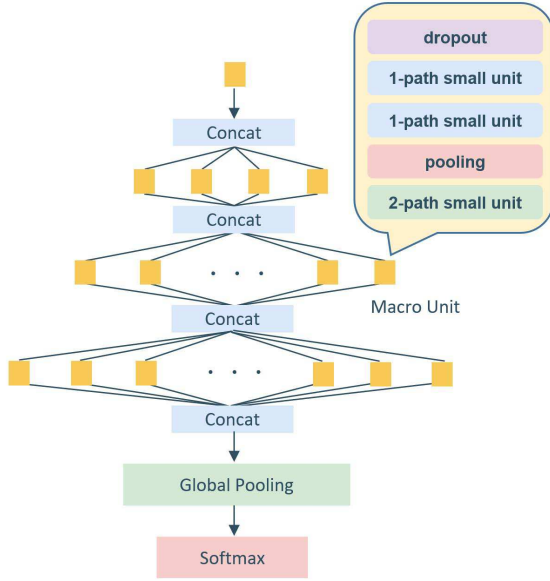


Figure 2: Overall framework of the proposed network.

corresponding data set, and the softmax function is taken and the final prediction is performed.

3.1. Macro unit

This section describes the structure of the MU, i.e., a key element of the proposed network. The overall MU configuration is shown in Fig. 3(a). After the network regularization process such as dropout [12], feature map learning is performed through multiple 1-path small units. The primary function of the 1-path small unit(s) is to increase the number of internal layers in the MU to prevent performance degradation. We then reduce the size of the feature map through the pooling process. Finally, the 2-path small unit is applied. Since the 2-path small unit is constructed based on the dynamic connections of 1D convolutions in different directions, which were also proposed in the inception module of [15], and symmetric convolution, it outputs feature maps with more abundant characteristics than 1-path small unit do. As a result, the MUNet having a small number of convolution filters can significantly reduce the computational complexity and the parameter size compared to the conventional GConv-based networks.

3.2. Small units

As in Fig. 3(a), the MU uses two kinds of small units. First, the 1-path small unit is used for learning a richer feature map by increasing the number of layers in the MU, and has the detailed structure of Fig. 3(b). The number of 1-path small units in the MU is considered a hyper parameter of the network, such as dropout, which is covered

in Section 4.6. 1D convolutions of $(1, N)$ and $(N, 1)$ are performed after the 2D convolution of (N, N) , and then the results are summed as shown in the figure. As mentioned above, the reason for using asymmetric convolution filters such as $(1, N)$ and $(N, 1)$ is to reduce the computational cost of the convolution operation. The computational cost of the standard 2D convolution is as follows.

$$C_1 \times C_2 \times N \times N \times D_1 \times D_2 \quad (1)$$

where C_1 and C_2 represent the number of input channels and the number of output channels, respectively, and $N \times N$ represents the convolution kernel size, and $D_1 \times D_2$ indicates the feature map size. The ratio of the computational complexity of the asymmetric 1D convolutions to the 2D convolution computation of Eq. (1) is calculated as follows.

$$\frac{C_1 \times C_2 \times (N + N) \times D_1 \times D_2}{C_1 \times C_2 \times N \times N \times D_1 \times D_2} = \frac{2N}{N^2} \quad (2)$$

That is, the parameter reduction effect is about $\frac{2}{N}$ times.

The structure of the 2-path small unit is shown in Fig. 3(c). The different thing from 1-path small unit is the number of output paths. The 2-path small unit has two output paths. Each output is to plus the asymmetric 1D convolution result and the standard 2D convolution result. Feature maps learned by 1D convolutions with different directions, i.e., asymmetric convolutions, are transferred to the MUs of the next layer, so that it is possible to generate feature maps with more abundant characteristics.

3.3. Comparison of GConv Block and MU Block

In this section, we compare the computational complexities of the GConv block and the corresponding MU-based block in Fig. 1. The specifications of the two network architectures are shown in Table 1. The computation cost of the GConv block is obtained from Table 1 and Eq. (1), which is as follows.

$$\begin{aligned} \text{Conv1} & : C_1 \times C_2 \times 3^2 \times D_1 \times D_2 + \\ \text{Conv2,3,4} & : (C_2 \times C_2 \times 3^2 \times D_1 \times D_2) \times 3 \end{aligned} \quad (3)$$

$$\begin{aligned} \text{Conv5,6} & : (C_2 \times C_2 \times 3^2 \times D_1/2 \times D_2/2) \times 2 \\ & C_2 \times C_2 \times 3^2 \times D_1 \times D_2 = \alpha \end{aligned} \quad (4)$$

If the common term α of Eq. (4) is substituted, the total computational complexity of the GConv block is expressed by Eq. (5).

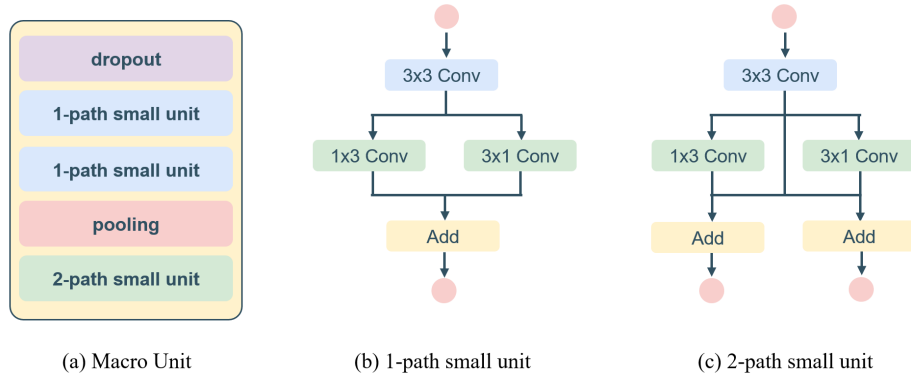


Figure 3: Architecture of macro unit (MU). (a) Overall architecture. (b) 1-path small unit. (c) 2-path small unit.

$$\left(\frac{7}{2} + \frac{C_1}{C_2}\right)\alpha \quad (5)$$

The computational complexity of the MU-based block corresponding to the GConv block is shown in Eq. (6).

$$\left(\frac{11}{64} + \frac{C_1}{4C_2}\right)\alpha \times 4 \quad (6)$$

Based on Eqs. (5) and (6), we can find that the MU-based block has about 70% reduction in computation cost over the GConv block, assuming that the number of convolution filters doubles in the general CNN layer hierarchy.

3.4. Interface of concatenation and distribution

Feature maps generated from several MUs of each layer are concatenated and uniformly distributed to the MUs of the next layer, as shown in Fig. 2, so that MUs with a small number of convolution filters can learn rich feature maps sufficiently. For example, at the bottom of Fig. 1, we can observe that the four MU outputs are concatenated and then delivered to the next layer.

Therefore, MUNet has an advantage that it has very small parameter size and computation amount with comparable to GConv-based networks.

3.5. Implementation details

As shown in Fig. 2, the proposed network consists of four MU-layers. The number of MUs from the top MU-layer is set to 1, 4, 8, and 16, respectively. The first MU-layer has only one MU, which is composed of only 2-path small unit, not dropout and pooling layer considering the input characteristics of the network. The second and third MU-layers have the structure shown in Fig. 3(a). Since the last fourth MU-layer no longer needs to increase the depth of the network, we only use the 1-path small unit once after the dropout.

Layer	GConv block	MU block
Input	$H \times W \times C_1$	$H \times W \times C_1$
Conv1	(3, 3) conv.	(3, 3) conv.
Conv2	(3, 3) conv.	(1, 3) & (3, 1) conv.
Conv3	(3, 3) conv.	(3, 3) conv.
Conv4	(3, 3) conv.	(1, 3) & (3, 1) conv.
Pooling	(2, 2) max pool	(2, 2) max pool
Conv5	(3, 3) conv.	(3, 3) conv.
Conv6	(3, 3) conv.	(1, 3) & (3, 1) conv.
Concat	-	Filters of 4 MUs
Output	$\frac{H}{2} \times \frac{W}{2} \times C_2$	$\frac{H}{2} \times \frac{W}{2} \times C_2$

Table 1: Network architecture between GConv block and MU. Here (3,1) & (3,1) conv. indicates asymmetric convolutions. Both GConv block and MU assume C_1 input feature maps. C_2 convolution filters are used for GConv block and $C_2/4$ convolution filters are used for each MU. After pooling, the feature map size is halved.

Meanwhile, max pooling is used in the pooling step, and the dropout rates of the remaining MU-layers except for the first MU-layer are 0.2, 0.5, and 0.5, respectively.

4. Experimental Results

We experimentally demonstrated the classification performance of MUNet for CIFAR, SVHN, and Tiny-ImageNet datasets, which are widely used as benchmarks for classification problems. We also verified the model adaptability beyond classification task.

Section 4.1 discusses the datasets used in the experiments and details of the learning method. Section 4.2 compares the performance of the proposed MUNet against the latest techniques via classification problem. Section 4.3 compares the computational complexity as well as the memory cost. Section 4.4 verified the adaptability of the proposed

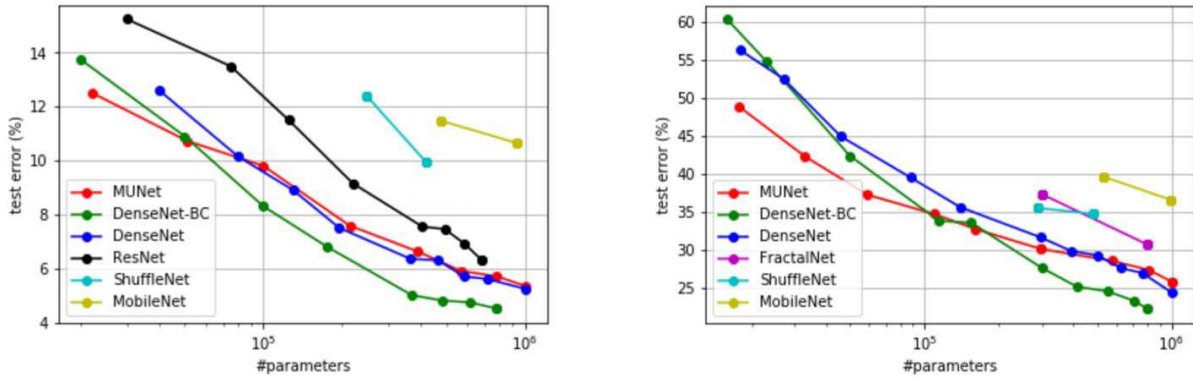


Figure 4: Comparison in terms of the parameter efficiency for CIFAR datasets. *Left*: CIFAR-100+. *Right*: CIFAR-10.

network to the surrounding environment. Section 4.5 experimentally justified the structural characteristics of the MUNet. Finally, Section 4.6 shows how to choose the optimal hyper-parameters of MUNet.

4.1. Datasets and training details

For fair evaluation of the proposed scheme, we used four popular datasets: CIFAR [23], SVHN [24], Tiny ImageNet [25], and down-sampled ImageNet [33]. The latter two datasets are from ILSVRC-2012 [34]. In order to evaluate the classification performance, CIFAR, SVHN, and Tiny ImageNet are used. The reason for including Tiny ImageNet here is to prove that MUNet is a network that is sufficiently applicable to the actual environment. Down-sampled ImageNet was adopted to verify fine tuning performance of the MUNet in Section 4.4.

As the augmentation technique of the datasets, the only horizontal mirroring and translation (uniform offsets in from -4 to 4) used in [7, 10, 11] was applied.

To optimize the proposed network model, the stochastic gradient descent (SGD) [27] method was used. Also, the weight decay of 10^{-4} and Nesterov momentum [28] of 0.9 without dampening were employed. For weight initialization, we adopted the distribution used in [29]. When learning with CIFAR and SVHN datasets, the total number of epochs was set to 300 and the batch size was set to 128. When learning with the Tiny ImageNet dataset, the total number of epochs was set to 300, and the batch size was set to 256, and the dropout was excluded here. Finally, we achieved the normalization effect of the feature map during learning process by using batch normalization [30] after the activation function.

The performance of a network was evaluated by the average of the results of three times experiments using the weights of the network after the final learning. We also employed the Tensorflow [21] based deep learning library, and used GPU named NVIDIA GeForce GTX 1080 Ti to finish the training process.

Model	Dataset	Params	Error rate (%)
ResNet [7]	CIFAR-10	0.8M	6.97
DenseNet [10]		1.0M	5.24
DenseNet (ours)		1.0M	5.12
DenseNet-BC [10]		0.8M	4.51
DenseNet-BC (ours)		0.8M	4.41
MUNet		1.0M	5.35
FractalNet [11]		0.8M	30.71
DenseNet [10]	CIFAR-100	1.0M	24.42
DenseNet (ours)		1.0M	24.04
DenseNet-BC [10]		0.8M	22.27
DenseNet-BC (ours)		0.8M	22.01
MUNet		1.0M	25.85

Table 2: Performance comparison of various CNN methods for the parameter size of about 1M in terms of error rates.

4.2. Classification performance evaluation

This section compared MUNet with state-of-the-art CNN techniques such as DenseNet, DenseNet-BC, ResNet, FractalNet, MobileNet, and ShuffleNet in terms of classification accuracy. The experimental results of MobileNet and ShuffleNet were from our own implementations, and those of the other networks were directly extracted from the corresponding papers. In case of DenseNet, because no numerical results for parameter sizes of about 500,000 or less are available in [10], we implemented the DenseNet and plotted the results in the right of Fig. 4. In order to ensure the reliability of our implementation(s), Table 2 compared it with the results from [10] for the parameter sizes of about 1M. The results from [7] and [11] were also compared together in the table. Note that according to [7, 10, 11], 1 M is almost the lower bound of the parameter size where the state-of-the-art

Model	Dataset	Params	Error rate (%)	Inference time (<i>Std</i>)	FLOPS
1.0 MobileNet	CIFAR-10	1.88M	10.03	37.42 (0.36)	51.44M
ShuffleNet × 1 ($g = 1$)		0.42M	9.93	39.17 (0.45)	43.39M
MUNet		0.10M	9.80	41.85 (0.35)	40.30M
1.0 MobileNet	CIFAR-100	1.97M	36.05	37.76 (0.19)	51.62M
ShuffleNet × 1 ($g = 1$)		0.47M	34.76	39.64 (0.18)	44.49M
MUNet		0.11M	34.72	41.39 (1.38)	40.32M
1.0 MobileNet	SVHN	1.88M	4.98	96.53 (0.26)	51.44M
ShuffleNet × 1 ($g = 1$)		0.42M	4.68	102.81 (0.24)	43.39M
MUNet		0.10M	4.41	104.72 (0.30)	40.30M
1.0 MobileNet	Tiny ImageNet	3.17M	56.48	45.40 (0.13)	79.33M
ShuffleNet × 1 ($g = 1$)		0.70M	56.18	59.49 (0.84)	59.49M
MUNet		0.13M	56.13	59.93 (0.49)	50.83M

Table 3: Comparison in terms of parameter size, inference time, and FLOPS assuming the same error rates. *Std* in the fifth column represents the standard deviation.

methods provide effective performance. We can say that our implementation is reliable. In addition, we can find that MUNet shows comparable results to the cutting-edge CNNs for CIFAR datasets.

Next, since MUNet pursues a network with a small parameter size, we evaluated its performance, focusing on a region where the learning parameter size is small to some extent. Because we can say that the parameter size range below 1M is a reasonable range based on the datasets we adopted (see Fig. 4), we concentrate on the comparison for the parameter size range of less than 100,000, i.e., an intermediate point that we pay attention to. As shown in Fig. 4, MUNet shows consistently superior performance to ResNet, FractalNet, and DenseNet. The performance of MUNet is somewhat lower than that of DenseNet-BC, but the smaller the parameter size, the more comparable the performance.

On the other hand, one of the points we should pay attention to is the pattern of performance change according to the parameter size. MUNet shows a much lower performance degradation rate than the other networks in the parameter size area below 1.0M. Especially, in case of very small parameter size area of less than 0.1M, the proposed network shows rather superior performance than DenseNet-BC. That is, when comparing MUNet and DenseNet-BC in a very small parameter size area, MUNet has a performance advantage of up to 1.5% based on CIFAR-10 and up to 10% based on CIFAR-100. For CIFAR-100, FractalNet showed an error rate of about 37% at parameter size of 0.3M, while MUNet showed about 37% error rate at parameter size of 0.059M. In other words, assuming the same classification accuracy, it is sufficient that MUNet is as small as 1/5 of the parameter size than FractalNet.

In addition, we can observe from Fig. 4 that MobileNet

and ShuffleNet perform significantly poorer than the other networks in terms of performance vs. parameter size. This is because MobileNet and ShuffleNet are the networks that are developed to reduce the amount of computation rather than to improve performance. In this aspect, the proposed MUNet should be noted that even if the parameter size decreases, the performance degradation rate is very slow.

4.3. Complexity analysis

The main purpose of MUNet is to derive a good tradeoff of performance and computation cost/parameter size to the extent that it can be implemented with embedded software or SoC for applications such as mobile devices. In this section, we analyze the computational cost and complexity of the proposed MUNet. In this experiment, the number of convolution filters in the first MU-layer of MUNet was set to 6 and the number of convolution filters of the remaining MU-layers was set to 4 for fair comparison. On the other hand, the benchmarking networks of MUNet were limited to MobileNet and ShuffleNet for convenience.

Table 3 shows the comparison of MUNet with three networks. We compared the parameter size (3rd column), inference time (5th column), and FLOPS (6th column) when they show experimentally similar classification accuracy, i.e., error rate (4th column). Inference times were measured in units of seconds on the desktop environment embedded with the Intel® i7-7700 CPU@3.6 GHz w/ 8 processors. FLOPS indicates the amount of computation, and it includes all operations such as convolution, pooling, and even multiplication, addition, and conditional statement for activation function. The inference time is measured as the sum of all test images in the dataset.

From Table 3, we can find that MUNet has a smaller parameter size than the other networks. For example, on a

Model	Dataset	Params	Error rate (fine tuned / original) (%)	FLOPS
1.0 MobileNet	CIFAR-10	1.88M	9.70 / 10.03	51.44M
MUNet		0.11M	9.03 / 9.54	49.58M
		0.10M	9.15 / 9.80	40.30M
1.0 MobileNet	CIFAR-100	1.98M	31.58 / 36.05	51.62M
MUNet		0.12M	32.10 / 34.58	49.59M
		0.11M	33.54 / 34.72	40.32M

Table 4: Fine tuning performance comparison between MobileNet and MUNet.

Method	Error rate (%)	Params	FLOPS [Mult. Add If]
Single MU per layer	9.01	0.49M	203.72M [101.75 101.16 0.81]
Multiple MUs per layer (MUNet)	9.80	0.10M	40.30M [20.11 19.88 0.31]

Table 5: Performance comparison according to MU-layer structure for CIFAR-10.

CIFAR dataset, MUNet achieves about four times the parameter reduction in comparison to ShuffleNet.

On the other hand, while MUNet has always smaller FLOPS than ShuffleNet and MobileNet for all datasets, its inference time is slightly larger than the other networks. This is because MUNet has a small number of filters used per convolution, but it performs many convolution operations due to the nature of the network. This issue will be touched again in Section 5.

4.4. Adaptability to the surrounding environment

MUNet is a CNN model designed for mobile device purpose. It is important that the CNN model fitted to a specific environment should maintain its performance even if the surrounding environment changes. Therefore, this section verifies whether MUNet conforms to the characteristics of the mobile device through the fine tuning experiment of MUNet with very small parameters. Unlike the experiment in the previous section, the reason for excluding ShuffleNet from the benchmarking group is that the fine tuning experiment was not performed in ShuffleNet.

First, we trained the MUNet by using a down-sampled ImageNet dataset, and performed a fine tuning process with the CIFAR dataset. Here, the reason for pre-training using down-sampled ImageNet is that as shown in [33], down-sampled ImageNet produces similar optimal hyper-parameters while exhibiting faster learning convergence speed than original ImageNet. In addition, since the image size of down-sampled ImageNet is the same as the image size to be fine-tuned, more accurate verification of fine tuning performance is possible. In both networks, the specification of Table 3 was used as is. In the case of MUNet, the number of convolution filters of small units was changed. As shown in Table 4, when fine tuning is performed, MUNet shows an additional performance

The number of 1-path small units	Error rate (%)	FLOPS [Mult. Add If]
1	36.82	27.00M [13.48 13.33 0.19]
2	34.72	40.32M [20.12 19.89 0.31]
3	34.40	53.69M [26.79 26.48 0.42]

Table 6: Performance comparison according to the number of 1-path small units for CIFAR-100.

improvement of about 2.4% based on CIFAR-100. On the other hand, MobileNet has a performance improvement of about 4.5%. It can be said that MobileNet, which has a larger number of parameters, could learn richer feature information than MUNet when pre-trained with a huge amount of datasets. Note that despite much smaller parameter size, MUNet shows successful fine tuning results. This means that even though the model size of the MUNet is very small, it can learn meaningful feature information when pre-trained with a large amount of datasets.

4.5. Performance evaluation according to MU-layer structure

The core concept of MUNet is to use multiple MUs with a small number of convolution filters at each MU-layer. In this experiment, the structural characteristics of such a MUNet are verified (see Table 5).

The first row of Table 5 shows the case where the MU-layer is designed with a single MU per layer structure. In this case, the number of convolution filters is increased enough to compensate for the reduced number of MUs per MU-layer for fair comparison. The result was compared with authentic MUNet (see the 2nd row in Table 5).

As can be seen from the experimental results in Table 5, the single MU per layer structure shows a minor accuracy improvement of 0.8% in comparison with the original MUNet, but it increases about 5 times both the parameter size and FLOPS. As a result, we can find that the architecture of MUNet is superior in performance.

4.6. Selection of the appropriate hyper-parameter

The core units of the MU architecture are 1-path small

unit and 2-path small unit. In Section 3.3, we emphasized that the 2-path small unit plays an essential role in MUNet performance, and also the reason for the existence of 1-path small unit is to increase the number of internal layers of MU model. So the number of 1-path small units can be regarded as a hyper-parameter.

In this section, we performed an experiment to find an optimal hyper-parameter for MUNet, i.e., the best number of 1-path small units. The result is shown in Table 6. As the number of 1-path small units increases, the FLOPS increases linearly, but the error rate decreases sharply and saturates. Therefore, this paper adopted two 1-path small units per MU, pursuing optimal trade-off performance. The number of MUs per MU-layer can also be regarded as a hyper parameter, although it has not been tested in this paper.

5. Discussion

As can be seen in Table 3, MUNet has significantly smaller parameter sizes and smaller FLOPS than the other networks. However, in terms of inference time, MUNet is somewhat inferior to the other networks unlike in terms of FLOPS. The reason is that MUNet has a small number of filters per convolution, but it must perform many convolution operations inherently. MUNet has a more complex structure than a simple group convolution-based network when forming the whole network structure under the same conditions as [21]. The more complicated the structure, the more burden can be placed on the inference time. Therefore, we can analyze that MUNet has less FLOPS but provides somewhat longer inference time than the other networks. In order to overcome this problem, shortening the computation time of the convolution as in [22, 32] can reduce the inference time of MUNet.

6. Conclusion

In this paper, we propose a new CNN architecture that escapes from the existing group convolution structure in order to design a deep neural network suitable for mobile devices. Since the conventional group convolution-based network uses a large number of convolution filters, the amount of computation is large and there may be a considerable redundant computation. In order to solve this computational cost problem fundamentally, this paper presented a MU-based CNN using only a small number of convolution filters where richer feature maps can be learned by using multiple MUs instead of a single GConv block. Experimental results show that the proposed MUNet consumes significantly less computation cost while maintaining comparable performance with existing networks including state-of-the-art CNNs.

Acknowledgements

This work was supported by the industrial Technology Innovation Program funded by the ministry of Trade, industry & Energy (MI, Korea) [10073154, Development of human-friendly human-robot interaction technologies using human internal emotional states].

References

- [1] R. Salakhutdinov, G. Hinton. Deep boltzmann machines. In *Artificial Intelligence and Statistics*, pages 448-455, 2009.
- [2] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th annual international conference on machine learning*, pages 609-616, 2009.
- [3] Y. L. Boureau, Y. L. Cun. Sparse feature learning for deep belief networks. In *Advances in neural information processing systems*, pages 1185-1192, 2008.
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097-1105, 2012.
- [5] K. Simonyan, A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [6] L. Sifre. Rigid-motion scattering for image classification. PhD thesis, Ph. D. thesis, 2014.
- [7] 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*, pages 770-778, 2016.
- [8] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. *arXiv preprint arXiv:1611.05431*, 2016.
- [9] F. Chollet. Xception: Deep Learning with Depthwise Separable Convolutions. *arXiv preprint arXiv:1610.02357*, 2016.
- [10] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten. Densely connected convolutional networks. *arXiv preprint arXiv:1608.06993*, 2016.
- [11] G. Larsson, M. Maire, and G. Shakhnarovich. Fractalnet: Ultra-deep neural networks without residuals. *arXiv preprint arXiv:1605.07648*, 2016.
- [12] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, pages 1929-1958, 2014.
- [13] X. Zhang, J. Zou, X. Ming, K. He, and J. Sun. Efficient and accurate approximations of nonlinear convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1984-1992, 2015.
- [14] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li. Learning structured sparsity in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 2074-2082, 2016.

- [15] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. *AAAI*. Vol. 4, 2017.
- [16] A. G. Howard, et al. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [17] X. Zhang, et al., ShuffleNet: An extremely efficient convolutional neural networks for mobile devices. *arXiv preprint arXiv:1707.01083*, 2017.
- [18] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger. Deep networks with stochastic depth. In *European Conference on Computer Vision*, pages 646-661, 2016.
- [19] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *European Conference on Computer Vision*, pages 630-645, 2016.
- [20] B. Zoph, and Q. V. Le. Neural Architecture Search with Reinforcement Learning. *arXiv preprint arXiv:1611.01578*, 2016.
- [21] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous systems, 2015. *Software available from tensorflow.org*, 1, 2015.
- [22] W. Xu, Z. Wang, X. You and C. Zhang. Efficient fast convolution architectures for convolutional neural network. In *ASIC (ASICON). IEEE 12th International Conference on*, pages 904-907, 2017.
- [23] A. Krizhevsky, G. Hinton. Learning multiple layers of features from tiny images. *Tech Report*, 2009.
- [24] I. J. Goodfellow, et al. "Multi-digit number recognition from street view imagery using deep convolutional neural networks." *arXiv preprint arXiv:1312.6082*, 2013.
- [25] A. Bastidas. Tiny ImageNet Image Classification. *url: https://tiny-imagenet.herokuapp.com/*.
- [26] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- [27] L. Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*. pp. 177-186, 2010.
- [28] I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139-1147, 2013.
- [29] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026-1034, 2015.
- [30] S. Ioffe, C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448-456, 2015.
- [31] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European Conference on Computer Vision*. pages 525-542, 2016.
- [32] M. Jaderberg, A. Vedaldi, and A. Zisserman. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014.
- [33] P. Chrabaszczyk, I. Loshchilov, and F. Hutter. A Downsampled Variant of ImageNet as an Alternative to the CIFAR datasets. *arXiv preprint arXiv:1707.08819*, 2017.
- [34] O. Russakovsky, J. Deng, H. Su, H. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211-252, 2015.
- [35] S. J. Hanson, and L. Y. Pratt. Comparing biases for minimal network construction with back-propagation. In *Advances in neural information processing systems*. Pages 177-185, 1989.