

TRk-CNN: Transferable Ranking-CNN for image classification of glaucoma, glaucoma suspect, and normal eyes

Tae Joon Jun^a, Youngsub Eom^b, Dohyeun Kim^a, Cherry Kim^c, Ji-Hye Park^b,
Hoang Minh Nguyen^a, Daeyoung Kim^{a,*}

^a *School of Computing, Korea Advanced Institute of Science and Technology,
34141 Daejeon, Republic of Korea*

^b *Department of Ophthalmology, Korea University College of Medicine
02841, Seoul, Republic of Korea.*

^c *Department of Radiology, Korea University College of Medicine
02841, Seoul, Republic of Korea.*

Abstract

In this paper, we proposed Transferable Ranking Convolutional Neural Network (TRk-CNN) that can be effectively applied when the classes of images to be classified show a high correlation with each other. The multi-class classification method based on the softmax function, which is generally used, is not effective in this case because the inter-class relationship is ignored. Although there is a Ranking-CNN that takes into account the ordinal classes, it cannot reflect the inter-class relationship to the final prediction. TRk-CNN, on the other hand, combines the weights of the primitive classification model to reflect the inter-class information to the final classification phase. We evaluated TRk-CNN in glaucoma image dataset that was labeled into three classes: normal, glaucoma suspect, and glaucoma eyes. Based on the literature we surveyed, this study is the first to classify three status of glaucoma fundus image dataset into three different classes. We compared the evaluation results of TRk-CNN with Ranking-CNN (Rk-CNN) and multi-class CNN (MC-CNN) using the DenseNet as the backbone CNN model. As a result, TRk-CNN achieved an average accuracy of 92.96%, specificity of 93.33%, sensitivity for glaucoma suspect of 95.12%

*Corresponding author

Email address: `kimd@kaist.ac.kr` (Daeyoung Kim)

and sensitivity for glaucoma of 93.98%. Based on average accuracy, TRk-CNN is 8.04% and 9.54% higher than Rk-CNN and MC-CNN and surprisingly 26.83% higher for sensitivity for suspicious than multi-class CNN. Our TRk-CNN is expected to be effectively applied to the medical image classification problem where the disease state is continuous and increases in the positive class direction.

Keywords: Glaucoma; Glaucoma suspect; Convolutional neural networks; Ranking classification

1. Introduction

The rapid development of deep learning technologies, especially convolutional neural network (CNN), is now considered to be a cutting-edge methodology for classifying medical images. The vast majority of recent medical image analysis literature uses CNN-based methodologies. The main reason CNN is effective in medical image analysis is that CNN is trained end-to-end. In other words, CNN's automated feature extraction process is more effective than traditional handcrafted feature extraction methods. However medical image classes are distinguished from general image classes. That is, the classes of medical images have a strong correlation with each other. In particular, there are innumerable intermediate states between the negative class, which is classified as normal, and the positive class, which is classified as disease. In addition, the negative class proceeds to a positive class in the direction of increasing the inherent characteristic. This characteristic depends on the type of medical image to be classified. For example, the cancer staging using the TNM system includes the size of the tumor [1], and in a cataract patient, the degree of turbidity of the ocular lens may be increased [2].

Thus, the actual disease state is continuous and increases in the positive class direction. However, when the medical image is taken and the physician makes a decision, the class of the medical image is determined based on a certain point on the continuous line. Therefore, depending on the disease, the intermediate class of the medical image may be defined by the physician, not

dichotomous separation into normal and disease. Glaucoma is a representative example of such diseases. The reason is that glaucoma should be treated appropriately before advanced stages where they are already positive, and the disease is worsening over a long period, it is necessary to observe persistent intermediate conditions. Glaucoma is an eye disease that causes narrowed vision and eventually leads to blindness, which is caused by various reasons such as elevated intra-ocular pressure (IOP) or blood circulation disorder [3]. Once glaucoma is diagnosed, it needs constant management for a lifetime, and the damaged vision is not restored. Therefore, early detection and treatment of glaucoma is the best prevention, but the optic nerve damage caused by glaucoma gradually develops, and when symptoms appear, the disease progresses considerably. In addition, since it is not easy to confirm glaucoma early, various tests including IOP measurement, optic nerve head examination, and anterior chamber angle examination are conducted and the results are combined to determine the existence of glaucoma.

As a result, recent glaucoma fundus image dataset includes the glaucoma suspect class and there are several existing studies that detect glaucoma using machine learning methods. Most of them use multi-class classification method that uses CNN as a classifier and utilizes the output values of softmax function. The literature on classification of glaucoma from fundus images will be discussed in more detail in the related work section. Although such machine learning based eye disease classification studies show reasonable performance, this multi-class classification method ignores inter-class information of eye diseases. In addition, in the binary classification problem of classifying normal and glaucoma, the addition of suspect class results in poor overall classification performance. In other words, in the case of diseases that show a sequential relationship among medical image classes, a method that can classify them considering the inter-class relationship is required.

Therefore, we propose a Transferable Ranking-CNN (TRk-CNN) for glaucoma detection considering information between three different fundus image classes. TRk-CNN consists of the following steps: primitive classification, re-

gion of interest (ROI) extraction, and final classification. Primitive classification follows the general Ranking-CNN [4] procedure. Ranking-CNN will be described in detail in the later sections. Briefly, it is a method of aggregating results of $N - 1$ binary classifiers to classify N number of ordinal classes. More specifically, when classifying N ordinal classes, k -th sub-classifier determine whether the predicted class is higher than the class k which ranges between 1 to $N - 1$. The difference from the original Ranking-CNN in primitive classification is that there are no fully-connected layers at the top-layers of the CNN classifier that performs binary classification. As a result of the primitive classification, we get the Class Activation Map (CAM) [5] for the predicted class. The CAM will also be discussed in detail later, but in a nutshell, it includes the importance of which spatial location in the input image highly affects the final prediction. The CAMs obtained from the $N-1$ sub-classifiers are combined into a single ROI based on the inter-class distance metrics definition, and the process of extracting the ROI and combining it with the original input is processed in the ROI extraction step. The new input, combined with the ROI, is used as an input to the final classification step. In this step, the final class is predicted through a sophisticated classification process including a fully-connected layer.

We evaluated TRk-CNN in glaucoma image dataset that was collected and labeled from Korea University Medical Center. Glaucoma dataset was labeled into three classes: normal, glaucoma suspect, and glaucoma eyes. Based on the literature we surveyed, this study is the first to classify three status of glaucoma fundus image dataset into three different classes. We compared the evaluation results of TRk-CNN with multi-class CNN (MC-CNN) and Ranking-CNN (Rk-CNN) using the DenseNet [6] as the backbone CNN model. As a result, TRk-CNN achieved an average accuracy of 92.96%, specificity of 93.33%, sensitivity for glaucoma suspect of 95.12% and sensitivity for glaucoma of 93.98%. Based on average accuracy, TRk-CNN is 8.04% and 9.54% higher than Rk-CNN and MC-CNN and surprisingly 26.83% higher for sensitivity for suspicious than MC-CNN.

The major contribution of this work is summarized as follows:

- Our proposed TRk-CNN is a method that can be effectively applied when the classes of images to be classified show a high correlation with each other. The multi-class classification method based on the softmax function, which is generally used, is not effective in this case because the inter-class relationship is ignored. Although there is the Ranking-CNN that takes into account the ordinal classes, it cannot reflect the inter-class relationship to the final prediction. TRk-CNN, on the other hand, combines the weights of the primitive classification model to reflect the inter-class information to the final classification phase. Through extensive experiments, we show that TRk-CNN is superior to both multi-class classification method and Ranking-CNN method.
- We evaluated TRk-CNN in glaucoma fundus images. Glaucoma can be labeled with suspicious states because it is important to find and take proper treatment before the condition becomes severe. We think that this is not a problem specific to glaucoma. Many diseases requiring medical imaging have intermediate states from negative class to positive class. Our TRk-CNN is expected to be effectively applied to those medical image classification problem using CNN.

The abstract version of this paper has been published in [7]. Compared with [7], this paper presents TRk-CNN as a general classification model that can be applied not only to three classes but also to N number of classes. We have also noticed that [7] showed an unusually high classification accuracy because the train-set and test-set of primitive and final classification steps are divided based on different random seeds. We have corrected the above error in this paper. In addition, a more robust evaluation was conducted to compare with the results of previous glaucoma detection studies. The rest of this paper is structured as followed. In Section 2, we review the literature using a machine-learning approach that includes deep-learning for glaucoma detection and also briefly review the multi-class classification and Ranking-CNN that is the background of this study. Section 3 explains in detail the three steps of TRk-CNN in the

general example of classifying N different classes. Section 4 describes the optimal TRk-CNN for glaucoma detection. In Section 5, we evaluate TRk-CNN in glaucoma dataset and compares the result with multi-class CNN and Ranking-CNN results. Finally, we conclude this study in Section 6 and discusses future plans.

2. Related Work

2.1. Glaucoma detection

Glaucoma is a disease in which the optic nerve and nerve fiber layers, which play an important role in delivering visual information received from the eye to the brain, are damaged and the visual field becomes narrower. Globally, glaucoma is a major cause of blindness, along with cataracts and diabetic retinopathy, and is one of the most common ophthalmic diseases, with a frequency of 2% of the total population [8] [9] [10]. In the past, glaucoma generally included increased intra-ocular pressure, but recently, normal tension glaucoma is a very common disease, and the definition of glaucoma has also changed. Primary open-angle glaucoma and normal-tension glaucoma, which account for the vast majority of glaucoma, chronically and slowly damage the optic nerve [11]. As a result, visual field damage progresses, damage to the peripheral vision first occurs, and central vision is often preserved until the end of the period. Therefore in the beginning, there is almost no subjective symptom and symptoms do not appear until glaucoma has progressed to advanced stages. As a result, most of the patients diagnosed with glaucoma are found incidentally through ophthalmologic examination or physical examination regardless of the glaucoma related symptoms. Figure 1 shows the progression of optic disc changes and visual field defects with normal, glaucoma suspect, and glaucoma eyes.

To overcome the difficulties in early diagnosis, applying machine learning methods to classify normal and glaucoma in fundus image have been proposed to play a supporting role in physician’s glaucoma diagnosis criteria.

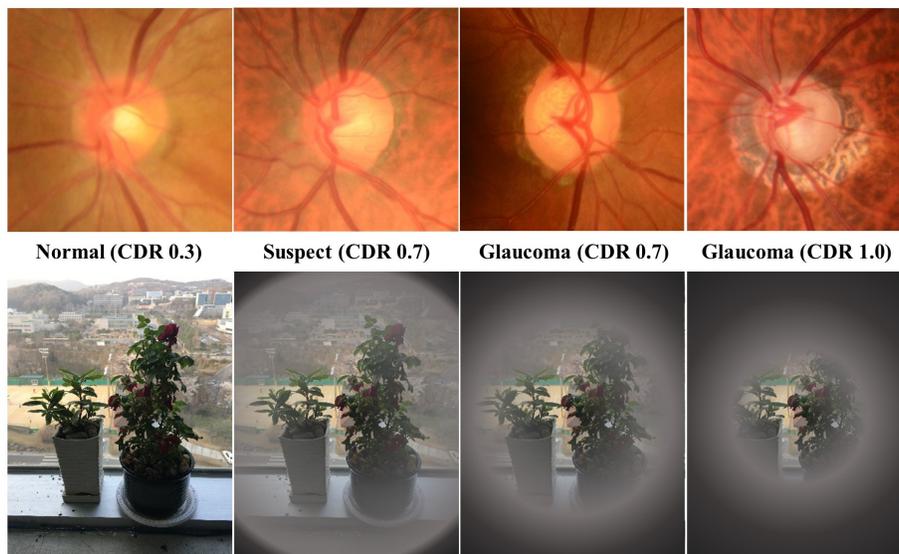


Figure 1: Optic disc changes and visual field loss with normal, glaucoma suspect, and glaucoma eyes

In 2009, Nayak proposed a method to classify normal and glaucoma with single hidden-layer neural network (ANN) by extracting features such as cup-to-disc ratio (CDR), optic nerve head shift, and ISNT ratio from the fundus image [12]. ISNT ratio is the total area of the blood vessels in the inferior and superior side of the optic disc to the total area of the blood vessels in the nasal and temporal area. Of the 24 normal and 37 glaucoma images, 5 normal and 10 glaucoma images were split into test-set. As a result, the specificity (Sp) was 80% and the sensitivity (Se) was 100%. Nayak's work is meaningful in that it extracts features and train them by the neural network, although the number of images is too small.

Bock proposed a method for extracting a probabilistic feature for glaucoma diagnosis from a fundus image called glaucoma risk index (GRI) in 2010 [13]. First, they perform pre-processing procedures such as illumination correction, vessel removal, and optic nerve head normalization. Then, Fourier analysis and spline interpolation are applied, and principal component analysis (PCA) is

performed to extract features. Finally, the features extracted by the PCA are passed into two-stage support vector machine (SVM) classifier and finally the classifier outputs a GRI indicating the probability for glaucoma. For 575 fundus images consisting of 336 normal and 239 glaucoma, the GRI method showed overall 80% accuracy with the area under the receiver operating characteristic (ROC) curve (AUC) of 0.88, sensitivity of 73%, and specificity of 85%. Bock work is also a representative approach to extract handcrafted features from images and use them as inputs for classifiers such as SVM and neural networks.

In 2011, Acharya proposed a method for extracting higher order spectra (HOS) parameter and texture descriptors from a fundus image and use them as inputs for four different classifiers [14]. Classifiers are SVM, sequential minimal optimization (SMO), naive Bayesian, and random-forest. As a result, random forest classifier showed the best performance with accuracy (Acc) of 91.7% in 60 fundus images composed of 30 normal and 30 glaucoma eyes. For the same dataset as Acharya's work, Dua proposed a method for extracting energy signatures as a feature by applying a 2-dimensional discrete wavelet transform to fundus images in 2012 [15]. Again for the four classifiers including SVM, SMO, naive Bayesian (NB), and random-forest (RF), Dua's work achieved the highest accuracy of 93.33% in both SVM and SMO classifiers.

From 2015, glaucoma detection studies based on convolutional neural networks have become mainstream with the rapid development of deep learning technology. Chen performed a classification of normal and glaucoma fundus images using CNN in 2015 [16]. Chen designed the AlexNet [17] based CNN model, and evaluated with the ORIGA [18] and SCES [19] fundus image dataset. The ORIGA dataset composed of 168 glaucoma and 482 normal fundus images and SCES dataset contains 1676 fundus images including 46 glaucoma cases. As a result, Chen obtained 0.831 and 0.887 AUC on ORIGA and SCES dataset. Chen's work is meaningful in that it is the first study which applied CNN's end-to-end training to glaucoma detection, deviating from the conventional manual feature extraction method. However, it did not perform better than the existing method because it simply applied CNN and did not refine the sophisticated

optimization process. In 2016, Li proposed a method to apply CNN models to the disc region and the original fundus image, respectively, and ensemble the predictions [20]. Li used four well known CNN models including AlexNet [17], GoogLeNet [21], 16-layer VGGNet [22], and 19-layer VGGNet [22]. Evaluated with ORIGA dataset, Li achieved AUC of 0.838. CNN is adopted and considering that the classification was binary classification, performance is not good, and similar to Chen’s work, there is a limitation that CNN model optimization is not sophisticated.

In 2018, Fu proposed a disc-aware ensemble network for glaucoma classification [23]. U-Net [24] was used for disc region segmentation and re-applied the resulting region to the original image to reduce the size of the input. Finally, 50-layer ResNet [25] was applied to fundus images of the various regions including disc region and original fundus images. The evaluation was performed in SCES and Singapore Indian Eye Study (SINDI) [23] dataset and showed 0.918 AUC and 0.817 AUC, respectively. SINDI dataset contains a total of 5783 fundus images including 113 glaucoma and 5670 normal eyes. Fu’s work has ensured the results by applying CNN to various regions similar to Li’s work [20], and the CNN model is well optimized. Also in 2018, Li classified the glaucoma eyes by applying the GoogLeNet to 48116 fundus images, which is the largest number of a dataset in the literature [26]. They also labeled the dataset as normal, glaucoma suspect, and glaucoma eyes, same as in our study. Dataset consists of 31745 train-set and 8000 test-set images. The train-set consists of 23433 normal, 2190 glaucoma suspect, and 6122 glaucoma eyes. The test-set consists of 6033 normal, 430 glaucoma suspect, and 1537 glaucoma eyes. However, the evaluation was performed as a binary classification to classify normal and abnormal (glaucoma suspect and glaucoma cases). As a result, they obtained 0.986 AUC, sensitivity of 95.6%, and specificity of 92%.

Overall, none of the studies described above take into account to classify the three continuous classes of normal, glaucoma suspect, and glaucoma eyes. Li’s work [26] is the only one that labels the fundus image in three classes but performs the binary classification by treating glaucoma suspect and glaucoma as

a single positive class. As we will see later in the evaluation, binary classification of fundus images with CNN models of the same structure is 10% higher overall accuracy than three class classification. Therefore, in order to improve the performance of the normal, glaucoma suspect, and glaucoma classification, TRk-CNN which considers inter-class information is necessary. In addition, TRk-CNN can be effectively applied to the classification of other medical images having intermediate stages between negative and positive cases.

2.2. Multi-class classification and Ranking-CNN

The multi-class classification is a method in which the size of the final prediction vector is N for N number of classes. In addition, the N different classes are converted to one-hot-encoding, where the index to which they belong is 1 and the remainder is 0. Generally, in deep learning, the softmax function is applied to the output vector to express as the probability between 0 and 1, although it is not the actual probability, and predict the class with the largest probability as the final class. In this case, the cross entropy of the probability of a class that is a true class becomes a loss, which is an error. Therefore, in the next epoch of training, gradient descent is processed in the direction of reducing this loss. However, when classes are highly related to each other, their inter-class relationship disappears because classes are one-hot-encoded in multi-class classification. Especially, the age prediction problem is where this problem is obvious. For example, in the case of classifying tree, truck, and cat images, there is no problem in classifying $[1,0,0]$, $[0,1,0]$, and $[0,0,1]$ through one-hot-encoding. However, when one-hot-encoding is used to classify 10-year-old, 11-year-old, and 12-year-old face images, the ordinal relationship of the age disappears.

Ranking-CNN was proposed by Chen in for age estimation from human face images [4]. Prior to Ranking-CNN, ranking algorithms for machine learning-based age estimation such as Ranking SVM [27], Rank-Boost [28] [29], and RankNet [30] were introduced. Ranking-CNN proposed a ranking algorithm suitable for CNN-based facial age estimation problem. In the case of classifying

N different ages from images, Ranking-CNN creates $N-1$ sub-CNN models, and each model performs binary classification with one age as a reference point. For example, when predicting the ages of 10 to 19-year-old faces, the first CNN model classifies whether the face age is older than 10 years or not. Similarly, the i -th sub-CNN model classifies facial images that are older than i years old and continues until the 9-th sub-CNN model. For a single facial image, nine different $[0,1]$ are output as the result, and the final age is determined based on the sum of these values. The major contribution of Ranking-CNN is that by taking the ordinal relation between ages into consideration, Ranking-CNN is more likely to get smaller estimation errors when compared with multi-class classification approaches [4].

However, since Ranking-CNN considers only the final binary value of the trained sub-CNN models, features extracted during the training of each sub-CNN model cannot be transferred. In addition, age is an ordinal relationship, but the classes of medical data like normal, glaucoma suspect, and glaucoma are not always directly proportional to the class relationship. Therefore, our proposed TRk-CNN can achieve higher accuracy by allowing each sub-CNN model to transfer the extracted high-dimensional features.

3. Transferable Ranking-CNN

TRk-CNN consists of the following steps: Primitive classification, ROI extraction, and Final classification. Figure 2 shows the overall structure of TRk-CNN including primitive classification, ROI extraction, and final classification steps.

Primitive classification step follows the general Ranking-CNN procedure and its purpose is to extract the major features of the reference class of each sub-CNN model. The major feature here is that each sub-CNN model should extract different features according to the result of performing a binary classification on a given input image. Therefore, we can not generally use the weight of the last convolutional layer of well known CNN such as VGGNet, GoogLeNet, and

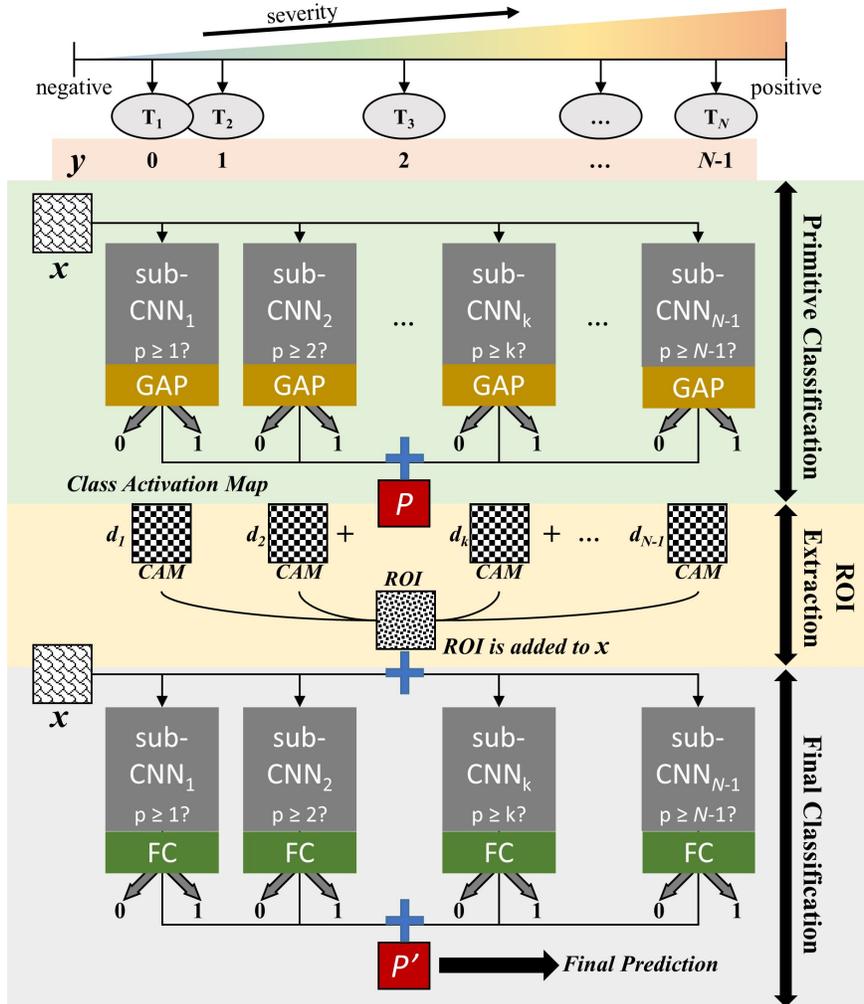


Figure 2: The overall architecture of TRk-CNN

ResNet. The reason is that the weight of the last convolutional layer contains the general characteristics of the entire dataset, but of each of the input image, the weight does not include the characteristics of the classification results the sub-CNN model. Features extracted from each sub-CNN model in TRk-CNN should individually represent the characteristics of N different classes. In order to satisfy these requirements, Class Activation Map (CAM) is extracted for each

input image as a transferable feature of each sub-CNN model. A more detailed description of CAM will be discussed in the later section. In other words, the purpose of primitive classification is to obtain the CAM for input image from each sub-CNN model through training.

In ROI extraction step, CAM extracted from each sub-CNN model is merged into a single ROI. However, when CAMs are combined through simple summation, the low relevant classes and the high relevant classes are treated equally. So we take into account the association between classes by including the distance function when merging the CAMs. In this paper, we define the distance function assuming that classes have a linear relationship. However, the distance function depends on how the domain expert defines the relationship between classes. For example, a linear relationship is reasonable for age prediction, but it is highly likely that it will not be linear in medical data. As a result, the ROI extraction step is to combine these distance functions with the CAM to create the final ROI of each input image and pass the generated ROI to the final classification step.

Final classification step combines the ROI, which received from the previous ROI extraction stage, with the original image to create a new input for classification. Although there are many possible ways to combine the ROI with the original image, we concatenate the ROIs on the additional channels of the input to preserve the information of the original image. In other words, if the original image has three channels, the number of channels for the new input is now four. We will explain other possible methods in more detail in the later section. Since this step leads to the final prediction, hyper-parameter tuning is strict and regularization is applied more strongly than primitive classification. In addition, the final classification also follows the Ranking-CNN structure and starts to converge from the earlier epoch by loading the pre-trained weights of the model from the primitive classification stage. A detailed description of each stage is provided in the following section.

3.1. Primitive Classification

Since primitive classification is almost similar to general Ranking-CNN, we will explain it with the notation from the original paper almost as it is. Let us first assume that the number of classes in the image dataset X we want to classify is N . Each class is labeled from 0 to $N-1$ depending on the direction in which the state of the class is increasing. Here, the examples of increasing refers to the age in the facial age estimation problem and severity of lesion in the medical image classification problem. When the arbitrary sample belonging to dataset X is x , the corresponding label of x is y , where $y \in \{0, 1, \dots, N-1\}$. As described in the related work section, Ranking-CNN creates $N-1$ number of sub-CNN models to classify dataset X . The role of k -th sub-CNN model is to perform binary classification in dataset X based on reference class k . If x is classified to be greater than or equal to k , the output is 1 and if it is classified to be smaller than k , the output is 0. After training k -th sub-CNN model, dataset X is divided into two subsets as shown below.

$$\begin{aligned} X_k^0 &= \{(x, 0) | y < k\} \\ X_k^1 &= \{(x, 1) | y \geq k\} \end{aligned} \tag{1}$$

Let the output value of the k -th sub-CNN model for arbitrary input x is $p_k(x)$ where the value is 0 or 1. The role of primitive classification here is to optimize each k -th sub-CNN model to minimize the binary classification error. After error is reduced enough, we aggregate the $p_k(x)$ of all sub-CNN models for arbitrary input x as follows.

$$P(x) = \sum_{k=1}^{N-1} p_k(x) \tag{2}$$

where $P(x)$ corresponds to the predicted class of primitive classification for arbitrary input x . The important point here is that the class we deliver to the ROI extraction step should be the predicted class $P(x)$, not the actual class y . The reason is that if the ROI is created through an actual class, we can not generate the ROI for test-set where the actual class is only available in the final evaluation phase. In other words, if ROI is created with an actual class, test-set

cannot be evaluated after the final classification step because the input is an original image without ROI. In the primitive classification, the fully-connected layer cannot come after the last convolutional layer, and the class classifier should follow immediately after the Global Average Pooling (GAP) layer. The reason for this will be explained in detail in the next section, ROI extraction step. Algorithm 1 provides the entire process of training and validation procedure of primitive classification step.

Algorithm 1 Primitive Classification

```

1: procedure TRAINING PROCEDURE
2:   for  $k = 1$  to  $N - 1$  do
3:     initialize  $k$ -th sub-CNN
4:   top:
5:     for  $k = 1$  to  $N - 1$  do
6:        $X^0_k = \{(x, 0) | y < k\}$ 
7:        $X^1_k = \{(x, 1) | y \geq k\}$ 
8:       fine-tune  $k$ -th sub-CNN
9:       if not converged then
10:        goto top
11: procedure PREDICTION PROCEDURE
12:   for  $k = 1$  to  $N - 1$  do
13:      $p^k(x) \leftarrow k$ -th sub-CNN
14:    $P(x) \leftarrow \sum_{k=1}^{N-1} p^k(x)$ 

```

3.2. ROI Extraction

The outputs from primitive classification step to ROI extraction step are the predicted value $P(x)$ for input x and the weights of trained sub-CNN models. In the previous section, we explained that the Global Average Pooling layer comes after the last convolutional layer of each sub-CNN model, and the fully-connected layers can not. The reason is that Class Activation Map is the feature of input x that we want to extract from each sub-CNN model and it requires

GAP layer directly after the last convolutional layer. CAM is a concept introduced by Zhou, and it schematically shows which spatial location of the input image played an important role when classified as the final class [5]. In general, combining CAM with original input in case of training a single CNN model is not expected to have a great effect on performance, but when combining results trained by multiple CNN models, such as TRk-CNN, CAM can be used to transfer important features between CNN models.

Lets assume that $f_m^k(i,j)$ is the activation result of filter $m \in \{1, 2, \dots, n\}$ in the last convolutional layer of k -th sub-CNN model at spatial location (i,j) of filter m . The size of filter m depends on the pooling policy of the sub-CNN model. Suppose the sub-CNN model performs stride 2 pooling, which is a general situation, for l number of times. When the size of input x is $h \times h$, the size of filter m becomes $h/2^l \times h/2^l$. Finally, the result F_m^k obtained from applying GAP layer to filter m can be expressed by the following equation.

$$F_m^k = \sum_{(i,j=1)}^{h/2^l} f_m^k(i,j) \quad (3)$$

From the primitive classification step, predicted class $p(x)$ is either 0 or 1. Thus, if the predicted class $p(x)$ is 1 in k -th sub-CNN model, the input S_1^k for the softmax layer as final prediction can be expressed by the following equation.

$$S_1^k = \sum_{m=1}^n w_m^1 F_m^k \quad (4)$$

where w_m^1 represents the weights between m -th node of GAP layer and class 1 node in softmax layer and n refers the total number of filters in the last convolutional layer. Substituting F_m^k with equation 3 into S_1^k yields the following

equation.

$$\begin{aligned}
S_1^k &= \sum_{m=1}^n w_m^1 F_m^k \\
&= \sum_{m=1}^n w_m^1 \sum_{(i,j=1)}^{h/2^l} f_m^k(i,j) \\
&= \sum_{(i,j=1)}^{h/2^l} \sum_{m=1}^n w_m^1 f_m^k(i,j) \\
&= \sum_{(i,j=1)}^{h/2^l} C_1^k(i,j)
\end{aligned} \tag{5}$$

where $C_1^k(i,j)$ is the Class Activation Map for (i,j) spatial location in k -th sub-CNN model for predicted class 1. Since the size of input x is $h \times h$, resizing $C_1^k(i,j)$ by h/\sqrt{n} gives the same size as input x and we can define it as $C_1^k(x)$. From the equation 5, it can be said that $C_1^k(i,j)$ indicates the importance of the activation at spatial location (i,j) leading to the classification to predicted class 1 in k -th sub-CNN model. Likewise, $C_1^k(x)$ represents which pixels of input x played an important role in classifying input x as a predicted class 1 in k -th sub-CNN model. Based on the equations described so far, $C_0^k(x)$ can be defined as the CAM for predicted class 0 in k -th sub-CNN model for given input x .

So far we have explained the CAM generation process for input x at each sub-CNN model. As a result, input x generates two types of CAMs, $C_0^k(x)$ and $C_1^k(x)$, in k -th sub-CNN model. The next thing to define is combining these $C_0^k(x)$ and $C_1^k(x)$ into unified feature for input x for aggregated predicted class $P(x)$ from primitive classification. This unified feature can be seen as Region of Interest (ROI) and defined as $R(x)$. When generating $R(x)$, we need to consider that the more distant $P(x)$ and k are, the lower the effect of $C_0^k(x)$ and $C_1^k(x)$. For example, if the predicted age at the facial age estimation problem is 20 years old, it is obvious that the sub-CNN model classified by age 19 has a higher influence than the sub-CNN model classified by age 50. Therefore, we introduce distance metric $D_P^k(x)$ to quantify this influence of

CAM for input x in k -th sub-CNN model. We can define $D_P^k(x)$ by directly applying background information of inter-class relation. If the actual class x has an ordinal relationship, such as an facial age estimation problem, $D_P^k(x)$ can be expressed by the following equation.

$$D_P^k(x) = \begin{cases} \frac{1}{P(x)-k+1}, & k \leq P(x) \\ \frac{1}{k-P(x)}, & k > P(x) \end{cases} \quad (6)$$

where $P(x) \in \{1, 2, \dots, N-2\}$. If $P(x)$ is 0 or $N-1$, $D_P^k(x)$ is not needed because $R(x)$ is defined differently. Combining $C_0^k(x)$ and $C_1^k(x)$ with $D_P^k(x)$, $M(x)$ can be defined as follows.

$$R(x) = \begin{cases} \sum_{k \leq P(x)} D_P^k(x) C_0^k(x) + \sum_{k > P(x)} D_P^k(x) C_1^k(x), & P(x) \in \{1, \dots, N-2\} \\ C_0^1(x), & P(x) = 0 \\ C_1^{N-1}(x), & P(x) = N-1 \end{cases} \quad (7)$$

From the equation 7, when k is less than or equal to $P(x)$, we multiply the distance metric $D_P^k(x)$ by the $C_0^k(x)$ of the k -th sub-CNN model. Otherwise, we multiply the $D_P^k(x)$ with $C_1^k(x)$. This part can be reversed according to the definition of the user, but from our experimental results, it was better to define it as above. A more intuitive reason is as follows. When $P(x)$ is aggregated with $p_k(x)$, $p_k(x)$ is likely to be 1 in the k -th sub-CNN model where k is less than or equal to $P(x)$. For facial age estimation example, if the predicted age is 20 years old, then it is likely that the sub-CNN model classified by age 15 is likely to have output 1 and the model by age 30 is likely to have output 0. In other words, it can be assumed that the abstract representation of $C_1^k(x)$ is already contained in $P(x)$ if k is less than or equal to $P(x)$. Therefore, if we create a $R(x)$ by aggregating the opposite class CAMs, it is presumed that final classification process can be trained with various information which is more likely to correct error of $P(x)$ with higher probability. We experiment on both combinations and compare the results later in the evaluation. In addition, when $P(x)$ is 0 or $N-1$, only the CAM from the first or the $N-1$ th sub-CNN model are

$R(x)$ without considering the other sub-CNN models. This is because Ranking-CNN performs one-to-all classification for classes at both ends. That is, when $P(x)$ is 0, the first sub-CNN model can be thought of as a model that directly classifies $P(x)$ equals 0 and vice versa in case of $P(x)$ equals $N-1$. Therefore, when $P(x)$ is 0, it is reasonable to set $R(x)$ directly with $C_0^1(x)$ and $C_1^{N-1}(x)$ when $P(x)$ is $N-1$.

The ROI extraction step can be summarized as generating $R(x)$ for the arbitrary input x from $P(x)$ and weights of sub-CNN models in the primitive classification and passing it to the final classification step. Algorithm 2 provide the entire process of ROI extraction step.

Algorithm 2 ROI Extraction

```

1: procedure CAM GENERATION PROCEDURE
2:   for  $k = 1$  to  $N - 1$  do
3:      $C_0^k(x) \leftarrow k$ -th sub-CNN
4:      $C_1^k(x) \leftarrow k$ -th sub-CNN
5: procedure ROI GENERATION PROCEDURE
6:    $P(x) \leftarrow$  Prediction Procedure in Algorithm1
7:   if  $P(x) = 0$  then
8:      $R(x) \leftarrow C_0^1(x)$ 
9:   else if  $P(x) = N - 1$  then
10:     $R(x) \leftarrow C_1^{N-1}(x)$ 
11:   else
12:     $R(x) \leftarrow \sum_{k \leq P(x)} D_P^k(x) C_0^k(x) + \sum_{k > P(x)} D_P^k(x) C_1^k(x)$ 

```

3.3. Final Classification

The role of the final classification step is to combine the $R(x)$ received from the ROI extraction step with the arbitrary input $x \in X$ to generate a new input $x' \in X'$ and perform strict training for final prediction. Algorithm 3 represents the overall process of final classification step from input x' generation to final prediction.

Algorithm 3 Final Classification

1: **procedure** GENERATE INPUT PROCEDURE

2: $R(x) \leftarrow$ *ROI Generation Procedure in Algorithm 2*

3: **for** $x \in X$ **do**

4: $x' \leftarrow x + R(x)$

5: $x' \in X'$

6: **procedure** FINAL TRAINING PROCEDURE

7: **for** $k' = 1$ *to* $N - 1$ **do**

8: **initialize** k' -th sub-CNN

9: *top:*

10: **for** $k' = 1$ *to* $N - 1$ **do**

11: $X^0_{k'} = \{(x', 0) | y < k'\}$

12: $X^1_{k'} = \{(x', 1) | y \geq k'\}$

13: **fine-tune** k' -th sub-CNN

14: **if** *not converged* **then**

15: **goto** *top*

16: **procedure** FINAL PREDICTION PROCEDURE

17: **for** $k' = 1$ *to* $N - 1$ **do**

18: $p^{k'}(x') \leftarrow$ k' -th sub-CNN

19: $P(x') \leftarrow \sum_{k'=1}^{N-1} p^{k'}(x')$

There are several ways to combine input x and $R(x)$, but we define input x' with additional channels for $R(x)$ to preserve the information of original x . That is, when input x is the size of $h \times h \times 3$, then the new input x' is the size of $h \times h \times 4$ with the $R(x)$ of size $h \times h$ appended. The advantage of this method is that even if the input x is augmented during the training, the spatial information of $R(x)$ can be maintained by applying same augmentation policy. In other words, if input x' is shifted, rotated, and resized, both input x and $R(x)$ are applied in the same way. The process of classifying the input x' is similar to the primitive classification, but there is no need to output CAM, so adding fully-connected layer after the last convolutional layer is no longer restricted. Once the training is finished, evaluation procedure is done with the test-set that was separated from the beginning. As we mentioned in the previous section, $R(x)$ for the $P(x)$ from the previous two steps should be combined with the test-set to be classified into the correct class.

4. TRk-CNN for glaucoma detection

In this section, we introduce the method of glaucoma detection based on the TRk-CNN. The fundus images we want to classify are labeled as normal, glaucoma suspect, and glaucoma eyes. Since glaucoma suspicious eyes can be seen as an intermediate stage between normal eyes and glaucoma eyes, better performance can be achieved by considering the inter-class relationship with TRk-CNN. The overall process of glaucoma detection is as follows. First, we perform pre-processing on fundus images. Then, the fundus images are augmented to perform primitive classification. Next, the ROI is generated from the predicted value and the weight of the sub-CNN model obtained as results of the primitive classification. Finally, the ROI is combined with the original fundus image to perform the final classification, and the aggregated predicted class is compared with the actual class. Figure 3 shows the overall process for classifying normal, glaucoma suspect, and glaucoma eyes with TRk-CNN.

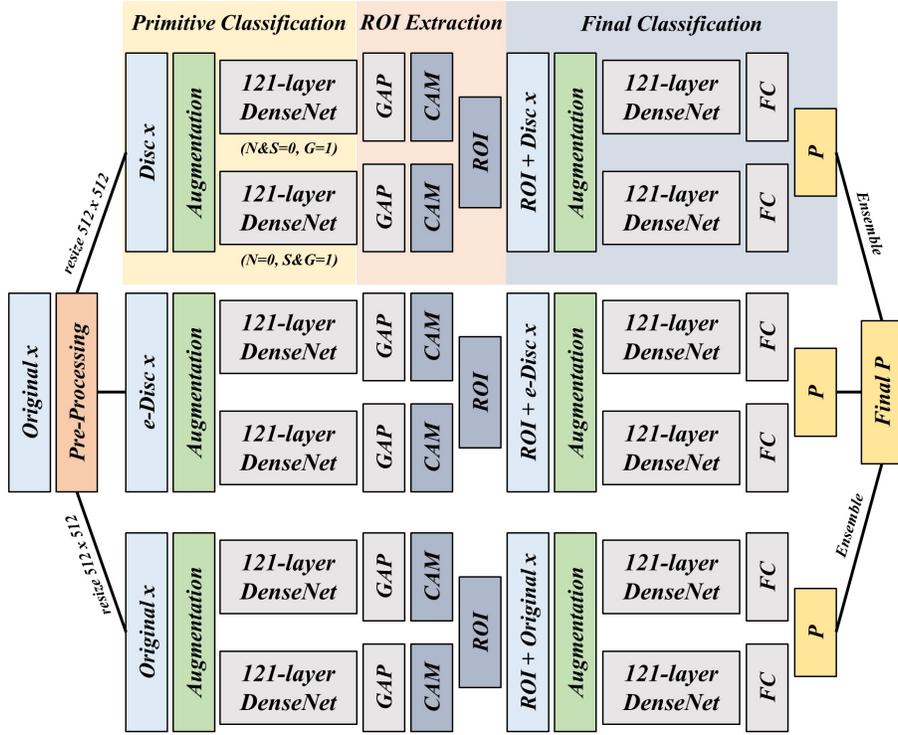


Figure 3: The overall overall process of TRk-CNN for glaucoma detection

4.1. Pre-processing

Although the resolution of the fundus image is very high, the area that plays an important role in the diagnosis of glaucoma is the disc/cup region. This is because cup-to-disc ratio (CDR) is one of the main criteria for discriminating glaucoma suspicious eyes. Therefore, in the pre-processing stage, we manually extract the disc/cup region of the fundus images. The optimized model will apply TRk-CNN models to the original fundus image, disc region image, and extended disc (e-disc) region image and then ensemble the results of the three models. The extended disc region is a region where the same range of pixels (t) is added to the top, bottom, left, and right sides of the disc region that we manually extracted. Therefore, the extended disc region can be regarded as the intermediate image between the disc region and the original image. There are

several previous studies that automatically segments the disc/cup region with machine learning approaches. However in this paper, we believe it is sufficient to draw it manually because the area we are interested in is a square box that contains disc/cup, not the exact pixel-by-pixel disc/cup region. And although our the evaluation results show that applying TRk-CNN to the disc region has the highest performance, it is not much different from the results of the other two images. Figure 4 shows images of the three different regions from the same fundus images obtained as a result of pre-processing.

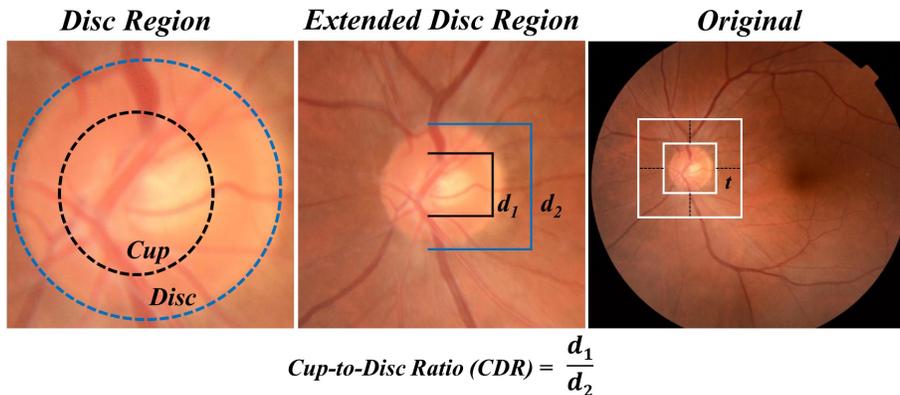


Figure 4: Three regions of pre-processing results

4.2. Data augmentation

Since our data set consists of about 1,000 fundus images, without augmentation the model will fall into overfitting problem shortly and it will be hard to expect reasonable performance for validation and test set. Fortunately, fundus images are not as varied as the general dataset such as ImageNet or Cifar10. In other words, normal, glaucoma suspicious, and glaucoma eyes are classified from fundus images with relatively similar class distribution, compared to a general image dataset with a heterogeneous class distribution. Therefore, even with a thousand number of images, the proper application of augmentation can yield acceptable classification accuracy. Our image augmentation policy is as follows. First, we zoom-in and zoom-out an image at a random ratio within $\pm 20\%$. And

the height and width of the image are shifted at a random ratio within $\pm 20\%$ of image size $h \times h$. Also, the image flips horizontally with a random probability, which has the effect of augmenting the right eye into the left eye and vice versa. Next, since the fundus image may have a different eye orientation depending on the angle of the screening, we rotate the image within $\pm 45^\circ$ at random rates. Finally, because the brightness of the fundus image is also different, the brightness is also changed within $\pm 40\%$ at random rates. Figure 5 shows images when each augmentation policy is applied to a single image at the maximum rate.

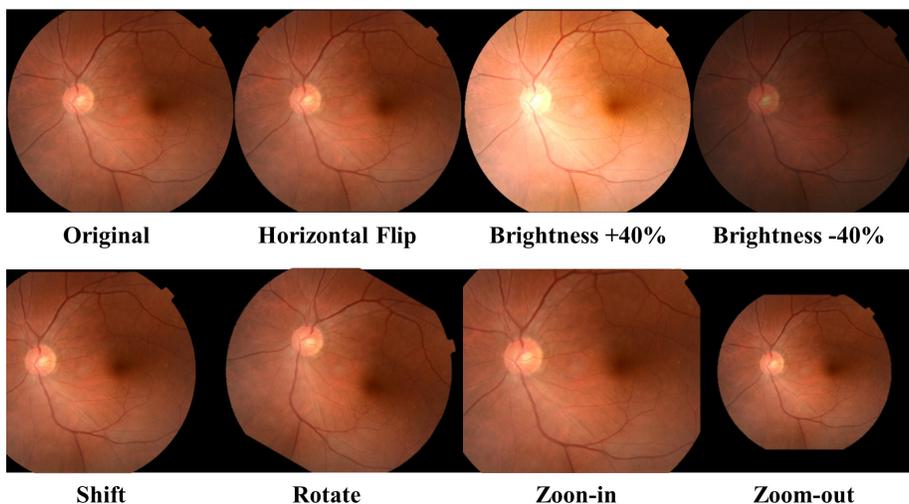


Figure 5: Example images from data augmentation

4.3. Primitive Classification

Starting from the primitive classification, the backbone structure of the CNN model to be used in the following steps is the DenseNet [6] with 121 number of layers. DenseNet extends ResNet’s [25] skip-connection concept and is characterized by a densely connected block. Dense connection encourages feature reuse and reduces the number of free parameters, thereby reducing overfitting in a relatively small train-set. Therefore, we judged that DenseNet as backbone CNN model is suitable for our fundus dataset which has smaller train-set than

general image dataset such as ImageNet [31] and Cifar10 [32]. However, the free parameters of DenseNet are still many to optimize with a thousand fundus images. So we started training by taking the weight of pre-trained 121-layer DenseNet in ImageNet and this method is generally called as transfer learning. Of course, because ImageNet images and fundus images are different types of images, we trained the entire weight from the beginning, unlike the general transfer learning method which trains only a few top layers.

Before starting training, the train-set is transformed according to the augmentation policy. And an input image is resized to $512 \times 512 \times 3$ in all regions including original, disc, and e-disc region. The reason for adjusting the image size to 512×512 , which is larger than common sizes 224×224 or 256×256 , is because the size of the original fundus image is very large, which has a minimum size of 3500×2500 . The resized and augmented train-set with the mini-batch size is now passed to the input of the 121-layer DenseNet to start training.

Since our fundus dataset has three classes, two sub-CNN models are required to perform Ranking-CNN in primitive classification. We labeled the actual class of normal eye as 0, glaucoma suspect eye as 1, and glaucoma eye as 2. Of course, the actual class of normal and glaucoma eye may be interchanged, but the existence of a glaucoma suspect eye between them should be maintained to perform Ranking-CNN. Let the 1st sub-CNN model as Sub^1 and the 2nd sub-CNN model as Sub^2 . Then input class of Sub^1 is 0 for normal eye, 1 for glaucoma suspect and glaucoma eyes. Likewise, in Sub^2 , normal and glaucoma suspect eyes become class 0, and glaucoma eye becomes class 1. After the input is passed to each sub-CNN model with the 121-layer DenseNet, the size of a final convolutional layer is $32 \times 32 \times 1024$. Applying the global average pooling results in a layer with a size of 1024, followed by a size 2 softmax layer for binary classification. The optimization parameters of the model will be explained more concretely in the final classification section. As a result, the weight of the model with a minimum loss for the validation set and the aggregated predicted class $P \in \{0,1,2\}$ for the input are passed to the ROI extraction step. The aggregated predicted class P can be obtained as $P = p_1 + p_2$, where $p_1 \in \{0,1\}$ is the

predicted class of Sub^1 and $p_2 \in \{0,1\}$ is the predicted class of Sub^2 for the input. Figure 6 shows the overall process of primitive classification step for glaucoma detection.

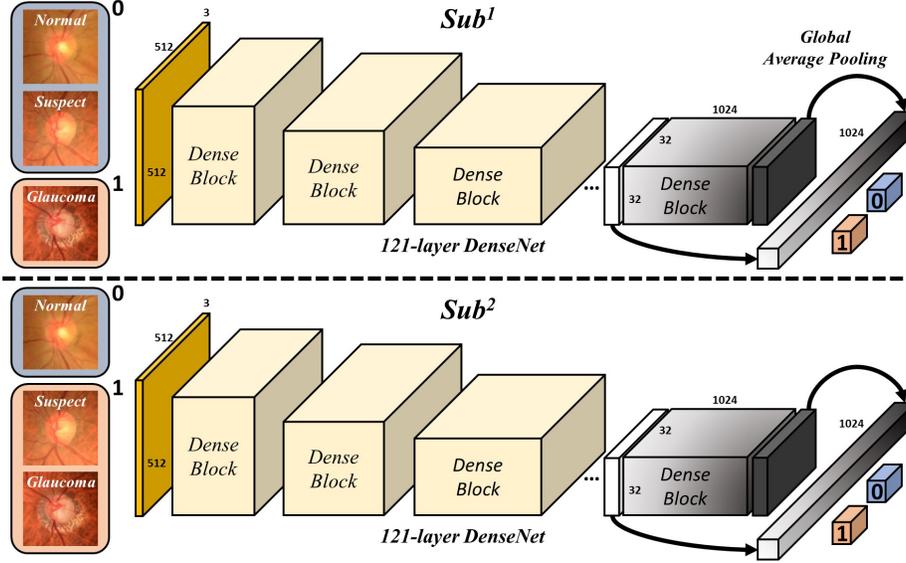


Figure 6: Primitive classification for glaucoma detection

4.4. ROI Extraction

The purpose of the ROI extraction step is to generate the region of interest R based on the P and model weight received earlier from the primitive classification. First, the Class Activation Maps for the binary classes of Sub^1 and Sub^2 for the given input are called Cam_0^1 , Cam_1^1 , Cam_0^2 , and Cam_1^2 , respectively. That is, Cam_0^1 is the CAM of Sub^1 as class 0 for the given input. Since the size of the input excluding the channel is 512 x 512 and the number of nodes in the GAP is 1024, we obtain the size 512 x 512 CAM by resizing the output by $512/\sqrt{1024}$ times, which is 16. Based on the equation 7 the ROI R for the

predicted class P can be expressed by the following equation.

$$R = \begin{cases} Cam_0^1, & P = 0 \\ Cam_0^1 + Cam_1^2, & P = 1 \\ Cam_1^2, & P = 2 \end{cases} \quad (8)$$

Finally we perform z-score normalization before passing the generated R to the final classification step. As we mentioned in section 2.3.2, we have also evaluated $R = Cam_1^1 + Cam_0^2$ when the P is 1 in Result section to compare the performance difference. Figure 7 shows the overall process of ROI extraction step for glaucoma detection.

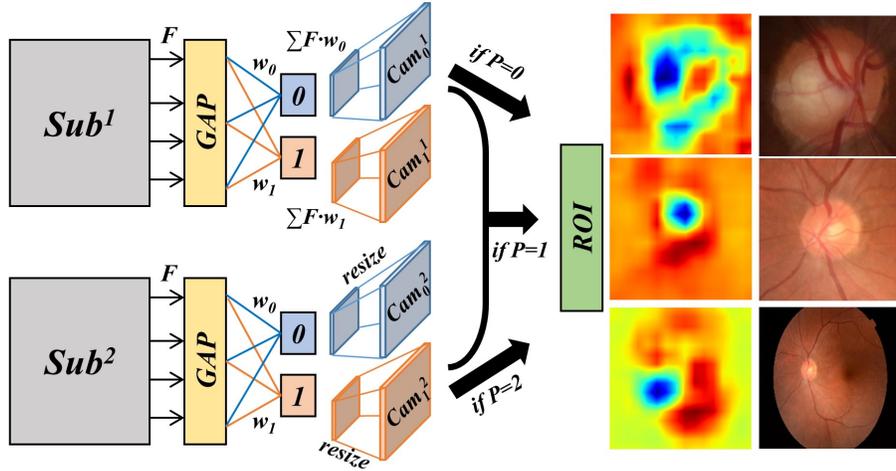


Figure 7: ROI extraction for glaucoma detection

4.5. Final Classification

The final classification begins by concatenating the input with R from the ROI extraction step. Since the train-set has a size of $512 \times 512 \times 3$ and R has a size of 512×512 , if we concatenate the two, the size of the new train-set is $512 \times 512 \times 4$. The image augmentation policy is the same as the primitive classification, but for brightness policy, it should only be applied to the original train-set of the new train-set. The reason is that the last channel, which is

R , should be transformed with the rotation, translation, and zooming policies because R represents the spatial characteristics of the given input, but it is not affected by brightness. The newly generated input is classified through Ranking-CNN based on 121-layer DenseNet as well as the primitive classification step. One difference now is that strict training is available by adding fully-connected layers after the GAP layer. From here, we will explain a specific description of the parameters applied to the final classification.

4.5.1. Loss Function

In general, categorical cross-entropy loss ($CELoss$) is used for the loss function in classification problems, but in our experience, using categorical cross-entropy (CE) alone increases the gap between minimum validation loss and maximum validation accuracy (Acc). As will be described later in the evaluation, intuitively, the gap between the softmax output vector and the predicted class vector occurs when argmax function is applied. Therefore, we use a loss function that combines both categorical cross-entropy loss and average accuracy. When we use categorical cross-entropy loss alone in the glaucoma detection problem, we confirmed that it converges at a validation loss of about 0.1 and that the validation accuracy converges to around 0.9. However, since the fluctuation of cross entropy per epoch is greater than the fluctuation of accuracy, we needed to adjust the scale of categorical cross-entropy loss from the final loss. As a result, the categorical cross-entropy loss with accuracy ($CEALoss$) for input $x = \{x_1, x_2, \dots, x_b\}$ with mini-batch size b is as follows.

$$\begin{aligned}
 CE(x) &= -\sum_{i=1}^c \ln s_i(x) \\
 CELoss &= \frac{1}{b} \sum_{j=1}^b CE(x_j) \\
 Acc &= \frac{1}{b} \sum_{j=1}^b y(x_j) \cdot p(x_j) \\
 CEALoss &= 1 + \alpha CELoss - Acc
 \end{aligned} \tag{9}$$

where c is the number of classes, $s_i(x)$ is the softmax output value for class $i \in \{1, 2, \dots, c\}$, $y(x)$ is the one hot encoded vector represents true class for input x , $p(x)$ is the one hot encoded vector represents predicted class for input x , and α is coefficient for adjusting the scale of $CEALoss$ which set to 0.1 in this paper. However, α can be intuitively changed depending on the classification problem. We compared the performance of the $CEALoss$ and the $CELoss$ in the evaluation, and as a result, the performance of the $CEALoss$ was better.

4.5.2. Activation and Optimizer Functions

The role of the activation function is to define the output value of kernel weights in the model. In modern CNN models, nonlinear activation is widely used, including rectified linear units (ReLU) [33], leakage rectified linear units (LReLU) [34], and exponential linear units (ELU) [35]. As we experimentally confirmed, we have applied the most commonly used ReLU because the three activation functions were not significantly different in performance.

The role of the optimizer function is to minimize the loss function through the stochastic gradient descent approach with learning rate. There are several well-known optimizer functions such as Adam [36], Adagrad [37], and Adadelta [38]. In general, Adam function converges faster than other functions. Therefore, we also used Adam for optimizer function and the initial learning rate was set to 0.0001. In addition, we reduced the learning rate by half if the validation loss does not improve for the last 10 epochs.

4.5.3. Regularization

Regularization is a method to reduce overfitting during the training phase. Overfitting is a problem especially when the size of the train-set is small and the free parameter of the model is large like our glaucoma detection problem. Image augmentation is also a regularization technique, which is not directly applied to the model, so it is described after the pre-processing section. Typical regularization methods are using L1 and L2 norm, however, it is common to apply Dropout [39] and Batch Normalization [40] in recent CNN models. In deep

learning, when a layer is deepened, a small parameter change in the previous layer can have a large influence on the input distribution of the later layer. This phenomenon is referred to as internal co-variate shift. Batch normalization has been proposed to reduce this internal co-variate shift, and the mean and variance of input batches are calculated, normalized, and then scaled and shifted. The location of Batch Normalization is usually applied just before the activation function and after the convolution layer. The 121-layer DenseNet [6] we used uses Batch Normalization by default, and it is also applied to the last two fully-connected layers.

Another popular regularization technique is Dropout, which stochastically participates in nodes in the same layer, reducing dependency between layers to prevent overfitting. In the training phase, Dropout intentionally excludes some networks, so the model can achieve the voting effect through a combination of partial models. In recent, however, only Batch Normalization is applied to the convolution layers, and Dropout has been selectively added to the fully-connected layer. We also apply Dropout of 0.5 probability to only the last two fully-connected layers.

Finally, an ensemble of several models can be regarded as regularization from the viewpoint of machine learning. In this paper, we use the ensemble method of voting the three prediction results of the trained models from different image regions including original, disc, and e-disc regions. Figure 8 shows the concrete process of the final classification together with optimization parameters.

5. Results

5.1. Data acquisition

This study included 1022 fundus images from 301 consecutive patients (582 eyes) who underwent fundus imaging with a non-mydratic fundus camera (TRC-NW8; Topcon, Oakland, NJ, USA), at Korea University Ansan Hospital between January 2016 and August 2017. During the study period, patient electronic medical records and fundus imaging were reviewed to determine the pres-

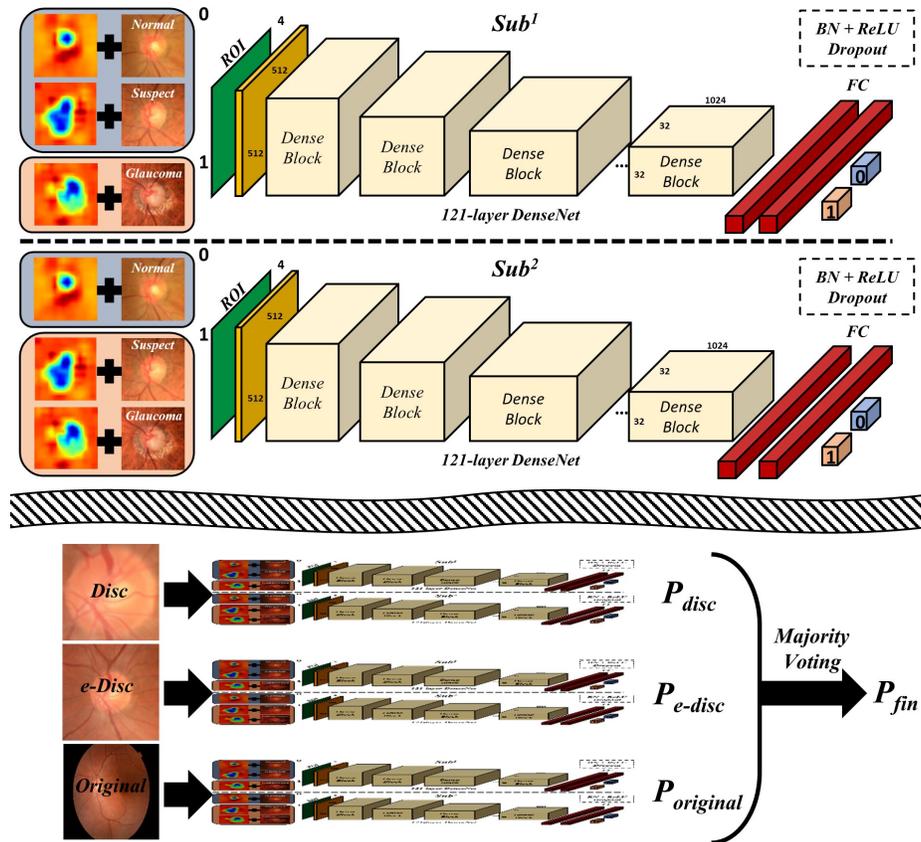


Figure 8: Final classification for glaucoma detection

ence of glaucoma by the glaucoma specialist. Based on fundus imaging and electronic medical records, 1022 fundus images were divided into three categories; normal, glaucoma suspect (suspicious), and glaucoma. Fundus images were classified as a glaucoma suspect when a vertical cup-to-disc ratio (CDR) is greater than 0.7 or the peripapillary retinal nerve fiber layer (RNFL) has a characteristic thinning (the presence of RNFL defect) but there is no glaucomatous visual field loss. Fundus images were classified as glaucoma when there is a RNFL defect or visual field loss with a corresponding glaucomatous optic disc change. When fundus images do not correspond to the two mentioned categories above, they were classified as normal. This study adhered to the Declaration of Helsinki

and approval for retrospective review of clinical records was obtained from Korea University Ansan Hospital Institutional Review Board (2017AS0036). The patient information was completely anonymized and de-identified prior to analysis. Of the 301 patients, 138 (45.8%) were men and 163 were women. The mean age (\pm SD) was 59.7 (\pm 15.4) years (range, 19-92 years). There were 291 right eyes (50.0%) and 291 left eyes. However, 992 images were used as the fundus image dataset of this study because 30 images of the wrong file format were excluded. Of the 922 fundus images, 403 (40.6%) were normal, 208 (21.0%) were glaucoma suspect, and 381 (38.4%) were glaucoma eyes. From the total number of 992 fundus images, 793 images (80%) were randomly split into train-set and 199 images into (20%) test-set with the similar class distribution. 199 test-set images consisted of 75 normal images, 41 glaucoma suspect images, and 83 glaucoma images. Validation-set consists of 119 images which correspond to 15% of the train-set, and also the class distribution is similar. As a result, 674 train-set images consisted of 272 normal images, 142 glaucoma suspect images, and 260 glaucoma images. Likewise, of the 119 validation-set images, 56 images are normal, 25 images are glaucoma suspect, and 38 images are glaucoma eyes.

5.2. Evaluation setup

The software and hardware environment for the evaluation are as follows. We tested on a 64GB server with two NVIDIA Titan X GPUs and an Intel Core i7-6700K CPU. The operating system is Ubuntu 16.04, and the development of the CNN model uses Python-based machine learning libraries including Keras [41], Scikit-learn [42], and TensorFlow [43].

We conducted the evaluation from two perspectives. The first is to compare TRk-CNN with Ranking-CNN (Rk-CNN) and multi-class CNN (MC-CNN) under the same conditions. Here, the same condition means that the region of fundus images and the structure of the model are the same. First, the fundus image with a disc region is only used because a disc region shows the best performance among the three regions. Experimental results in three regions are shown by applying TRk-CNN. The same augmentation policy was then applied

to train-set images, which is described in detail in the previous section. Rk-CNN and MC-CNN have a 121-layer DenseNet as a basic structure, and two fully-connected layers are added after the last convolutional Layer. In other words, the structure of Rk-CNN and MC-CNN is the same as that of TRk-CNN’s final classification step which is shown in Figure 9. The only difference is that MC-CNN classifies normal, glaucoma suspect, and glaucoma at once, so the last prediction layer consists of three nodes. Figure 9 outlines the structural differences between TRk-CNN, Rk-CNN, and MC-CNN. Finally, these three models are trained for 100 epochs with the Adam optimizer function with an initial learning rate set to 0.0001, and the learning rate is halved if there is no improvement in validation loss over 10 epochs. The loss function used for comparison is the $CELoss$, not the $CEALoss$. $CEALoss$ is used in the optimal TRk-CNN model for glaucoma detection.

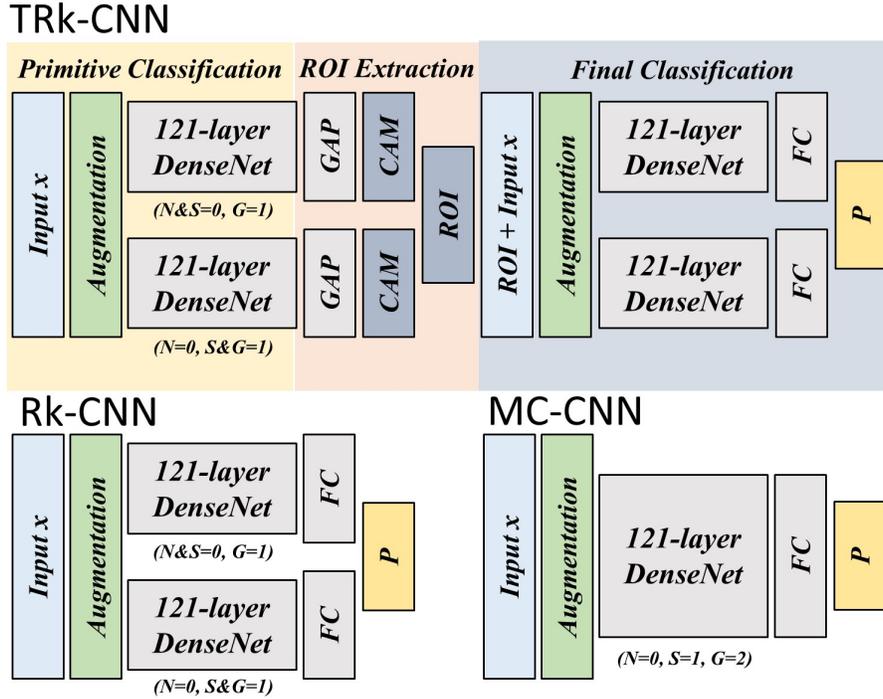


Figure 9: Structural differences between TRk-CNN, Rk-CNN, and MC-CNN

The second is to evaluate the different TRk-CNN models, which is optimized for glaucoma detection. We applied the *CEALoss* described in the previous section, and the final prediction was obtained by an ensemble of the three models trained in the original, disc, and e-disc regions of fundus image. The results of the binary classification of each sub-CNN model from the optimized TRk-CNN model are also compared with the performance of the existing literature. Other training parameters and model structure are the same as the three models described above.

5.3. Evaluation metrics

The evaluation of the glaucoma classification was based on the following four metrics: average accuracy (Acc), specificity (Sp), sensitivity (Se^S for glaucoma suspect, Se^G for glaucoma), precision (Pr^S , Pr^G), and F1 score ($F1^S$, $F1^G$). Average accuracy means a correctly predicted percentage of the total data. Specificity, also known as the true negative rate, measures the percentage of negatives that are correctly identified as normal. Sensitivity, also known as the true positive rate or recall, measures the percentage of positives that are correctly identified as glaucoma suspect or glaucoma. Precision measures the percentage of positives that are predicted as glaucoma suspect or glaucoma. F1 score is a harmonic mean of sensitivity and precision. These metrics are defined with the following four terminologies.

- True Positive(TP): The number of fundus images correctly identified as glaucoma suspect or glaucoma.
- False Positive(FP): The number of fundus images incorrectly identified as glaucoma suspect or glaucoma.
- True Negative(TN): The number of fundus images correctly identified as normal
- False Negative(FN): The number of fundus images incorrectly identified as normal

$$\begin{aligned}
Accuracy(Acc) &= \frac{TP + TN}{TP + TN + FP + FN} \times 100(\%) \\
Specificity(Sp) &= \frac{TN}{FP + TN} \times 100(\%) \\
Sensitivity(Se) &= \frac{TP}{TP + FN} \times 100(\%) \\
Precision(Pr) &= \frac{TP}{TP + FP} \times 100(\%) \\
F1 - score(F1) &= \frac{2 \times Pr \times Se}{Pr + Se}
\end{aligned} \tag{10}$$

5.4. Evaluation results of TRk-CNN, Rk-CNN, and MC-CNN

Table 1 shows the evaluation results of TRk-CNN, Rk-CNN, and MC-CNN models experimented under the same condition explained in the previous section.

Method	Acc(%)	Sp(%)	Se ^S (%)	Se ^G (%)	Pr ^S (%)	Pr ^G (%)	F1 ^S (%)	F1 ^G (%)
TRk-CNN	88.94	89.33	85.37	90.36	74.47	94.94	79.55	92.59
Rk-CNN	84.92	85.33	85.37	84.34	60.34	100.0	70.71	91.50
MC-CNN	83.42	85.33	68.29	89.16	75.68	85.06	71.79	87.06
MC-CNN ¹	91.46	89.33	92.74		93.50		93.12	
MC-CNN ²	92.96	97.41	86.75		96.00		91.14	

Table 1: Comparison results between TRk-CNN, Rk-CNN, and MC-CNN

Overall, TRk-CNN showed higher performance in all metrics except precision. In terms of accuracy, TRk-CNN achieved 88.94%, which is 4.02% higher than Rk-CNN and 5.52% higher than MC-CNN. From the specificity perspective, TRk-CNN was the highest at 89.33%, which is 4% higher than Rk-CNN and MC-CNN. The sensitivities of glaucoma suspect for TRk-CNN and Rk-CNN were 85.37%, which is 17.08% higher than MC-CNN. The precision of glaucoma suspect for TRk-CNN achieved 74.47%, which is 14.13% higher than Rk-CNN and 1.21% lower than MC-CNN. Since sensitivity and precision have a trade-off relation, it is better to consider F1 score together. In terms of F1-score for glaucoma suspect, TRk-CNN was the highest at 79.55% which is 7.84% higher than Rk-CNN and 6.76% higher than MC-CNN. The sensitivity of glau-

coma for TRk-CNN was 90.36% while Rk-CNN was 84.34% and MC-CNN was 89.16%. The precision of glaucoma for TRk-CNN achieved 94.94% while Rk-CNN was 100.00% and MC-CNN was 85.06%. Finally, F1-score of glaucoma for TRk-CNN was the highest at 92.59% which is 1.09% higher than Rk-CNN and 5.53% higher than MC-CNN.

The reason why MC-CNN has the lowest overall performance is that MC-CNN assumes three classes as independent classes without considering the inter-class relationship. For a more precise description, we evaluated MC-CNN to perform binary classification. MC-CNN¹ classifies normal eye as 0, glaucoma suspicion and glaucoma eyes as 1. Likewise, MC-CNN² classifies normal and glaucoma suspect eyes as 0, glaucoma eye as 1. From Table 1, we can observe that the overall performance is improved despite the same structure as MC-CNN. This shows that our classification problem is a difficult problem compared to the binary classification problem that classifies the normal eyes and glaucoma eyes in the previous studies. We will show the results in comparison with the previous studies in the following section.

Figures 10 show the training loss and validation accuracy of the sub-CNN models of Rk-CNN and TRk-CNN during the 100 epochs, along with those of MC-CNN. Since the number of classes to classify is different, there is a limit to directly comparing MC-CNN with sub-CNN models of Rk-CNN and TRk-CNN, in terms of validation accuracy. However, looking at the tendency of the graphs, TRk-CNN's training loss and validation accuracy converge from earlier epochs than the Rk-CNN and MC-CNN. In other words, by exchanging the ROI extracted from different models, additional information on the input is obtained, so that training with lower loss becomes possible. This explains why TRk-CNN performs better than Rk-CNN considering that the total error of Rk-CNN is bound to the max error of sub-model. As a result, the validation accuracy of the sub-CNN model in TRk-CNN is higher than that of Rk-CNN.

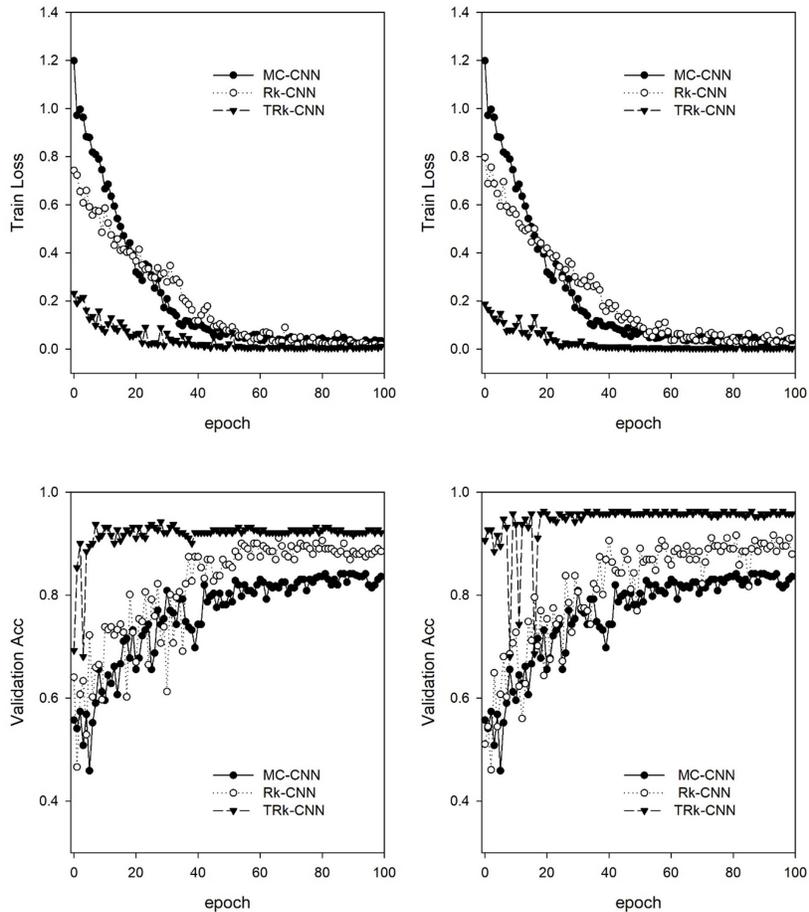


Figure 10: Training loss and validation accuracy of TRk-CNN, RK-CNN, and MC-CNN

5.5. Evaluation results of optimized TRk-CNN for glaucoma detection

Table 2 shows the results of optimized TRk-CNN models trained in several different conditions. DISC, EDISC, and ORIGINAL represent three different regions of the fundus image, all of which were trained using *CEALoss*. ENSEMBLE is the result of majority voting on the predicted classes of three models, and if there is no dominant class it follows the result of DISC model. DISC¹ and DISC² are models for comparison, and DISC¹ shows the result of training using

$CELoss$ instead of $CEALoss$ in disc region. $DISC^2$ is the result of training glaucoma suspect’s ROI with $Cam_1^1 + Cam_0^2$ instead of $Cam_0^1 + Cam_1^2$. Since the loss of $DISC^2$ uses $CEALoss$, only the difference of performance according to ROI is compared.

Method	$Acc(\%)$	$Sp(\%)$	$Se^S(\%)$	$Se^G(\%)$	$Pr^S(\%)$	$Pr^G(\%)$	$F1^S(\%)$	$F1^G(\%)$
ENSEMBLE	92.96	93.33	95.12	91.57	81.25	98.70	87.64	95.00
DISC	91.46	89.33	90.24	93.98	78.72	96.30	84.09	95.12
EDISC	88.94	90.67	90.24	86.75	68.52	98.63	77.89	92.31
ORIGINAL	90.45	92.00	92.68	87.95	77.55	97.33	84.44	92.41
$DISC^1$	88.94	89.33	85.37	90.36	74.47	94.94	79.55	92.59
$DISC^2$	86.43	86.67	78.05	90.36	86.49	86.21	82.05	88.24

Table 2: Evaluation results of TRk-CNN models for glaucoma detection

Since DISC have the highest performance among the models studied in the three regions, we used the disc region in the comparison of TRk-CNN, Rk-CNN, and MC-CNN. The best overall performance was the ENSEMBLE model which achieved the highest results for all metrics except sensitivity and F1-score for glaucoma. In terms of accuracy, ENSEMBLE achieved 92.96%, which is 1.50% higher than DISC, 4.02% higher than EDISC, and 2.51% higher than ORIGINAL. From the specificity perspective, ENSEMBLE was the highest at 93.33%, which is 4% higher than DISC, 2.66% higher than EDISC, and 1.33% higher than MC-CNN. The sensitivities of glaucoma suspect for ENSEMBLE was 95.12%, which is 4.88% higher than both DISC and EDISC, and 2.44% higher than ORIGINAL. The precision of glaucoma suspect for ENSEMBLE achieved 81.25%, which is 2.53% higher than DISC, 12.73% higher than EDISC, and 3.70% higher than ORIGINAL. Considering the trade-off between sensitivity and precision, F1-score for glaucoma suspect in ENSEMBLE was the highest at 87.64% which is 3.55% higher than DISC, 9.75% higher than EDISC, and 3.20% higher than ORIGINAL. The sensitivity of glaucoma for ENSEMBLE was 91.57% while DISC was 93.98%, EDISC was 86.75%, and ORIGINAL was 87.95%. The precision of glaucoma for TRk-CNN achieved 98.70% while 96.30%

for DISC, 98.63% for EDISC, and 97.33% for ORIGINAL. Finally, F1-score of glaucoma for DISC was the highest at 95.12% which is 0.12% higher than ENSEMBLE, 2.81% higher than EDISC, and 2.71% higher than ORIGINAL. The results show that referring to the disc region is the best performance for specifying glaucoma. However, in the case of detecting normal and glaucoma suspect eyes, it is better to refer to a wider area, and as a result, the ENSEMBLE model that combines all of these is the best.

From the results of DISC and DISC¹, using *CEALoss* instead of *CELoss* showed higher performance in all metrics. This means that lower *CELoss* does not necessarily result in higher accuracy as we explained earlier. Also, if we use a metric other than accuracy as an evaluation, many variations are possible. For example, since the Dice Similarity Coefficient score (*DCS*) is the main metric for the segmentation problem, the combined loss of *CELoss* and *DCS* may show better performance.

The results of DISC and DISC² showed that the performance of DISC was higher in all the indicators except precision for glaucoma suspect. However, the F1-score for glaucoma suspect was higher on the DISC, so overall it was better to use ROI as $Cam_0^1 + Cam_1^2$. As described in earlier section, defining ROI as $Cam_0^1 + Cam_1^2$ is considered to contain information that is likely to be the opposite of the prediction in each sub-model. In other words, to output the glaucoma suspect class in the primitive classification step of TRk-CNN, the probability of predicting 1 in Sub^1 and 0 in Sub^2 is higher than in the opposite case. Therefore, it is expected that Cam_0^1 and Cam_1^2 are highly contrary to predicted class information, and combining these two can transfer more features to the final classification step.

Table 3 compares the results of previous studies with the results of our proposed TRk-CNN model for glaucoma detection. However, since previous studies were binary classifications that classify normal and glaucoma instead of three classes, we included the binary classification results of the proposed model. Proposed¹ classifies normal eye as 0, glaucoma suspicious and glaucoma eyes as 1. In other words, Proposed¹ is the ensemble of Sub^1 models from DISC,

EDISC, and ORIGINAL.

Method	$Acc(\%)$	$Sp(\%)$	$Se^S(\%)$	$Se^G(\%)$	$Pr^S(\%)$	$Pr^G(\%)$	$F1^S(\%)$	$F1^G(\%)$
Proposed	92.96	93.33	95.12	91.57	81.25	98.70	87.64	95.00
Rk-CNN	84.92	85.33	85.37	84.34	60.34	100.0	70.71	91.50
MC-CNN	83.42	85.33	68.29	89.16	75.68	85.06	71.79	87.06

	Year	Data	AUC	$Acc(\%)$	$Sp(\%)$	$Se(\%)$	$Pr(\%)$	$F1(\%)$
Proposed ¹	2019	992	0.974	95.48	93.33	96.77	96.00	96.39
Li [26]	2018	48116	0.986	-	92	95.6	-	-
Fu [23]	2018	SCES	0.918	-	-	-	-	-
		SINDI	0.817	-	-	-	-	-
Li [20]	2016	ORIGA	0.838	-	-	-	-	-
Chen [16]	2015	ORIGA	0.831	-	-	-	-	-
		SCES	0.887	-	-	-	-	-
Dua [15]	2012	60	-	93.33	-	-	-	-
Acharya [14]	2011	60	-	91.7	-	-	-	-
Bock [13]	2010	575	0.88	-	85	73	-	-
Nayak [12]	2009	61	-	-	80	100	-	-

Table 3: Result table including comparison with results of previous studies

Since TRk-CNN is a model for considering the inter-class relationship, it can be seen that there is no significant difference from using MC-CNN in case of binary classification. This can be seen from the fact that the work of Li [26] and the performance of Proposed¹ do not differ greatly. However, when performing three-class classification, the performance difference between MC-CNN and Proposed is large, because the classes of normal, glaucoma suspect, and glaucoma have a high relation with each other. Therefore, when multi-class classification is performed considering the inter-class relationship, using TRk-CNN can be expected to perform better than the multi-class classification approach.

Figures 11 show the training loss and validation accuracy of 1st and 2nd sub-CNN models of DISC, EDISC, and ORIGINAL, respectively. One notable difference is that the overall validation accuracy of the 1st sub-CNN model was the highest in EDISC, but the results in test-set were the highest in DISC. This

implies that the best performance in one sub-model may not necessarily be the best for the aggregated result.

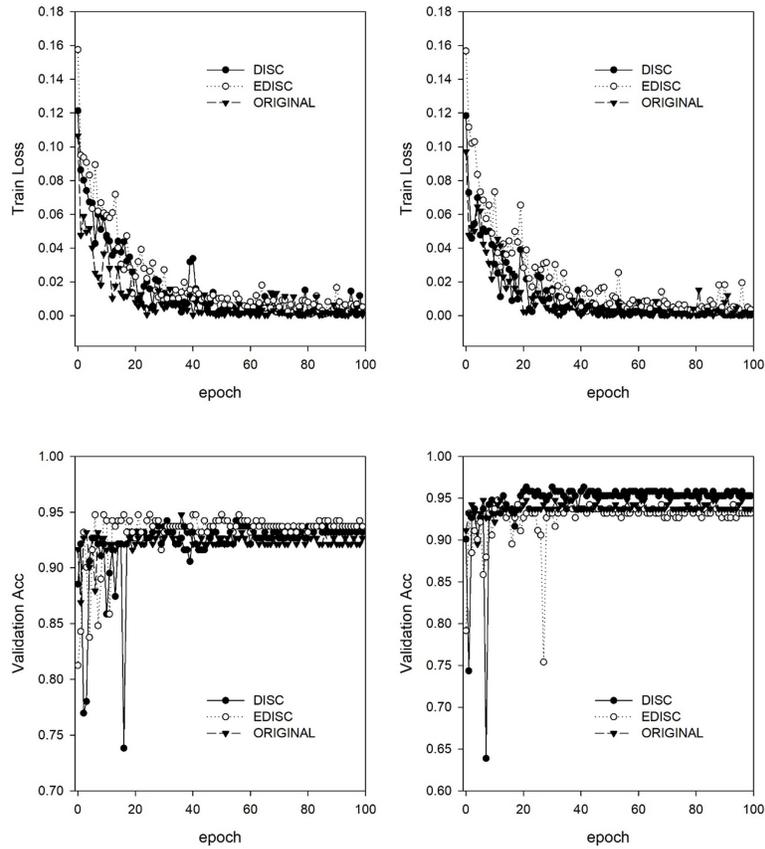


Figure 11: Training loss and validation accuracy of DISC, EDISC, and ORIGINAL

Figures 12 show how the activation of each convolutional layer is visualized where the input image is the three regions of the same fundus image. The top left image in each figure represents disc, e-disc, and original region for the same fundus image.

The six images on the bottom left of each figure are visualizations of the activation in the pooling layer of each model. The six images show the deepening of the model from top to bottom, highlighting the retinal blood vessel and

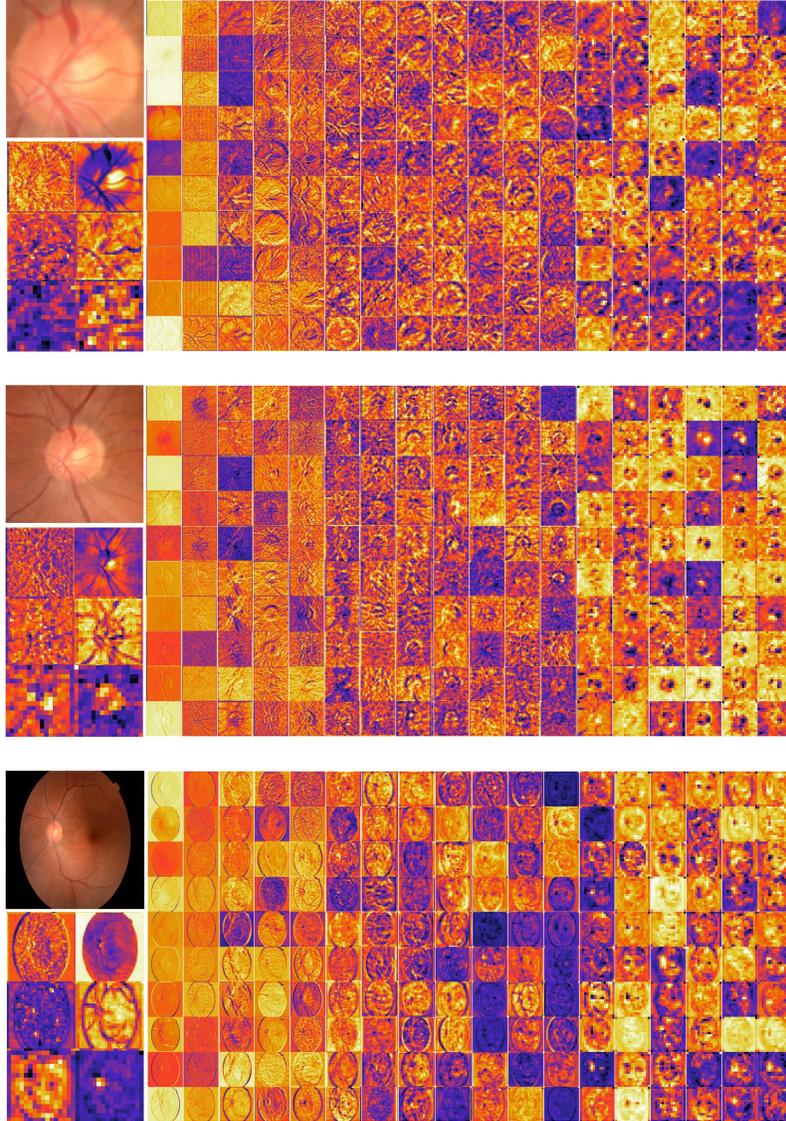


Figure 12: Visualization of the convolution layer in the DISC

disc/cup regions. This can be seen more clearly in DISC and EDISC. Although we manually draw the region box contains only the disc region, we can observe that the model automatically emphasizes the cup region. Other small patch images are from the left to the right in the direction of deepening the model, all

visualizing the activation of the convolutional layer. The patches in the same column represent the first 10 filters of the convolutional layer. As shown in the figures, the early part of the convolutional layer extracts low-level feature such as image outline and contrast. As the model deepens, we can see that high-level features are extracted. One peculiar point is that in the case of EDISC and ORIGINAL, the disc region is still emphasized even though the depth of the model is deep enough.

6. Conclusion

Our proposed TRk-CNN is a method that can be effectively applied when the classes of images to be classified show a high correlation with each other. The multi-class classification method based on the softmax function, which is generally used, is not effective in this case because the inter-class relationship is ignored. Although there is a Ranking-CNN that takes into account the ordinal classes, it cannot reflect the inter-class relationship to the final prediction. TRk-CNN, on the other hand, combines the weights of the primitive classification model to reflect the inter-class information to the final classification phase. Through extensive experiments, we show that TRk-CNN is superior to both the multi-class classification method and Ranking-CNN method.

We evaluated TRk-CNN in glaucoma image dataset that was collected and labeled from Korea University Medical Center. Glaucoma dataset was labeled into three classes: normal, glaucoma suspect, and glaucoma eyes. Based on the literature we surveyed, this study is the first to classify three status of glaucoma fundus image dataset into three different classes. We compared the evaluation results of TRk-CNN with multi-class CNN (MC-CNN) and Ranking-CNN (Rk-CNN) using the DenseNet as the backbone CNN model. As a result, TRk-CNN achieved an average accuracy of 92.96%, specificity of 93.33%, sensitivity for glaucoma suspect of 95.12% and sensitivity for glaucoma of 93.98%. Based on average accuracy, TRk-CNN is 8.04% and 9.54% higher than Rk-CNN and MC-CNN and surprisingly 26.83% higher for sensitivity for suspicious than multi-

class CNN.

Our TRk-CNN is expected to be effectively applied to the medical image classification problem where the disease state is continuous and increases in the positive class direction. Therefore, we will apply TRk-CNN to medical images with the above characteristics in future work.

Acknowledgments

This research was supported by International Research & Development Program of the National Research Foundation of Korea(NRF) funded by the Ministry of Science, ICT&Future Planning of Korea(2016K1A3A7A03952054) and by a Korea University Grant (K1625491, K1722121, and K1811051) and by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Science and ICT (2018R1C1B6002794). The second funding source had no role in the design or conduct of this research.

References

References

- [1] L. H. Sobin, M. K. Gospodarowicz, C. Wittekind, TNM classification of malignant tumours, John Wiley & Sons, 2011.
- [2] L. T. Chylack, J. K. Wolfe, D. M. Singer, M. C. Leske, M. A. Bullimore, I. L. Bailey, J. Friend, D. McCarthy, S.-Y. Wu, The lens opacities classification system iii, *Archives of ophthalmology* 111 (6) (1993) 831–836.
- [3] R. N. Weinreb, P. T. Khaw, Primary open-angle glaucoma, *The Lancet* 363 (9422) (2004) 1711–1720.
- [4] S. Chen, C. Zhang, M. Dong, Deep age estimation: From classification to ranking, *IEEE Transactions on Multimedia* 20 (8) (2018) 2209–2222.

- [5] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, A. Torralba, Learning deep features for discriminative localization, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 2921–2929.
- [6] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely connected convolutional networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 4700–4708.
- [7] T. J. Jun, D. Kim, H. M. Nguyen, D. Kim, Y. Eom, 2sranking-cnn: A 2-stage ranking-cnn for diagnosis of glaucoma from fundus images using cam-extracted roi as an intermediate input, arXiv preprint arXiv:1805.05727.
- [8] G. A. Stevens, R. A. White, S. R. Flaxman, H. Price, J. B. Jonas, J. Keeffe, J. Leasher, K. Naidoo, K. Pesudovs, S. Resnikoff, et al., Global prevalence of vision impairment and blindness: magnitude and temporal trends, 1990–2010, *Ophthalmology* 120 (12) (2013) 2377–2384.
- [9] R. R. Bourne, G. A. Stevens, R. A. White, J. L. Smith, S. R. Flaxman, H. Price, J. B. Jonas, J. Keeffe, J. Leasher, K. Naidoo, et al., Causes of vision loss worldwide, 1990–2010: a systematic analysis, *The lancet global health* 1 (6) (2013) e339–e349.
- [10] S. Kingman, Glaucoma is second leading cause of blindness globally, *Bulletin of the World Health Organization* 82 (2004) 887–888.
- [11] E. D. P. R. Group, et al., Prevalence of open-angle glaucoma among adults in the united states, *Archives of ophthalmology* 122 (4) (2004) 532.
- [12] J. Nayak, R. Acharya, P. S. Bhat, N. Shetty, T.-C. Lim, Automated diagnosis of glaucoma using digital fundus images, *Journal of medical systems* 33 (5) (2009) 337.
- [13] R. Bock, J. Meier, L. G. Nyúl, J. Hornegger, G. Michelson, Glaucoma risk index: automated glaucoma detection from color fundus images, *Medical image analysis* 14 (3) (2010) 471–481.

- [14] U. R. Acharya, S. Dua, X. Du, C. K. Chua, et al., Automated diagnosis of glaucoma using texture and higher order spectra features, *IEEE Transactions on information technology in biomedicine* 15 (3) (2011) 449–455.
- [15] S. Dua, U. R. Acharya, P. Chowriappa, S. V. Sree, Wavelet-based energy features for glaucomatous image classification, *Ieee transactions on information technology in biomedicine* 16 (1) (2012) 80–87.
- [16] X. Chen, Y. Xu, D. W. K. Wong, T. Y. Wong, J. Liu, Glaucoma detection based on deep convolutional neural network, in: 2015 37th annual international conference of the IEEE engineering in medicine and biology society (EMBC), IEEE, 2015, pp. 715–718.
- [17] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [18] Z. Zhang, F. S. Yin, J. Liu, W. K. Wong, N. M. Tan, B. H. Lee, J. Cheng, T. Y. Wong, Origa-light: An online retinal fundus image database for glaucoma analysis and research, in: 2010 Annual International Conference of the IEEE Engineering in Medicine and Biology, IEEE, 2010, pp. 3065–3068.
- [19] C. C. Sng, L.-L. Foo, C.-Y. Cheng, J. C. Allen Jr, M. He, G. Krishnaswamy, M. E. Nongpiur, D. S. Friedman, T. Y. Wong, T. Aung, Determinants of anterior chamber depth: the singapore chinese eye study, *Ophthalmology* 119 (6) (2012) 1143–1150.
- [20] A. Li, J. Cheng, D. W. K. Wong, J. Liu, Integrating holistic and local deep features for glaucoma classification, in: 2016 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), IEEE, 2016, pp. 1328–1331.
- [21] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich, Going deeper with convolutions, in:

Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 1–9.

- [22] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.
- [23] H. Fu, J. Cheng, Y. Xu, C. Zhang, D. W. K. Wong, J. Liu, X. Cao, Disc-aware ensemble network for glaucoma screening from fundus image, *IEEE transactions on medical imaging* 37 (11) (2018) 2493–2501.
- [24] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: *International Conference on Medical image computing and computer-assisted intervention*, Springer, 2015, pp. 234–241.
- [25] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [26] Z. Li, Y. He, S. Keel, W. Meng, R. T. Chang, M. He, Efficacy of a deep learning system for detecting glaucomatous optic neuropathy based on color fundus photographs, *Ophthalmology* 125 (8) (2018) 1199–1206.
- [27] R. Herbrich, Large margin rank boundaries for ordinal regression, *Advances in large margin classifiers* (2000) 115–132.
- [28] Y. Freund, R. Iyer, R. E. Schapire, Y. Singer, An efficient boosting algorithm for combining preferences, *Journal of machine learning research* 4 (Nov) (2003) 933–969.
- [29] P. Yang, L. Zhong, D. Metaxas, Ranking model for facial age estimation, in: *2010 20th International Conference on Pattern Recognition*, IEEE, 2010, pp. 3404–3407.
- [30] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, G. N. Hullender, Learning to rank using gradient descent, in: *Proceedings of the*

22nd International Conference on Machine learning (ICML-05), 2005, pp. 89–96.

- [31] O. Russakovsky, J. Deng, H. Su, J. 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) (2015) 211–252.
- [32] A. Krizhevsky, G. Hinton, Learning multiple layers of features from tiny images, Tech. rep., Citeseer (2009).
- [33] V. Nair, G. E. Hinton, Rectified linear units improve restricted boltzmann machines, in: *Proceedings of the 27th international conference on machine learning (ICML-10)*, 2010, pp. 807–814.
- [34] A. L. Maas, A. Y. Hannun, A. Y. Ng, Rectifier nonlinearities improve neural network acoustic models, in: *Proc. icml*, Vol. 30, 2013, p. 3.
- [35] D.-A. Clevert, T. Unterthiner, S. Hochreiter, Fast and accurate deep network learning by exponential linear units (elus), *arXiv preprint arXiv:1511.07289*.
- [36] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, *arXiv preprint arXiv:1412.6980*.
- [37] J. Duchi, E. Hazan, Y. Singer, Adaptive subgradient methods for online learning and stochastic optimization, *Journal of Machine Learning Research* 12 (Jul) (2011) 2121–2159.
- [38] M. D. Zeiler, Adadelta: an adaptive learning rate method, *arXiv preprint arXiv:1212.5701*.
- [39] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, Dropout: a simple way to prevent neural networks from overfitting, *The Journal of Machine Learning Research* 15 (1) (2014) 1929–1958.

- [40] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, arXiv preprint arXiv:1502.03167.
- [41] F. Chollet, keras, <https://github.com/fchollet/keras> (2015).
- [42] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al., Scikit-learn: Machine learning in python, *Journal of machine learning research* 12 (Oct) (2011) 2825–2830.
- [43] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al., Tensorflow: A system for large-scale machine learning, in: 12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16), 2016, pp. 265–283.