

Received January 20, 2020, accepted February 4, 2020, date of publication February 17, 2020, date of current version March 2, 2020.

Digital Object Identifier 10.1109/ACCESS.2020.2974242

Learning to Recognize Chest-Xray Images Faster and More Efficiently Based on Multi-Kernel Depthwise Convolution

MENGJIE HU^{ID}, HEZHENG LIN, ZIMENG FAN, WENJIE GAO, LU YANG^{ID}, CHUN LIU^{ID},
AND QING SONG^{ID}

Pattern Recognition and Intelligent Vision Laboratory, Beijing University of Posts and Telecommunications, Beijing 100876, China

Corresponding author: Qing Song (priv@bupt.edu.cn)

ABSTRACT The development of convolutional neural networks has promoted the progress of computer-aided diagnostic systems. Details in medical image, such as the texture and tissue structure, are crucial features for diagnosis. Therefore, large input images combined with deep convolution neural networks are adopted to boost the performance in recent research of chest X-ray diagnosis. Meanwhile, due to the variable sizes of thoracic diseases, many researchers have worked to introduce additional module to capture multi-scale feature of images in CNN. However, these efforts hardly consider the computational costs of large inputs and introduced additional modules. This paper aims to automatically diagnose diseases on chest X-rays images quickly and effectively. We propose the multi-kernel depthwise convolution(MD-Conv) which contains depthwise convolution kernels with different filter sizes in one depthwise convolution layer. MD-Conv has high calculation efficiency and few parameters. Because its ability to learn multi-scale feature based on the multi-size kernels, it is appropriate for medical images diagnosis tasks in which abnormalities varied in sizes. In addition, larger depthwise convolution kernels are adopted in MD-Conv to obtain a larger receptive field efficiently, which can ensure sufficient receptive field for high resolution inputs. MD-Conv can be easily applied in modern lightweight networks to replace the normal depthwise convolution layer. We conduct experiments on the Chest X-ray 14 Dataset, which is the largest available chest x-ray dataset, and obtain competitive results. We also evaluate the MD-Conv on the new released dataset for pediatric pneumonia diagnosis. We obtain a better performance of 98.3% AUC than original paper (96.8%) for recognize pneumonia versus normal. Meanwhile we compare the FLOPs and Params of different models to show their efficiency for chest X-rays recognition.

INDEX TERMS Chest x-ray recognition, lightweight networks, multi-kernels depthwise convolution.

I. INTRODUCTION

The development of convolutional neural networks(CNN) has made a dramatic breakthrough in a series of computer vision tasks, which has also promoted the computer-aided diagnosis system. Medical images have grown exponentially in hospital, and disease screening is a time consuming task for radiologists. The computer-aided diagnosis system can help to do preliminary screening and reduce the burden of radiologists.

Chest X-ray is one of the most accessible radiology examinations in the world. Research on chest X-ray includes

The associate editor coordinating the review of this manuscript and approving it for publication was Yongqiang Zhao^{ID}.

thoracic disease identification and localization [34], lung regional segmentation [19] and diseases report generation [35]. Among all these studies, the details in medical images, such as textures and structure of lung tissue, are crucial features for diagnosis. That is the reason that the Chest X-ray 14 dataset keeps 1024×1024 bitmap images [34] to preserve details, which exceeding the 512×512 images in OpenI dataset [3]. Similarly, most global image-based CNN methods adopt large images as inputs, 512×512 or even 1024×1024 in [2], [18], [37]. And some local image-based CNN methods use 224×224 inputs [7], because the 224×224 inputs are large enough to take local details for local images. On the other hand, the pathologies in chest X-ray images are highly varied in their shapes and sizes. We conduct

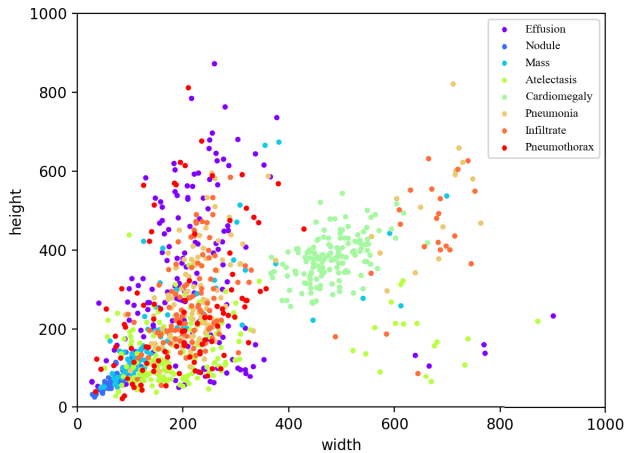


FIGURE 1. Statistics of the different sizes of eight common thoracic pathologies in Chest X-ray 14 dataset.

statistical analysis on the sizes of eight common thoracic pathologies based on the 984 boxes provided by Chest X-ray 14 dataset. As shown in FIGURE 1, the size of eight common thoracic pathologies varies in a wide range, and even different instances of one thoracic pathology, such as Infiltrate, have different sizes. This requires learning multi-scale convolutional features in CNN. Some recent work attempts to solve this problem by fusing the information from multiple resolutions [39], [41].

Due to the enlarged input as mentioned above, a deeper network is always adopted to ensure the network receptive field is large enough. Many works choose ResNet-50 [13] and DenseNet-121 [5] to extract CNN features [18], [22]. Although this improves the performance, the large inputs combined with deep networks brings quite huge computational costs and parameters, increasing time for network training and optimization. For example, doubled input size will lead to four times training time. And the multiple resolutions feature fusion also cost computation time and store space. Thus this is not conducive to future deployments on mobile and embedded systems.

In this paper, we focus on increasing network receptive field and efficiently learning multi-scale feature. We firstly leverage the lightweight networks. There are many excellent lightweight networks, such as MobileNet [10], ShuffleNet [40], MobileNet-v2 [23], ShuffleNet-v2 [21] and MobileNet-v3 [1]. These networks take the network structures of VGG [28] and ResNet [14] for reference, and replace normal convolution with depthwise convolution to reduce parameters and FLOPs while maintaining the accuracy. By fully considering the balance between computation and accuracy during the design process, these networks have achieved good performance on ImageNet dataset. In addition, we propose the multi-kernel depthwise convolution (MD-Conv), which can capture the multi-scale feature in one convolution layer without introducing extra layers or blocks. Meanwhile the larger depthwise convolution kernels, 5×5 kernels, are adopted in MD-Conv to efficiently obtain a larger

receptive field. MD-Conv is appropriate for medical images with abnormalities in various sizes.

We replace the normal depthwise convolution with the proposed MD-Conv in popular lightweight networks, and evaluate the modified model on two public datasets: the Chest X-ray14, which is the largest available chest x-ray dataset, and the Chest X-ray2017, which is a recently released chest x-ray dataset for pediatric pneumonia diagnosis. We achieve state-of-the-art results on both datasets. The modified MD-Conv can successfully identify the chest X-rays quickly and effectively.

The contributions of this paper are as follows.

(1) Compared with modern methods adding complex network and additional block to improve performance, we adopt a lightweight network to quickly recognize chest X-rays which requires small model parameters, and is suitable to employ on embedded systems.

(2) The problem of multi-scale feature learning is studied. Based on the various sizes of thorax diseases and enlarged inputs, we propose MD-Conv, which is conducive to learning the multi-scale feature of different thorax diseases and improving network performance.

(3) The modified MobileNet-v2 with MD-Conv achieves competitive results on the Chest X-ray 2017 dataset and the Chest X-ray 14 dataset.

II. RELATED WORK

A. DEEP LEARNING FOR CHEST X-RAY DIAGNOSIS

Wang *et al.* [34] release the largest chest X-ray dataset and utilize different models to recognize and locate thorax diseases. Since then, a series of studies have been explored based on the large dataset, such as image classification, weakly supervised localization, medical report generation for medical image. In [34], four classic CNN models, AlexNet [16], GoogleNet [31], VGG16 [28], ResNet-50 [13] are compared in the proposed DCNN framework for disease localization. In [22], CheXNet is proposed and demonstrate that DenseNet [5] performs much better on chest x-ray images recognition. And then cascade ConvNet [17], global local fusion method [7] and multi-scale feature are proposed to improve recognition performance, and all these works use ResNet and DenseNet as the basic network for feature extraction. Among all these works, ResNet50 and DenseNet121 are most widely used models. Reference [17] also uses DenseNet161, and [38] reduces the Conv-Block number within a DenseBlock to four to get a light model.

In [15], authors try to utilize CNN to diagnose pediatric pneumonia on chest X-ray images. Based on transfer learning algorithm, it reaches a AUC of 96.8% for recognize pneumonia from normal on chest X-rays dataset. It adopts efficient Inception-v3 [32] as the basic network.

B. MULTI-SCALE METHODS IN COMPUTER VISION TASKS

Though CNN is robust to do recognition on images with objects of different sizes, how to obtain a multi-scale feature

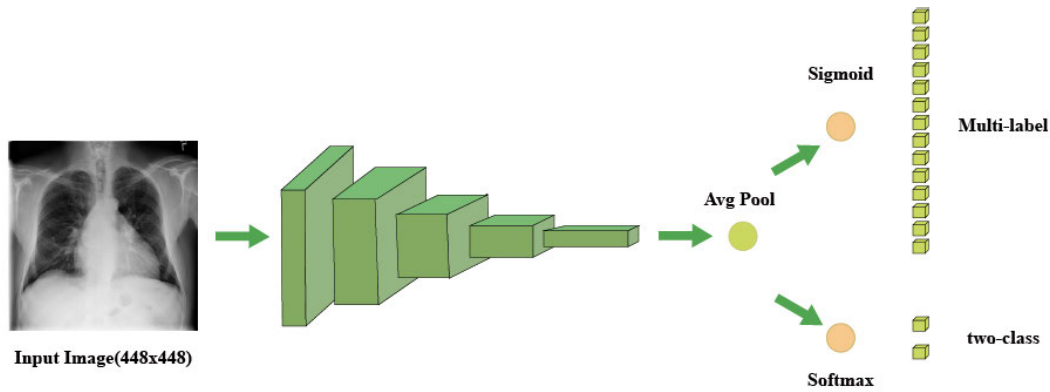


FIGURE 2. The framework of Chest-Xray image recognition.

representation is an important issue in many computer vision tasks.

Traditional approaches use image pyramids to get a more accurate results, for example the multi-scale test in [13]. While most state-of-the-art methods utilize features from different layers to obtain inherent multi-scale in network. FPN [33] uses upsample and latency to generate feature pyramid, and SSD [36] reuses the multi-scale feature maps from different layers. FCN [20] is one of the earliest methods to fuse the multi-scale representations in semantic segmentation. Recently, HRNet [30] for human pose estimation performs repeated multi-scale fusion to achieve state-of-the-art results.

Therefore, in the field of medical deep learning, many multi-scale CNNs have also been proposed to learn multi-scale feature for abnormalities. In [24], [39], additional blocks are added to fuse the multi-scale feature from different layers.

C. EFFICIENT NETWORK DESIGN

With the development of convolution neural networks, researchers have become interested in efficient model design. GoogleNet [31] is one of the earliest networks which is designed for computational efficiency. And since depthwise convolution is proposed [27], depthwise convolution is utilized in modern lightweight network to replace the normal convolution for its efficiency and effectiveness. There are five popular lightweight networks, MobileNet [10], ShuffleNet [40], MobileNet-v2 [23], ShuffleNet-v2 [21] and MobileNet-v3 [1]. Among these five lightweight networks, MobileNet modify VGG structure with depthwise convolution, while others adopt ResNet-like structure. The lightweight network CondenseNet [11] is created based on the extensive and explicit feature reuses structure of DenseNet.

The recent work [29] which is named as HetConv is similar to our work. However, their work focuses on efficient convolution computation, and they design the heterogeneous convolutional filter by 1×1 conv and 3×3 conv. While our work focus on the multi-scale feature learning ability and the large receptive field of the large convolution kernels.

Therefore, 3×3 depthwise conv and 5×5 depthwise conv are adopted in our proposed MD-Conv, and 1×1 kernel is not recommended in MD-Conv. The separated channels of depthwise convolution is suitable for implementation of multi-kernel depthwise convolution.

III. METHOD

We aim to build a CNN to quickly and efficiently recognize chest X-ray images. Compared with the widely used networks, ResNet50 and DenseNet121, we adopt the lightweight network MobileNet-v2. The framework of chest X-ray image recognition is shown in FIGURE. 2. In order to learn the multi-scale feature, we propose the MD-Conv to replace the depthwise convolution in MobileNet-v2.

The MD-Conv compared together with normal convolution and depthwise convolution are illustrated in FIGURE. 3. The MobileNet-v2 block with MD-Conv is shown in FIGURE. 3(e).

A. THE MULTI-KERNEL DEPTHWISE CONVOLUTION

For a standard convolution layer, it takes an input feature map $\mathbf{I} \in \mathbb{R}^{H \times W \times M}$, and outputs a feature map $\mathbf{O} \in \mathbb{R}^{H \times W \times N}$, here we assume that the feature map spatial width(W) and height(H) are constant, and M is the number of input channels, N is the number of output channels. For depthwise convolution layer(DW), each of these filters only connects to one input channel. Then an additional layer pointwise convolution(PW) is added after DW to calculate a linear combination of the output of DW. The feature map operations of standard convolution layer and DW + PW are shown in FIGURE. 3(a)(b) respectively. For a standard convolution layer with kernel size of $K \times K$, the computational cost and parameters are computed as:

$$Cost_{norm} = K \times K \times M \times N \times H \times W. \quad (1)$$

$$Params_{norm} = K \times K \times M \times N. \quad (2)$$

For a depthwise convolution of $K \times K$ paired with 1×1 pointwise convolution, the computational costs and parameters are

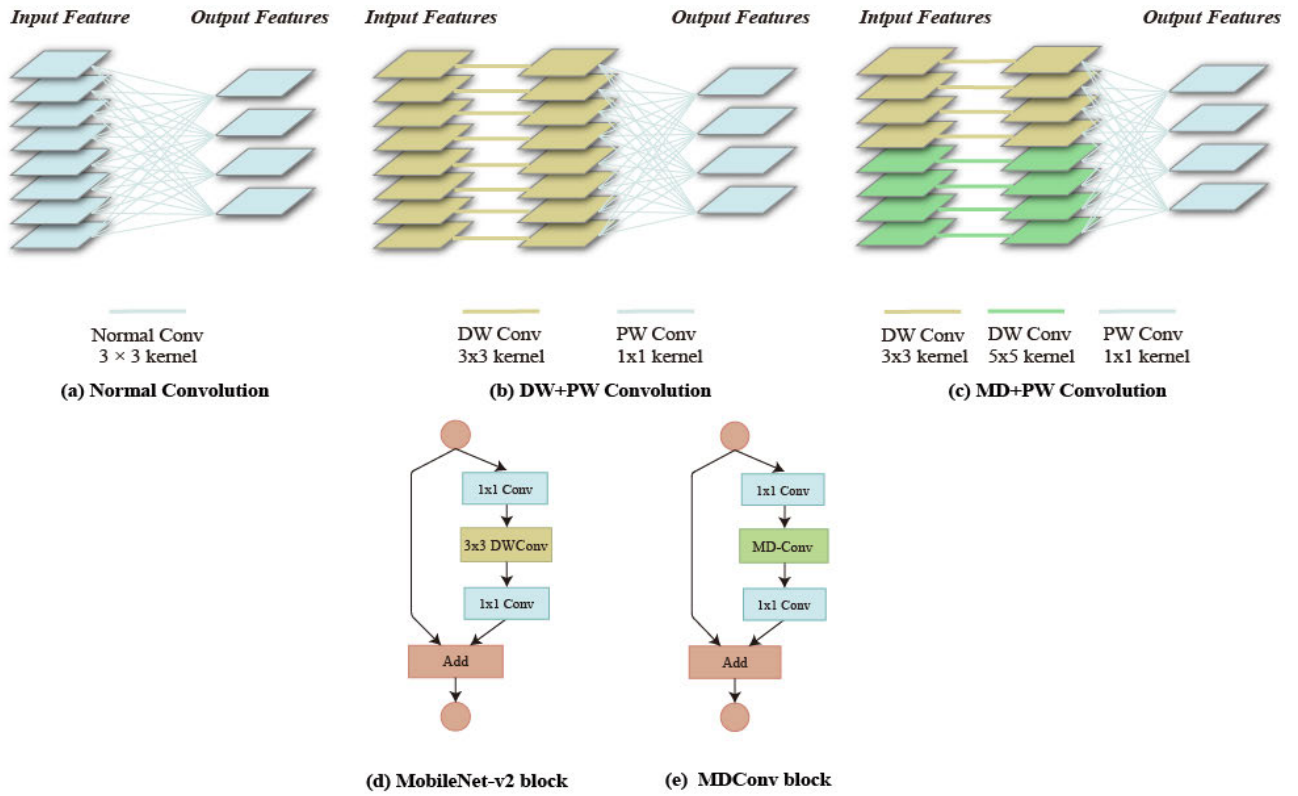


FIGURE 3. (a)(b)(c) show comparison of the different feature operations of normal convolution, depthwise convolution and MD-Conv. (d)(e) show the MobileNet-v2 block with depthwise convolution and MD-Conv.

computed as:

$$Cost_{DW} = K \times K \times M \times H \times W + M \times N \times H \times W. \quad (3)$$

$$Params_{DW} = K \times K \times M + M \times N. \quad (4)$$

Because the output channel N is much larger than K^2 , the computational cost and parameters of $DW + PW$ is $\frac{1}{N}$ times the normal convolution.

We introduce the multi-kernel depthwise convolution(MD-Conv), which contains both 3×3 DWConv and 5×5 DWConv in a multi-kernel depthwise convolution layer. The feature map computation of MD-Conv + PW is shown in FIGURE. 3(c). For each depthwise convolution kernel in the MD-Conv layer, it corresponds to only one input channel, thus we can easily implement the MD-Conv by channel split operation. Here, we consider an alternative slicing of the input feature $\mathbf{I} = [\mathbf{i}^{3 \times 3}, \mathbf{i}^{5 \times 5}]$, where $\mathbf{i}^{3 \times 3}$ corresponding to input feature for 3×3 DWConv and $\mathbf{i}^{5 \times 5}$ for 5×5 DWConv. As shown in FIGURE. 3(c), the yellow feature maps correspond to $\mathbf{i}^{3 \times 3}$, and the green feature maps correspond to $\mathbf{i}^{5 \times 5}$. Therefore, the output of DWConv is

$$\mathbf{O}_{MD-Conv} = [W_{3 \times 3} \mathbf{i}^{3 \times 3}, W_{5 \times 5} \mathbf{i}^{5 \times 5}]. \quad (5)$$

And then the PW is followed to fuse the separable channels of MD-Conv.

For the MD-Conv with an input feature map of $H \times W \times M$ and same size output feature map, the

computation cost and parameters are:

$$Cost_{MD} = (9 \times \mathbf{i}^{3 \times 3} + 25 \times \mathbf{i}^{5 \times 5}) \times H \times W. \quad (6)$$

$$Params_{MD} = 9 \times \mathbf{i}^{3 \times 3} + 25 \times \mathbf{i}^{5 \times 5}. \quad (7)$$

For the 3×3 DWConv with the same inputs and outputs, the number of its parameters is $9 \times (\mathbf{i}^{3 \times 3} + \mathbf{i}^{5 \times 5}) \times H \times W$. MD-Conv only adds tiny amounts of parameters compared to 3×3 DWConv.

1) THE EFFICIENT WAY TO IMPROVE RECEPTIVE FIELD

Receptive field plays an important role in CNN. For visual tasks with low resolution inputs, the receptive field of standard 3×3 filter is sufficient. While for medical images with high resolution, increased receptive field is needed to retain more details. The network receptive field is computed as:

$$\mathbf{RF}_n = \mathbf{RF}_{n-1} + (k_n - 1) \times s_n. \quad (8)$$

where n represents the network layer, k_n is the kernel size of layer n , s_n is the stride of layer n .

To increase the receptive field, a simple method is to increase the depth of the network, such as densenet121. The receiving area of the two stacked 3×3 convolutions is 5, which is the same as the 5×5 convolutions, and has less computation and parameters. In [32], two stacked 3×3 convolution are used to replace 5×5 convolution for efficiency.

For lightweight networks, an efficient way to increase the receptive field is to adopt a larger kernel depthwise convolution. A 5×5 depthwise convolution has the same receptive field as 5×5 convolution, and it requires much less extra computational costs and parameters. Two stacked 3×3 depthwise convolution is not adopted due to no channel cross talk [40].

Additionally, two stacked 3×3 convolutions have a computation cost of $18 \times M \times N \times H \times W$, while the 5×5 depthwise convolution has a computation cost of $25 \times M \times H \times W$. The 5×5 depthwise convolution has much smaller computation cost while maintaining the same receptive field.

Therefore, the part of 5×5 depthwise convolution can more efficiently obtain a large receptive field in MD-Conv.

2) COMPARED WITH THE MULTI-SCALE FEATURE FUSION METHOD

The image-level classification always requires coarse-scale features with high semantic information and context, while detection and segmentation need fine-scale features to capture detailed appearance information. Therefore, in [39], feature fusion at different resolutions is used. In [41], a multi-resolution CNN is adopted to recognize nodules of different sizes. These methods all require additional convolution blocks and downsample or upsample to integrate information from multiple scales.

The MD-Conv is proposed to replace the normal depthwise convolution. In one MD-Conv block as shown in FIGURE. 3(e), the MD-Conv can extract the multi-scale feature based on the multi-scale depthwise convolution kernels, and then the 1×1 convolution is followed to fuse information of different scales. This introduces no addition convolution layers or operations, and adds only slight parameters of 5×5 depthwise convolution. Moreover, the MD-Conv can be easily integrated into modern networks.

B. HOW AND WHERE TO USE MD-CONV

In this subsection, we will discuss how and where to use MD-Conv, which is the ratio of different kernels ($i^{3 \times 3}$, $i^{5 \times 5}$) in one MD-Conv layer, and where the MD-Conv should be used in a network.

1) BEST $i^{3 \times 3}/i^{5 \times 5}$ VALUE

As mentioned above, the MD-Conv consists of 3×3 depthwise convolution and 5×5 depthwise convolution. The parameters $i^{3 \times 3}$ and $i^{5 \times 5}$ control the number of different types of kernels in one MD-Conv layer. To find the best ratio of $i^{3 \times 3}/i^{5 \times 5}$, we do experiments on different ratios of $i^{3 \times 3}/i^{5 \times 5}$ in MD-Conv. In these experiments, we replace all the normal depthwise convolution layer in MobileNet-v2 with MD-Conv layer and perform the experiments on the Chest-Xray 14 Dataset. The results are shown in FIGURE. 4. According to it, we choose $i^{3 \times 3}/i^{5 \times 5}$ ratio of 1:1 to achieve the best cost and accuracy trade-off.

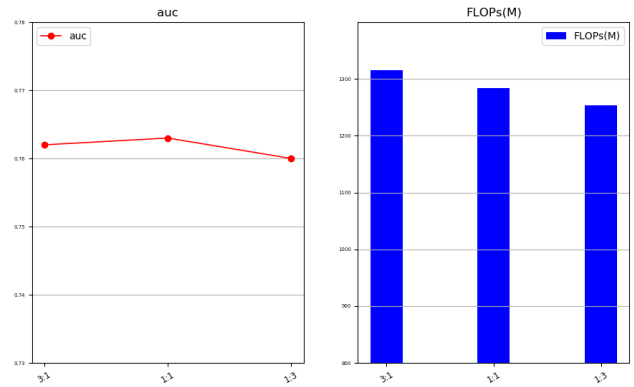


FIGURE 4. Comparison results of different ratios of $i^{3 \times 3}/i^{5 \times 5}$. The left one shows mean AUC score of 14 diseases, the right one shows the computation cost.

TABLE 1. The network architecture of modified MobileNet-v2 with MD-Conv. MD-Conv denotes whether MD-Conv is used in bottleneck, n denotes the repeated times of bottleneck in stage, and s denotes stride for the first layer of each stage.

Layer	Operator	n	MD-conv	s
conv	conv2d	-	-	2
stage0	bottleneck	1	-	1
stage1	bottleneck	2	-	2
stage2	bottleneck	3	-	2
stage3	bottleneck	4	-	2
stage4	bottleneck	3	-	1
stage5	bottleneck	3	✓	2
stage6	bottleneck	1	-	1
conv	conv2d	1	-	1
pool	avgpool	-	-	-
conv	conv2d	-	-	-

2) WHERE TO USES MD-CONV

MD-Conv is proposed to efficiently obtain the multi-scale feature and improve receptive field. We find that it is not necessary to use MD-Conv in all layers. Placing MD-Conv in right location can improve performance while saving FLOPs and Params. We conduct experiments to explore which layer the MD-Conv should be used. As shown in TABLE 4, we achieve the best AUC score replacing all the normal depthwise convolution layer in layer2 with MD-Conv. The final network architecture of modified MobileNet-v2 is shown in TABLE 1.

IV. EXPERIMENT

A. CHEST X-RAY DATASETS

Chest X-ray is the most common and efficient technique for screening and diagnosis of lung-related diseases, such as pneumonia, cardiomegaly, lung node. Several chest X-rays datasets have released for study. Early datasets, such as [12], [26], contain only hundreds of chest X-ray images, which are too few for deep learning. OpenI [3] is a publicly available dataset collected by Indiana University, it contains

TABLE 3. Mean auc of five popular lightweight networks on Chest X-ray14.

Network	mean_auc	FLOPs(M)
MobileNet	0.740	2291
MobileNet-v2	0.772	1223
MobileNet-v3	0.766	865
ShuffleNet	0.769	521
ShuffleNet-v2	0.758	2369

TABLE 4. Mean auc vs. FLOPs for modified MobileNet-v2 with MD-Conv in different location.

Layer	mean_auc	FLOPs(M)
all	0.777	1278
layer1	0.779	1235
layer2	0.777	1247
layer3	0.775	1236
layer5	0.782	1228
layer1,2	0.776	1260
layer1,3	0.778	1249
layer1,5	0.781	1241
layer2,3	0.779	1260
layer2,5	0.776	1252
layer3,5	0.781	1241
layer1,2,3	0.770	1273
layer1,2,5	0.776	1265
layer1,3,5	0.775	1254
layer2,3,5	0.774	1265
basic	0.772	1223

because the depthwise convolution layer in MobileNet-v2 has increased channels which are enough for each depthwise convolution group with different kernel size in MD-Conv that ShuffleNet-v1 dose not have. Meanwhile, MobileNet-v2 has the best accuracy in chest X-rays recognition as shown in TABLE 3.

3) ABLATION STUDY ON WHERE AND WHETHER TO USE MD-CONV

To explore where and whether we should use MD-Conv in MobileNet-v2, we conduct ablation experiments. In MobileNet-v2, the stride of layer1, layer2, layer3 and layer5 is 2, and using MD-Conv in these layers can make receptive field increase in multiples. Therefore, we use MD-Conv only in these layers, and the results are shown in TABLE 4. We can find that using MD-Conv in all layers achieve the low AUC score and its FLOPs are the largest. Using MD-Conv in layer5 achieves the highest AUC and the fewest FLOPs. Therefore, we use MD-Conv only in layer5, as shown in TABLE 1. We guess MobileNet-v2 with MD-Conv in layer5 can obtain a suitable receptive field for input size of 448×448 . Using MD-Conv in all layers is not necessary.

TABLE 5. The mean auc score and FLOPs(M) with or without MD-Conv in MobileNet-v2 and MobileNet-v3.

Network	mean_auc	FLOPs(M)
MobileNet-v2	0.772	1223
MobileNet-v2-mdconv	0.777	1278
MobileNet-v3	0.766	865
MobileNet-v3-mdconv	0.785	885

In order to explore the effect of MD-Conv, we also conduct additional ablation experiments using MobileNet-v3, MD-Conv is used in all layers, the results are shown in TABLE 5. As can be seen, MD-Conv improves the performance of both networks.

4) CLASSIFICATION PERFORMANCE COMPARED WITH THE STATE-OF-THE-ART APPROACH

We compare our results with Resnet50 [34], DenseNet + LSTM [38], DenseNet121 [8] based on the official dataset split. The results are shown in TABLE 6. As we can see, our method performs better than [34] which only adopts ResNet-50. While for [38] which introduces LSTM in network, and [8] which adopts DenseNet121 and uses additional data PLCO Dataset [6]. Although our results are not the best, we provide a baseline for latter research studying lightweight networks in chest X-ray images recognition task and the probability to identifying diseases in embedded and mobile devices. Meanwhile, we only provide a basic network for chest X-ray recognition, which can achieve competitive results with less computation costs and parameters and can be easily used in any other networks with depthwise convolution.

The comparison of computational costs of ResNet50, DenseNet121, and our modified MobileNet-v2 is shown in TABLE 7. As the table shows, our modified model is much lighter on FLOPs and Params. Thus we think there is more space to make improvements based on our basic network. And with such small computation cost and parameters, we still outperform the results of [34].

D. CHEST X-RAY 2017 DATASET

To verify the generalization ability of our model, we also do experiments on the ChestX-ray2017 Dataset released by [15]. We use MobileNet-v2 with MD-Conv to recognize pneumonia versus normal, bacterial versus viral pneumonia. The results are shown in Table 8.

For Chest-Xray recognition of pneumonia versus normal, we achieve an accuracy of 93.4%, outperforms [15] by 0.6%. And the sensitivity outperforms [15] by 4.2%. Though the specificity is 2.3% lower, we achieve an AUC of 98.3%, while is higher than 96.8% in [15]. The AUC curve is shown in FIGHRE 6. Besides, our modified MobileNet-v2 with MD-Conv is lighter than Inception-v3 adopted in [15].

TABLE 6. Comparison results of different methods. we list fourteen abnormalities and their AUCs, we also list the method adopted in every paper.

Abnormality	Wang et al. [34]	Yao et al. [38]	Gündel et al. [8]	Ours
Atelectasis	0.716	0.772	0.767	0.767
Cardiomegaly	0.807	0.904	0.883	0.896
Effusion	0.784	0.859	0.828	0.828
Infiltration	0.609	0.695	0.709	0.696
Mass	0.706	0.792	0.821	0.752
Nodule	0.671	0.717	0.758	0.733
Pneumonia	0.633	0.713	0.731	0.725
Pneumothorax	0.806	0.841	0.846	0.815
Consolidation	0.708	0.788	0.745	0.742
Edema	0.835	0.882	0.835	0.857
Emphysema	0.815	0.829	0.895	0.829
Fibrosis	0.769	0.767	0.818	0.816
PT	0.708	0.765	0.761	0.749
Hernia	0.767	0.914	0.896	0.735
mean auc	0.738	0.803	0.804	0.782
method	Resnet50	DenseNet+LSTM	DenseNet121	modified MobileNet-v2

TABLE 7. The computation cost and parameters of some basic networks.

Method	FLOPs(M)	Params(M)
ResNet50	16395	23.5
DenseNet121	11405	7.0
ours	1241	2.3

TABLE 8. Comparison results between ours and [15]. pne represents pneumonia and bac represents bacterial.

Class	Method	accuracy(%)	sensitivity(%)	specificity(%)
pne vs. norm	[15]	92.8	93.2	90.1
	ours	93.4	97.4	86.8
bac vs. viral	[15]	90.7	88.6	90.9
	ours	91.0	88.5	92.6

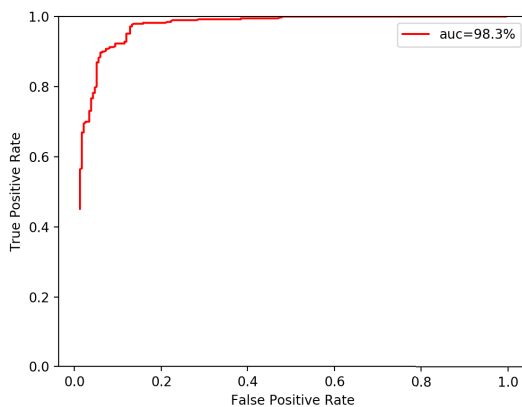


FIGURE 6. The ROC curve of our method for detecting pneumonia versus normal. The area under the ROC curve(AUC) is 98.3%.

E. QUALITATIVE ANALYSIS

We obtain the discriminative regions for each disease by Grad-CAM to show the visual explanation of how CNN recognize Chest-Xray images. As shown in FIGURE 7, different

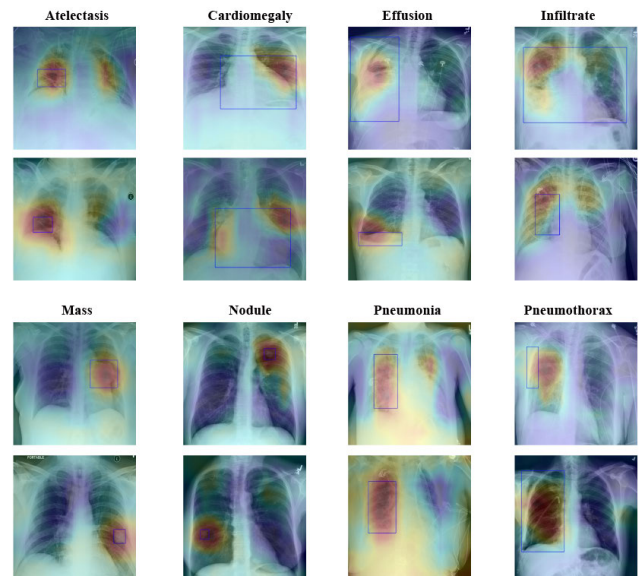


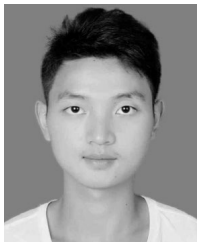
FIGURE 7. Examples of some network attention visualization on test image in Chest X-ray14 dataset. The attention map is generated by Grad-CAM [25], the red region represents the place where disease appears probably. The blue boxes is the ground-truth bounding box provided in the dataset.

diseases have different sizes. The Mass and Nodule are small and have equal length and width, concluded in size distribution in FIGURE 1 and the visualization results in FIGURE 7, the location heatmaps concentrate on an approximate circular area. And for Pneumonia, the location heatmaps focus on most chest region, and the area with most attention is in the blue ground-truth box. For Cardiomegaly, the ground-truth boxes include all the heart region, while the heatmaps pay attention to the enlarged heart margin which can also recognize Cardiomegaly effectively. Since that, the model can locate the disease region no matter the shapes and sizes of the diseases.

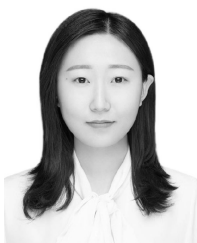
- [35] X. Wang, Y. Peng, L. Lu, Z. Lu, and R. M. Summers, "TieNet: Text-image embedding network for common thorax disease classification and reporting in chest X-rays," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 9049–9058.
- [36] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C. Fu, and A. Berg, "SSD: Single shot multibox detector," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 21–37.
- [37] L. Yao, E. Poblenz, D. Dagunts, B. Covington, D. Bernard, and K. Lyman, "Learning to diagnose from scratch by exploiting dependencies among labels," 2017, *arXiv:1710.10501*. [Online]. Available: <https://arxiv.org/abs/1710.10501>
- [38] L. Yao, E. Poblenz, D. Dagunts, B. Covington, D. Bernard and K. Lyman, "Learning to diagnose from scratch by exploiting dependencies among labels," 2017, *arXiv:1710.10501*. [Online]. Available: <https://arxiv.org/abs/1710.10501>
- [39] L. Yao, J. Prosky, E. Problez, B. Covington, and K. Lyman, "Weakly supervised medical diagnosis and localization from multiple resolutions," 2018, *arXiv:1803.07703*. [Online]. Available: <https://arxiv.org/abs/1803.07703>
- [40] X. Zhang, X. Zhou, M. Lin, and J. Sun, "ShuffleNet: An extremely efficient convolutional neural network for mobile devices," 2017, *arXiv:1707.01083*. [Online]. Available: <http://arxiv.org/abs/1707.01083>
- [41] W. Zuo, F. Zhou, Z. Li, and L. Wang, "Multi-resolution CNN and knowledge transfer for candidate classification in lung nodule detection," *IEEE Access*, vol. 7, pp. 32510–32521, 2019.



MENGJIE HU received the Ph.D. degree from Beihang University, Beijing, China, in 2017. She is currently a Lecturer with the Beijing University of Posts and Telecommunications. Her current research interests are primarily in computer vision and machine learning, especially object detection, visual tracking, and visual geometry.



HEZHENG LIN received the B.Sc. degree from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China, in 2019, where he is currently pursuing the master's degree with the Pattern Recognition and Intelligent Vision Laboratory (PRIV). His research interests include lightweight networks, video classification, and object detection.



ZIMENG FAN received the B.Sc. degree from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China, in 2017, where she is currently pursuing the master's degree with the Pattern Recognition and Intelligent Vision Laboratory (PRIV). Her research interests cover lightweight networks, image classification, and human pose estimation.



WENJIE GAO was born in Henan, China, in 1994. He received the B.Sc. degree from the Nanyang Institute of Technology. He is currently pursuing the M.Sc. degree with the Beijing University of Posts and Telecommunications. His research interests include image classification, 3-D image detection, and human pose estimation in computer vision.



LU YANG received the bachelor's degree from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China, in 2012, where he is currently pursuing the Ph.D. degree with the Automated Institute. He has participated in the MSCOCO 2018 Challenge, one of the top competitions in the field of computer vision, and won the DensePose champion. He has been involved in research work with the Pattern Recognition and Intelligent Vision Laboratory, since 2012. His research interests include the fields of artificial intelligence, computer vision, and machine learning.



CHUN LIU received the Ph.D. degree from the University of Kassel, Germany, in 2014. She is currently a Lecturer with the Beijing University of Posts and Telecommunications (BUPT). Her main research interests are primarily in intelligent computation and optimization, especially evolutionary algorithms in solving optimization problems, e.g., planning and scheduling.



QING SONG received the Ph.D. degree from Tianjin University, Tianjin, China, in 2006. She is currently a Scientific Researcher with the Beijing University of Posts and Telecommunications (BUPT), where she is also involved in computer vision technology study. She is also the Founder of the Pattern Recognition and Intelligent Vision Laboratory (PRIV) and led the PRIV Team to the Championship of COCO2018-DensePose Challenge. She is also in charge of many national, provincial and ministerial projects, and enterprise cooperation projects. She has published more than 70 academic articles in international journals and conferences.

...