# Discriminative Kernel Convolution Network for Multi-Label Ophthalmic Disease Detection on Imbalanced Fundus Image Dataset

Amit Bhati, Neha Gour, Pritee Khanna and Aparajita Ojha

*Abstract*—It is feasible to recognize the presence and seriousness of eye disease by investigating the progressions in retinal biological structure. Fundus examination is a diagnostic procedure to examine the biological structure and anomaly of the eye. Ophthalmic diseases like glaucoma, diabetic retinopathy, and cataract are the main reason for visual impairment around the world. Ocular Disease Intelligent Recognition (ODIR-5K) is a benchmark structured fundus image dataset utilized by researchers for multi-label multi-disease classification of fundus images. This work presents a discriminative kernel convolution network (DKCNet), which explores discriminative region-wise features without adding extra computational cost. DKCNet is composed of an attention block followed by a squeeze and excitation (SE) block. The attention block takes features from the backbone network and generates discriminative feature attention maps. The SE block takes the discriminative feature maps and improves channel interdependencies. Better performance of DKCNet is observed with Inception-Resnet backbone network for multi-label classification of ODIR-5K fundus images with 96.08 AUC, 94.28 F1-score and 0.81 kappa score. The proposed method splits the common target label for an eye pair based on the diagnostic keyword. Based on these labels oversampling and undersampling is done to resolve class imbalance. To check the biasness of proposed model towards training data, the model trained on ODIR dataset is tested on three publicly available benchmark datasets. It is found to give good performance on completely unseen fundus images also.

*Index Terms*—Multi-Label Classification, Channel Shuffle, Discriminative Kernel Convolution, Fundus Image, ODIR-5K.

Fig. 1. Fundus image of an eye with different kinds of abnormalities [4].

## I. INTRODUCTION

Ophthalmic diseases are leading cause of blindness worldwide. A report published by the world health organization (WHO) in 2021 says that around 2.2 billion people are visually impaired, and almost half of these are preventable based on timely detection and treatment [1]. The human retina is a light-sensitive layer of tissues in the rear end of the eye. The incident light is converted into neural signals through the receptors on retina and handled by the brain's visual cortex to generate a picture. The retina gets affected by different abnormalities which influences vision [2]. Fundus, fluorescein angiography, and optical coherence tomography (OCT) are standard modalities used by experts to investigate ophthalmic diseases [3]. Fundus imaging is the primary image modality utilized for clinical examination of ophthalmic diseases.

The presence of ophthalmic diseases can be recognized by observing abnormalities close to different retinal areas like optic nerve, veins, macula, optic plate, etc. Early identification of retinal abnormalities depicted in Fig 1 is crucial, yet difficult, as a few signs appear in the beginning stage. Deep Neural Networks (DNN) are effectively utilized for retinal vessel segmentation, lesion detection [5], [6], and glaucoma or diabetic retinopathy stage classification [7], [8]. Traditional single-label classification, also known as multi-class classification, includes a single class label for each instance.

However, ophthalmic disease classification is a complex problem as multiple labels may be associated with a single instance. Among various publicly available fundus image datasets, only Ocular Disease Intelligent Recognition (ODIR-5K) dataset [9] presents the real-life challenge of multi-class multi-label ophthalmic disease detection. However, ODIR-5K is an imbalanced dataset. The biasing towards the majority class in an imbalanced dataset affects models' training and classification accuracy.

This work presents a DNN based framework for multi-label ophthalmic diseases classification of the fundus image. The proposed architecture improves multi-label classification accuracy by handling the issue of class imbalance in the fundus image dataset. The proposed dilated convolution based attention network named Discriminative Kernel Convolution Network (DKCNet) can simultaneously detect multiple lesion parts related to ophthalmic diseases appearing in a fundus image. Different backbone models are also used to evaluate their efficiency for multi-label classification. The proposed DKCNet achieves better performance for ophthalmic disease classification as compared to the methods proposed in the literature.

The work is organized as follows. Section II connects retinal abnormalities in fundus images with ophthalmic diseases. State-of-the-art techniques for fundus image classification are discussed in Section III. Section IV describes the proposed transfer learning based DKCNet model for fundus image classification. The dataset used for experimental evaluation is also discussed here. The results are discussed in Section V and the work is concluded in Section VI.

## II. OPHTHALMIC DISEASES IDENTIFIED THROUGH RETINAL ABNORMALITIES IN FUNDUS IMAGES

Fig. 2 shows fundus images belonging to different disease classes from the ODIR-5K dataset. Glaucoma is an eye condition that affects the optic nerve, whose strength is crucial for good vision. This harm is typically caused by an unusually high pressure in the eye [10]. Glaucoma can be distinguished by noticing changes in the proportion of the optic disc cup and neuro-retinal edge surface region known as the cup-to-disc ratio. Diabetic retinopathy is a diabetes intricacy that influences the tissues inside the retina. Diabetic retinopathy may not show any adverse indications, and patients may have minor issues related to clear vision. However, it may lead to visual impairment if not treated at an earlier stage. The macula is liable for clear central vision, and the distortion in vision starts if liquids collect in it. Age-related Macular Degeneration (AMD) can be distinguished by noticing the growth of fresh blood vessels or the presence of dead

Amit Bhati, Neha Gour, Pritee Khanna, and Aparajita Ojha are with the Computer Science and Engineering Discipline, PDPM Indian Institute of Information Technology, Design and Manufacturing, Jabalpur, India (e-mail: pkhanna@iiitdmj.ac.in).
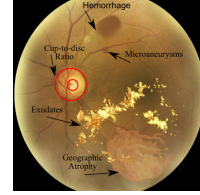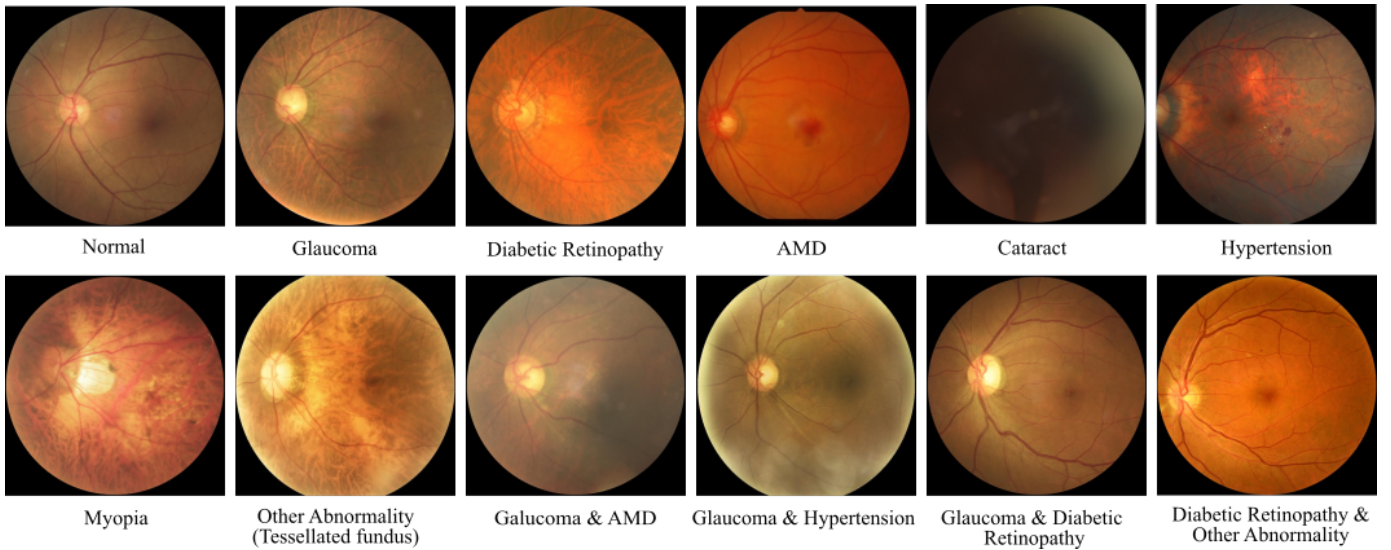
Fig. 2. Fundus images of different disease classes from ODIR-5K dataset [9].

retinal cells known as neovascularization and geographic atrophy, respectively [11]. The cause of cataracts includes the development of patches that make vision difficult [12]. The presence of cataracts can be observed as optic disc, fovea, and other parts of the eye become hazy. Hypertension is a silent illness. This disease changes the biological shapes of veins, like length and thickness, and results in cardiovascular disease, stroke, and respiratory failures over the long run [13]. Myopia is a kind of eye issue that causes critical vision loss because of diminishing epithelial tissues and change in eye color. It fundamentally modifies the visualization of any object from a certain distance by making them blurred. The fundus images can also be used to detect some other kind of anomalies like pigment epithelium proliferation, the epiretinal membrane, tessellated fundus, and vitreous macular degeneration.

## III. RELATED WORKS

Deep learning in ocular imaging can be used in conjunction with telemedicine as a possible solution for selecting, diagnosing, and controlling ophthalmic diseases for patients in primary care [14]. Recent advances in neural network approaches are at the forefront of state-of-the-art disease recognition systems [15], [16]. Many researchers have made significant efforts to resolve the multi-label classification problem of ophthalmic disease. All these works are simulated on the publicly available ODIR-5K dataset.

Islam et al. [17] proposed a shallow CNN-based model trained from scratch for classification of fundus images of ODIR-5K dataset. The left and right eye fundus images are input to the CNN model independently, and the disease label is assigned accordingly. Their approach made the disease classification model less complex, but their model is not able to distinguish multiple disease. Wang et al. [18] preprocessed fundus images using gray and color histogram equalization. Various data augmentation techniques are also used. The preprocessed gray and colored images are applied to two parallel EfficientNet models, and feature concatenation is done at the last layer for final classification. But they are able to achieve only 73% AUC and 88 % F1-Score on ODIR-5K dataset.

Li et al. [19] proposed a dense correlation network (DCNet) using transfer learning based ResNet architecture. The spatial correlation module (SCM) is the basic building block of this network architecture. The SCM block defines pixel-wise dense correlation between features extracted from color fundus images. These correlated features are fused to create the final feature map for classifying ophthalmic disease classes of ODIR-5K dataset with 93% AUC and 91.3% F1-score. Similarly, Gour and Khanna [20] proposed a pre-trained, two-input CNN architecture for the ODIR-5K dataset. They applied left and right eye fundus images to two parallel pre-trained VGG-16 simultaneously to extract the features [21], which are concatenated to create a final feature map. In spite of the use of VGG model, they failed to beat the performance of [19].

Li et al. [22] chose VGG-16, ResNet, Inception-v4, and Densenet [23] architectures with the sum, multiply, and concatenate operations on features extracted from the baseline model. They found that element-wise sum operation on feature maps yields better abnormality detection compared to the other two methods. Lin et al. [24] proposed a graph convolution network (GCN) based self-supervising learning model known as MGC-Net. GCN is utilized to capture contextual information for multi-label fundus images whereas self-supervising learning is used for generalization of the network. In comparison to the backbone network, their model showed performance enhancement for fundus disease classification on ODIR dataset. Ou et al. [25] proposed two input CNN based attention model with multi-scale module for multi-label fundus image classification. Multi-scale module utilized $3 \times 3$ and $1 \times 1$ dilated convolution filters to capture multi-scale features. A spatial attention module is used for feature enhancement and learning inter-dependency between global and local information. The model is found computationally efficient but could not beat Li et al. [19] performance-wise.

ODIR-5k has a common target label for a pair of fundus image. Due to high variation in the sample counts of each class, it suffers from the class imbalance problem. None of the state-of-the-art methods attempt to address this issue. Class imbalance is a well-known issue in medical image classification, yet limited research is available on it in the context of deep learning methods [26]. Pratt et al. [27] implemented a CNN based model for the five-class classification of diabetic retinopathy disease. The authors demonstrated that the class-weight strategy can be used to resolve over-fitting and class imbalance issues. Buda et al. [26] examined the impact of class imbalance on classification problems using CNN. In their experimentation, CIFAR-10, MNIST, and ImageNet datasets are sub-sampled to construct artificially balanced datasets.

In CNN, the receptive field can be made large with increasing kernel size, but this usually results in an increasing number of

learning parameters, which may lead to over-fitting problem [28]. To overcome this issue, Yu et al. [28] proposed dilated convolution operation which can enlarge the receptive field without adding additional computational cost. Since large receptive fields may not be able to recognise small objects, multi-scale feature extraction can be utilized to improve the image classification. Qi Zhang [29] proposed a dilated convolution based network to extract multi-scale features with increased receptive field size. The model showed improved performance by extracting broader and deeper semantic information for liver tumour classification. Similarly, Tao et al. [30] proposed a multi-scale hybrid dilated convolution module for segmentation. They used multiscale dilated convolution with variable dilation rate in the encoder-decoder architecture. Simulated on several backbone, CNN based model achieved improved performance for object segmentation. In this manuscript, a novel attention-based model is proposed with improved multi-label classification accuracy by resolving the class imbalance issue of the dataset.

## IV. Methodology

### A. ODIR-5K Dataset

ODIR-5K dataset used for experimentation in this work is made available online through a grand challenge by Peking University. The dataset contains around 5000 organized fundus images of the left and right eyes of patients. These fundus images were captured using different fundus cameras, like Kowa, Zeiss, and Cannon, having different image resolutions. Disease diagnostic keywords were assigned to these images from eye specialists [9]. Based on these diagnostic keywords, the disease classification labels are assigned to each pair of fundus images. This visual pathology dataset is unique in comparison to other publicly accessible data sets as it contains color fundus images of both left and right eyes of a patient with single/multiple abnormalities in a single image. The dataset contains a common target label for the pair of eye images. The images are grouped into eight disease classes, Normal, Diabetes, Glaucoma, Cataract, AMD, Hypertension, Myopia, and Others. Patients' age and gender are also included. Another challenging part of the ODIR-5K dataset is that the other class images contain lesions related to 12 different ophthalmic diseases. It is not easy to learn appropriate features in such a case. Also, the dataset is highly imbalanced, considering the number of images in each of the eight classes. These issues negatively affect the accuracy and loss of the trained models for classification. Although ODIR-5K dataset is more applicable to real-life clinical situations as the images are captured with different cameras in different illumination conditions, this imposes a great challenge for any disease identification model.

### B. Pre-processing

Most deep neural networks require dimension of input images in 1:1 aspect ratio. Images in the ODIR-5K dataset have different resolutions as these are captured from different cameras. Image crop operation is utilized to make these images appropriate for model training. In image crop operation, the field of view is identified by seeking the start position of non-black pixel in the input fundus image. This position is used to identify the image mask border for crop operation. To support different DNN models, the image size needs to be adjusted explicitly. Therefore, the size of cropped image is kept to $224 \times 224$ pixels as it is commonly accepted size by most of the DNN models [21], [23], [31]–[33].

### C. Class Balancing

The left and right fundus images are treated individually as input to the CNN architecture in this work. Based on the disease description

**TABLE I**
ODIR-5K DATASET CLASS STATISTICS IN ORIGINAL, AFTER OVERSAMPLING, AND AFTER UNDERSAMPLING OPERATION.

| Class (Label) | Samples | Oversampling | | Undersampling | |
|---|---|---|---|---|---|
| | | CBF | Samples | CBF | Samples |
| Normal (N) | 1135 | 0 | 1135 | 12 | 95 |
| Diabetes (D) | 1131 | 0 | 1131 | 11 | 103 |
| Glaucoma (G) | 207 | 5 | 1035 | 2 | 104 |
| Cataract (C) | 211 | 5 | 1055 | 2 | 106 |
| AMD (AMD) | 171 | 7 | 1197 | 2 | 85 |
| Hypertension (H) | 94 | 12 | 1128 | 1 | 94 |
| Myopia (M) | 177 | 6 | 1062 | 2 | 86 |
| Others (O) | 944 | 0 | 944 | 10 | 95 |

appearing for a particular eye part, the disease label is being assigned to each image. A fundus image having multiple diseases is treated separately with individual disease class labels. The images with artifacts like "low-quality image", "optical disk photographically invisible," "lens dust," and "pimage offset" are removed from the final dataset to reduce the false recognition rate. After separating left and right eye images with their corresponding ground truth, the class-wise distribution is shown in Table I. It is observed that three classes, normal, diabetes, and other classes have a significantly greater number of images in comparison to other disease classes. Random sampling techniques are used in this work to address the class balancing issue. Oversampling and undersampling are done to create two different versions of the datasets for implementation of the proposed CNN model.

*1) Oversampling:* The most popular approach for producing synthetic image samples is to generate images with random attributes for minority classes. It can be seen in Table I, classes Hypertension, Myopia, Cataract, AMD, and Glaucoma have a lesser number of samples compared to other classes. New samples for these classes are synthetically generated using different data augmentation strategies. The new sample size is calculated as per equation (1).

$$M_{minor} = N_{minor} \times (1 + k) \tag{1}$$

where $N_{minor}$ and $M_{minor}$ represent the total number of samples in a minor class and synthetically generated samples for that class, respectively. Here, $k$ is defined as a class balancing factor (CBF) having a value ranging from 1 to 13, which represent the number of augmentation operations. For data augmentation, flip, re-scaling with scaling ratio (0.5, 0.7, 0.8 and 0.9), crop, rotation, contrast change, hue, saturation, and gamma value change operations are used.

*2) Undersampling:* In this case, the whole dataset is shuffled and random images from majority classes have been selected to limit majority classes within the range of minority classes as per equation (2).

$$M_{major} = \left\lfloor \frac{N_{major}}{k} \right\rfloor \tag{2}$$

where $N_{major}, M_{major}$ represent the total number of samples in a major class and random samples for that class within a range, respectively. Again, $k$ is defined as class balancing factor (CBF).

In this manner, image statistics for all 8 disease classes get equalized.

### D. Overview of the Proposed DKCNet Architecture

As shown in Fig. 3, DKCNet is comprised of the backbone, attention block, and squeeze-excitation (SE) block. The backbone network
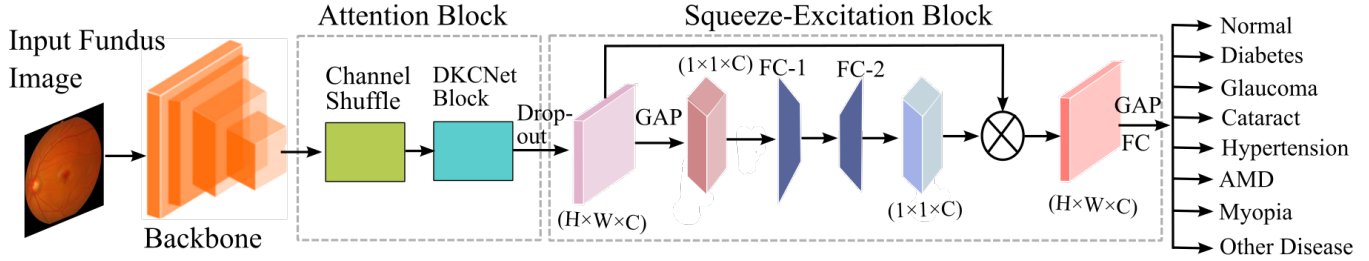
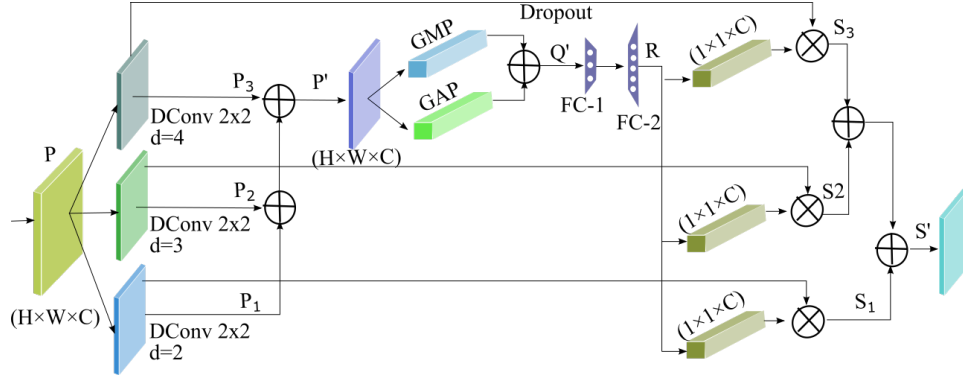Fig. 3. Block diagram of the proposed DKCNet architecture.



Fig. 4. Description of DKCNet Block.

is utilized to acquire the global feature maps. Any CNN based deep-learning model pre-trained on ImageNet dataset can be used as the backbone network to extract feature maps from the last layer of the model, which contains high-level semantic features of fundus image. These feature maps are $F_m \in R_m(H \times W \times C)$, where W, H, and C are width, height, and the number of channels in the feature maps. The output of the backbone network is fed into the attention block, which learns more region-wise features to discriminate lesion parts. A dropout layer follows this to reduce overfitting with a drop rate of 0.3. These discriminative features are processed by SE block, which dynamically re-calibrate channels. Finally, a global average pooling layer followed by a fully connected layer performs disease classification by predicting class label probability.

### E. DKCNet Block

The standard convolution with a fixed kernel captures contextual information by sliding on the feature maps. In this case, features of a similar group of pixels may have a different representation in other regions, resulting in intra-class inconsistency. It is widely accepted that more contextual information can be captured by generating multi-scale features using different receptive field sizes [34], [35]. Dilated convolution captures multi-scale information by varying kernel sizes known as dilation rate. With a large receptive field size, more semantic information can be captured. As depicted in fig. 4, dilated convolution is utilized in the proposed model to capture multi-scale features without adding extra computational cost.

The extracted feature maps from the backbone are passed to the channel shuffle block, which permutates the channel and permits data stream across feature channels. The DKCNet Block takes input from the channel shuffle block. It applies dilated convolution operation using $2 \times 2$ kernel with dilation rate ranging from 1 to 3. Dilated convolution (DConv) is followed by batch normalization ($f_{BN}$) and ReLu activation ($f_{ReLu}$) function as shown in equation (3).

$$P_i = f_{ReLu}(f_{BN}(DConv_d(P))), i = \{1, 2, 3\}, d = \{2, 3, 4\} \tag{3}$$

Features obtained after dilated convolution operation are grouped together by element wise addition operation as per equation (4).

$$P' = P_1 \oplus P_2 \oplus P_3 \tag{4}$$

Now the spatial information is squeezed from the feature maps by performing global average pooling ($f_{GAP}$) and global max-pooling ($f_{GMP}$) operations followed by element-wise addition to get global spatial information as per equation (5).

$$Q = f_{GMP}(P') \oplus f_{GAP}(P') \tag{5}$$

This squeezed spatial information is passed through two fully connected layers for channel dimension reduction by a factor $r$. Some of the features are then dropped by introducing a feature drop layer with dropout rate of 0.25 followed by sigmoid ($\sigma$) activation function as per equation (6) & (7).

$$Q' = f_{Drop=0.25}(Q) \tag{6}$$

$$R = f_{Sigmoid}(f_{FC}(Q')) \tag{7}$$

After that, obtained squeezed information vector is element-wise multiplied with feature maps obtained by dilated convolution as shown in equation (8).

$$S_i = R \otimes P_i, i = \{1, 2, 3\} \tag{8}$$

$$S' = S_1 \oplus S_2 \oplus S_3 \tag{9}$$

Finally, as per equation (9), the obtained features are added together to get the attention map.

## F. Loss Function

This work deals with a multi-class multi-label classification problem; one or more disease labels are required as the output for each left and right eye image input. For the computation of the difference between predicted target labels and actual labels, the Binary Cross-Entropy (BCE) loss function is used, which is given as:

$$BCE(y, \hat{y}) = -\frac{1}{M}\sum_{i=1}^{M} y_i \log(\hat{y}) + (1 - y_i)\log(1 - \hat{y}) \quad (10)$$

where $M$ is the number of samples in the training set, $y$ is the actual label, and $\hat{y}$ is the predicted label.

## G. Experimentation Setup

In this study, a backbone network is selected via experimentation with pre-trained ResNet, InceptionV3, and InceptionResNet architectures. The ODIR-5K dataset is split into 80% and 20% for the training and validation set, respectively. For oversampling, CBF value is selected as 12, 6, 5, 7, 5 for hypertension, myopia, cataract, AMD and glaucoma disease classes, respectively; and 0 for normal, diabetes, and other disease classes. Similarly, for under-sampling, CBF value is selected as 12, 11, 10, 1 for normal, diabetes, others, and hypertension classes; and 2 for glaucoma, cataract, AMD, and myopia class. Samples synthetically generated by both oversampling and undersampling are depicted in Table I. Model is optimized with SGD optimizer. The initial learning rate is set to 0.0005 with decay factor $1e^{-6}$ along with the batch size as 16. All the models are trained for 100 epochs with BCE loss function. All experimentation work is performed on NVIDIA T4 GPU with 16 GB memory. Two experimentation scenarios have been implemented to investigate the performance of DKCNet. In the first scenario, the model is trained with a backbone network for the classification of eight ophthalmic diseases. In the second scenario, training is done on the fusion of the backbone with DKCNet.

## V. RESULTS & DISCUSSION

The performance of the proposed model is evaluated using Area Under receiver operating Curve (AUC), F1-Score, and kappa score. AUC curve is a performance estimator for classification problems. The AUC value closer to 1 means better performance of the model to classify the disease labels. Similarly, F1-score is defined as the harmonic mean of recall and precision values as shown in equation (13). Kappa score is a proportion of how intently the instance classified results matched with the ground truth label. It is calculated as per equation (14).

$$Precision = \frac{T_p}{F_p + T_p} \quad (11)$$

$$Recall = \frac{T_p}{F_n + T_p} \quad (12)$$

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (13)$$

$$k = \frac{(p_o - p_e)}{1 - p_e} \quad (14)$$

where $T_p, F_p, F_n$ are true positive, false positive, and false negative, respectively; $p_o$ is the empirical probability of agreement on the label assigned to the sample and $p_e$ is the expected agreement when both annotators assign labels randomly.

TABLE II
OPHTHALMIC DISEASE CLASSIFICATION BY USING DKC BLOCK WITH STATE-OF-THE-ART BACKBONE NETWORKS ON ODIR-5K DATASET.

| Model | No Sampling | | Undersampling | | Oversampling | |
|---|---|---|---|---|---|---|
| | AUC | F1-Score | AUC | F1-Score | AUC | F1-Score |
| ResNet-101 | 70.12 | 78.53 | 87.04 | 89.56 | 94.33 | 92.28 |
| ResNet-101+ DKC Block | 70.26 | 79.28 | 86.53 | 88.15 | 93.52 | 91.62 |
| InceptionV3 | 72.94 | 78.17 | 87.45 | 88.93 | 86.31 | 87.37 |
| InceptionV3+ DKC Block | 73.86 | 79.44 | **88.24** | **88.93** | 88.59 | 89 |
| InceptionResNet | 74.22 | 80.44 | 86.94 | 88.68 | 94.24 | 91.53 |
| InceptionResNet + DKC Block | **74.55** | **80.96** | 88.05 | 88.93 | **95.4** | **93.18** |

TABLE III
OPHTHALMIC DISEASE CLASSIFICATION RESULTS OBTAINED WITH DIFFERENT METHODS ON ODIR-5K DATASET.

| Author | Method | AUC | F1-Score | Params (M) | Flops (G) |
|---|---|---|---|---|---|
| Islam et al. [17], 2019 | Shallow CNN | 80.5 | 85 | 1.1 | - |
| Wang et al. [18], 2020 | EfficientNet | 73 | 88 | - | - |
| Gour and Khanna [20], 2020 | Two Input VGG-16 | 84.93 | 85.57 | 15.2 | 80.2 |
| Li et al. [19], 2020 | ResNet-101 | 93 | 91.3 | 74.2 | 68.7 |
| Ning Li et. al. [22], 2021 | Inception-v4 | 88 | 85.93 | - | - |
| Lin et. al. [24], 2022 | Graph Conv. Network | 78.16 | 89.66 | - | - |
| Ou et. al. [25], 2022 | ResNet-50 | 90.3 | 88.6 | 82.6 | 67 |
| **Proposed Method** | **InceptionResnet + DKC Block** | **96.08** | **94.28** | **87.7** | **13.4** |

The effectiveness of DKCNet is experimentally investigated, and results for 10-fold cross validation are shown in Table II. The proposed method can be applied to a wide range of backbone networks to improve ophthalmic disease classification. ResNet-101, Inception V3, and InceptionResNet backbone networks are implemented with the three different baselines (without sampling, oversampling and undersampling) in the proposed approach. The increase in performance can be seen in both over-sampled and under-sampled datasets compared to without sampling. Results reported in Table II indicate that the backbone network integrated with DKCNet achieves a significant performance improvement (96.08 AUC, 94.28 F1-Score, and 0.81 kappa-score) for ophthalmic disease classification with over-sampling of minority classes for training. For each test image, we assigned target labels with a confidence greater than 0.5 to be positive, and compared them with a ground truth labels. Fig 5 shows the AUC curves for each disease class classification performance.

The proposed model is also compared with five recent multi-class, multi-label ophthalmic disease classification models to show its efficacy. The results are reported in Table III. As mentioned in Section III, Islam et al. [17] considered single shallow CNN model which could not perform well for multi-label classification. Wang et al. [18] obtained a better F1-score with EfficientNet model but lacked in AUC compared to other methods. Furthermore, Gour and Khanna [20] achieved slightly better performance. But they utilized a heavy VGG16 model, which does not contain a batch normalization layer, making it hard to converge. The performance obtained by Ning Li et.
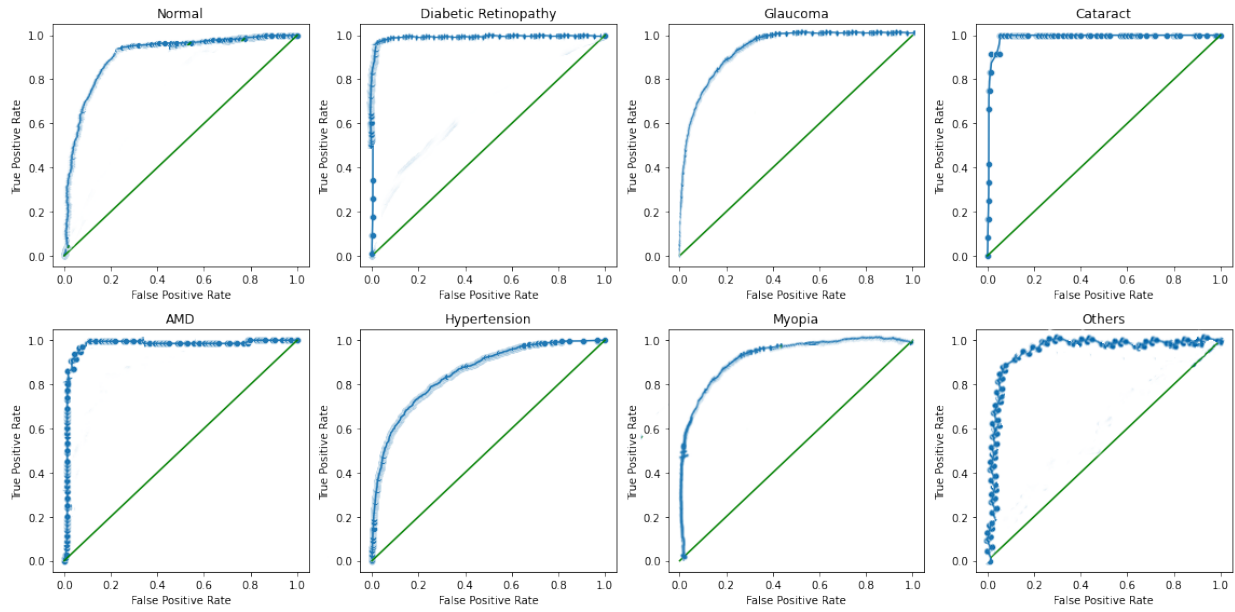
Fig. 5. AUC curves obtained with the proposed model for multi-class multi-label classification of retinal disease.
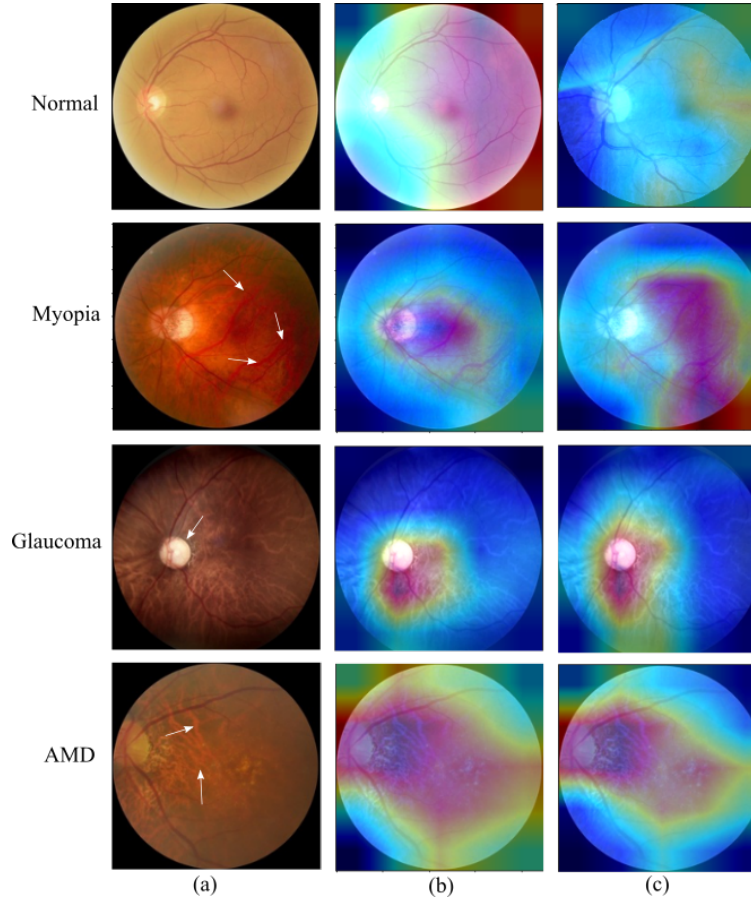


Fig. 6. Class activation maps obtained with DKCNet for single diseased fundus images from ODIR-5K dataset (a) original fundus images with lesion parts highlighted using white arrow, (b) heatmaps generated by backbone network, and (c) heatmaps refined by DKCNet.

al. [22] by using inception-v4 model with element-wise sum feature fusion is comparable with that obtained by Gour and Khanna [20]. The spatial corelation model proposed by Li et al. [19] delivered better outcomes for different mixes of ResNet structures. The model performs better among counterpart models in terms of F1-score and AUC. However, the proposed DKCNet model achieves improvement in AUC and F1-Score by 2.5% and 2.05%, respectively as compared to those achieved by Li et al. [19]. Table III also shows that the proposed network requires less number of floating-point operations per second (FLOPS) compared to state-of-the-art methods.
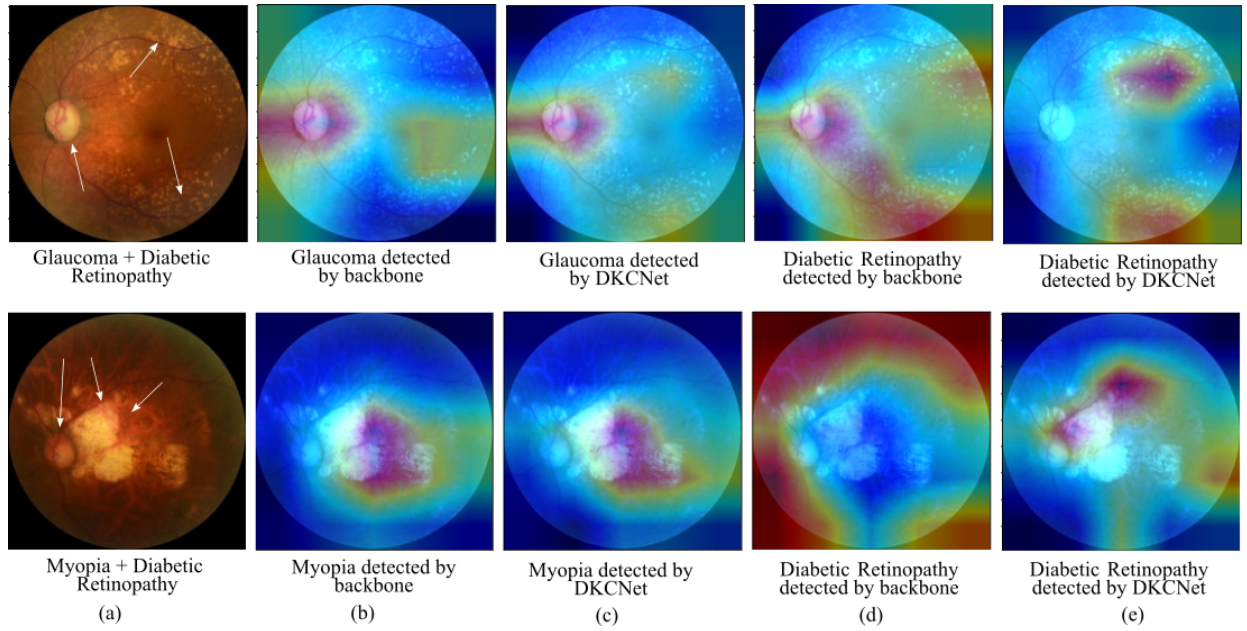
Fig. 7. Class activation maps obtained with DKCNet for multiple diseased fundus images from ODIR-5K dataset (a) original fundus images with lesion parts highlighted using white arrows, (b) heatmaps generated by backbone network, and (c) heatmaps refined by DKCNet.

TABLE IV
CROSS DATASET PERFORMANCE EVALUATION OF THE PROPOSED MODEL.

| Testing Dataset | AUC | F1-Score |
|---|---|---|
| Messidor (Diabetic Retinopathy) | 89.37 | 87.75 |
| G1020 (Glaucoma) | 93.14 | 91.42 |
| Joint Shantou International Eye Centre (Multi class) | 94.18 | 91.15 |

To check the biasness of proposed model towards training data, it is further tested on three publicly available benchmark datasets: Messidor (Diabetic Retinopathy), G1020 (Glaucoma), Joint Shantou International Eye Centre (Multi class). From the Table IV, it can be observed that the proposed DKCNet can predict retinal diseases effectively on completely unseen fundus image datasets.

Qualitative analysis of results is performed by visualizing activation maps using Grad-CAM [36]. In Fig. 6, the first column shows fundus images containing a single disease class with lesion part marked with a white arrow. The results of the backbone network, i.e., InceptionResnet are visualized in the second column. The third column shows refined activation maps obtained by using DKCNet with the backbone network. In the first row, the results generated by the backbone network show false detection of lesion part in a normal eye image, whereas the proposed model is able to discriminate such situations effectively. Similarly, the backbone network failed to detect some lesion parts in the input image in the second row, but the proposed model can identify those efficiently. The third row corresponds to the cases where the prediction results obtained with the backbone network and the proposed model are similar. In the fourth row, the backbone network falsely highlights the lesion part in the larger portion of the eye, whereas the proposed model refines the detection and is near the ground truth.

Similarly, multi-class classification performance can be visualized in Fig 7. Here, the first column shows input fundus images of patients' eyes having multiple diseases. The second and third columns correspond to class activation maps generated by the backbone network for predicted class, whereas the last two columns demonstrate a refined class activation map obtained with DKCNet with better visual classification. In the fourth column of the first row of Fig 7, it can be seen that the backbone network failed to detect some of the lesion parts, whereas the proposed method detects those well. In the second row, the backbone network shows false detection for lesion parts, whereas the proposed DKCNet highlights those parts more accurately, as seen in the fourth and the last column of Fig 7.

## VI. CONCLUSION AND FUTURE WORK

The DKCNet architecture proposed in this work enables CNN based model to learn discriminative features with an attention module without introducing extra cost. The novelty of this work lies in a model which improves ophthalmic disease classification performance and solves the class balancing issue of the highly imbalanced ODIR-5K dataset having multiple common labels for fundus image pair of a patient's left and right eye. DKCNet is composed of an attention block followed by a SE block. The attention block takes features from the backbone network and generates discriminative feature attention maps. The SE block takes the discriminative feature maps and performs channel wise attention almost at no computational cost. The experimentation has been done with three backbone CNN-based architectures, and it is found that the proposed DKCNet shows superior performance compared to counterpart methods with InceptionResnet model. As it is expensive to obtain good quality and labeled fundus images, it is planned to use generative adversarial networks to generate samples for minority classes artificially. Also, the model can be explored to localize and find the types of lesions.

## REFERENCES

[1] "World health organization - vision impairment and blindness," https://www.who.int/news-room/fact-sheets/detail/blindness-and-visual-impairment, accessed: 2021-12-18.

[2] R. Kawasaki and J. Grauslund, "Clinical motivation and the needs for ria in healthcare," in *Computational Retinal Image Analysis*. Elsevier, 2019, pp. 5–17.

[3] H. Raja, T. Hassan, M. U. Akram, and N. Werghi, "Clinically verified hybrid deep learning system for retinal ganglion cells aware grading of glaucomatous progression," *IEEE Transactions on Biomedical Engineering*, vol. 68, no. 7, pp. 2140–2151, 2021.

[4] L. Lin, M. Li, Y. Huang, P. Cheng, H. Xia, K. Wang, J. Yuan, and X. Tang, "The sustech-sysu dataset for automated exudate detection and diabetic retinopathy grading," *Scientific Data*, vol. 7, no. 1, pp. 1–10, 2020.

[5] T. B. Sekou, M. Hidane, J. Olivier, and H. Cardot, "From patch to image segmentation using fully convolutional networks–application to retinal images," *arXiv preprint arXiv:1904.03892*, 2019.

[6] H. Fu, J. Cheng, Y. Xu, D. W. K. Wong, J. Liu, and X. Cao, "Joint optic disc and cup segmentation based on multi-label deep network and polar transformation," *IEEE transactions on medical imaging*, vol. 37, no. 7, pp. 1597–1605, 2018.

[7] U. Raghavendra, H. Fujita, S. V. Bhandary, A. Gudigar, J. H. Tan, and U. R. Acharya, "Deep convolution neural network for accurate diagnosis of glaucoma using digital fundus images," *Information Sciences*, vol. 441, pp. 41–49, 2018.

[8] H. Fu, J. Cheng, Y. Xu, C. Zhang, D. W. K. Wong, J. Liu, and X. Cao, "Disc-aware ensemble network for glaucoma screening from fundus image," *IEEE transactions on medical imaging*, vol. 37, no. 11, pp. 2493–2501, 2018.

[9] "Peking university international competition on ocular disease intelligent recognition (odir-2019)," https://odir2019.grandchallenge.org/, accessed: 2022-02-10.

[10] Y. Jiang, H. Xia, Y. Xu, J. Cheng, H. Fu, L. Duan, Z. Meng, and J. Liu, "Optic disc and cup segmentation with blood vessel removal from fundus images for glaucoma detection," in *2018 40th annual international conference of the ieee engineering in medicine and biology society (EMBC)*. IEEE, 2018, pp. 862–865.

[11] L. Giancardo, F. Meriaudeau, T. P. Karnowski, Y. Li, S. Garg, K. W. Tobin Jr, and E. Chaum, "Exudate-based diabetic macular edema detection in fundus images using publicly available datasets," *Medical image analysis*, vol. 16, no. 1, pp. 216–226, 2012.

[12] E. Peli and T. Peli, "Restoration of retinal images obtained through cataracts," *IEEE transactions on medical imaging*, vol. 8, no. 4, pp. 401–406, 1989.

[13] P. Antonio, P. Marta, D. Luís, D. Antonio, S. Manuel, M. Rafael, G. Sonia, G. Manuel, M. Isabel, E. Carlos *et al.*, "Factors associated with changes in retinal microcirculation after antihypertensive treatment," *Journal of human hypertension*, vol. 28, no. 5, pp. 310–315, 2014.

[14] Z. Wang, P. A. Keane, M. Chiang, C. Y. Cheung, T. Y. Wong, and D. S. W. Ting, "Artificial intelligence and deep learning in ophthalmology," *Artificial Intelligence in Medicine*, pp. 1–34, 2020.

[15] A. Kwasigroch, B. Jarzembinski, and M. Grochowski, "Deep cnn based decision support system for detection and assessing the stage of diabetic retinopathy," in *2018 International Interdisciplinary PhD Workshop (IIPhDW)*. IEEE, 2018, pp. 111–116.

[16] M. Akil, Y. Elloumi, and R. Kachouri, "Detection of retinal abnormalities in fundus image using cnn deep learning networks," in *State of the Art in Neural Networks and their Applications*. Elsevier, 2021, pp. 19–61.

[17] M. T. Islam, S. A. Imran, A. Arefeen, M. Hasan, and C. Shahnaz, "Source and camera independent ophthalmic disease recognition from fundus image using neural network," in *2019 IEEE International Conference on Signal Processing, Information, Communication & Systems (SPICSCON)*. IEEE, 2019, pp. 59–63.

[18] J. Wang, L. Yang, Z. Huo, W. He, and J. Luo, "Multi-label classification of fundus images with efficientnet," *IEEE Access*, vol. 8, pp. 212 499–212 508, 2020.

[19] C. Li, J. Ye, J. He, S. Wang, Y. Qiao, and L. Gu, "Dense correlation network for automated multi-label ocular disease detection with paired color fundus photographs," in *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2020, pp. 1–4.

[20] N. Gour and P. Khanna, "Multi-class multi-label ophthalmological disease detection using transfer learning based convolutional neural network," *Biomedical Signal Processing and Control*, vol. 66, p. 102329, 2021.

[21] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[22] N. Li, T. Li, C. Hu, K. Wang, and H. Kang, "A benchmark of ocular disease intelligent recognition: one shot for multi-disease detection," in *International Symposium on Benchmarking, Measuring and Optimization*. Springer, 2020, pp. 177–193.

[23] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.

[24] J. Lin, Q. Cai, and M. Lin, "Multi-label classification of fundus images with graph convolutional network and self-supervised learning," *IEEE Signal Processing Letters*, vol. 28, pp. 454–458, 2021.

[25] X. Ou, L. Gao, X. Quan, H. Zhang, J. Yang, and W. Li, "Bfenet: A two-stream interaction cnn method for multi-label ophthalmic diseases classification with bilateral fundus images," *Computer Methods and Programs in Biomedicine*, p. 106739, 2022.

[26] M. Buda, A. Maki, and M. A. Mazurowski, "A systematic study of the class imbalance problem in convolutional neural networks," *Neural Networks*, vol. 106, pp. 249–259, 2018.

[27] H. Pratt, F. Coenen, D. M. Broadbent, S. P. Harding, and Y. Zheng, "Convolutional neural networks for diabetic retinopathy," *Procedia computer science*, vol. 90, pp. 200–205, 2016.

[28] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.

[29] Q. Zhang, "A novel resnet101 model based on dense dilated convolution for image classification," *SN Applied Sciences*, vol. 4, no. 1, pp. 1–13, 2022.

[30] T. Ku, Q. Yang, and H. Zhang, "Multilevel feature fusion dilated convolutional network for semantic segmentation," *International Journal of Advanced Robotic Systems*, vol. 18, no. 2, p. 17298814211007665, 2021.

[31] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.

[32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[33] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.

[34] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.

[35] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 801–818.

[36] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.