# LABEL UNCERTAINTY-GUIDED MULTI-STREAM MODEL FOR DISEASE SCREENING

*Chi Liu*[1]  *Zongyuan Ge*[2]  *Mingguang He*[1,3]  *Xiaotong Han*[1] *

[1] Zhongshan Ophthalmic Center, Sun Yat-sen University, Guangzhou, China
[2] Monash eResearch Centre, Monash University, Melbourne, Australia
[3] Centre for Eye Research, Melbourne University, East Melbourne, Victoria, Australia

## ABSTRACT

The annotation of disease severity for medical image datasets often relies on collaborative decisions from multiple human graders. The intra-observer variability derived from individual differences always persists in this process, yet the influence is often underestimated. In this paper, we cast the intra-observer variability as an uncertainty problem and incorporate the label uncertainty information as guidance into the disease screening model to improve the final decision. The main idea is dividing the images into simple and hard cases by uncertainty information, and then developing a multi-stream network to deal with different cases separately. Particularly, for hard cases, we strengthen the network's capacity in capturing the correct disease features and resisting the interference of uncertainty. Experiments on a fundus image-based glaucoma screening case study show that the proposed model outperforms several baselines, especially in screening hard cases.

***Index Terms***— Label uncertainty, disease screening.

## 1. INTRODUCTION

Deep learning (DL) models for image-based disease screening heavily rely on large-scale datasets consisting of $(x_i, y_i)$ pairs, where $x_i$ is the image data instance and $y_i$ is the corresponding label of disease severity. Practically, the label annotation is performed by multiple human experts (e.g., doctors or trained graders) in a collaborative pattern, where the consistent decision from the majority is often regarded as the "ground truth" of $y_i$ [1, 2, 3, 4]. Given individual differences in e.g., domain expertise, judgement and bias, intra-observer variability always persists in the annotation process [5, 6]. However, defining ground truth by majority consistency is somehow oversimplified, underestimating the influence of intra-observer variability and may even mislead the ground truth. For instance, two images, one is labelled as "positive" by $100\%$ of graders while the other one by $51\%$ of graders, will have the same final annotation as "positive", despite the grading difficulties being intuitively unequal.

In fact, the label uncertainty resulted from intra-observer variability contains important prior information, which im-
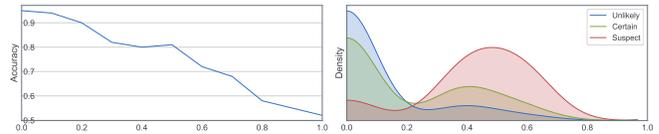
---
*indicates the corresponding author. Email: lh.201205@aliyun.com



**Fig. 1**. (a) A DL model's screening accuracy in different label uncertainty groups; (b) GON severity distributions on different label uncertainty. Both x-axes are label uncertainty scores.

plies how difficult an image is in identifying its disease severity. Previous studies exploited such information for quantifying imaging quality [7], or identifying hard patient cases that may require a medical second opinion [8]. Different from previous viewpoints, we believe that label uncertainty can provide prior guidance to improve DL models' decisions, and should be carefully considered in the initial model design. We have two observations in a fundus image-based glaucomatous optic neuropathy (GON) screening task:

***O1:*** *A DL model generally performs worse on samples with higher label uncertainty.* Fig.1.a shows the performances of a screening model in different uncertainty score (computed by Eq.1) groups. The average group accuracy degrades with the uncertainty score increases.

***O2:*** *Disease severity distribution is potentially correlated with label uncertainty.* Fig.1.b shows the distributions of GON severity in images with different uncertainty scores. The label uncertainty of suspect GON samples are generally higher than those of unlikely and certain GON samples.

The insights stemming from the observations is that DL screening models should be improved on the samples with higher label uncertainty with identifying "purer" disease features in the presence of label uncertainty. To this end, we first model label uncertainty as scalar scores using the empirical distribution of intra-observer variability. The images are accordingly divided into simple and hard cases. Then a multi-stream disease screening model is proposed, comprising two main streams, one for simple case screening (SC-Net) and one for hard case screening (HC-Net), and an auxiliary stream (US-Net) for extracting uncertainty-associated features and predicting uncertainty scores. The label uncertainty information is elaborately embedded into the HC-Net's learning

process in various ways to offer guidance for classifying hard cases. Specifically, in the training phase, the variability-based encoding and uncertainty-guided joint loss are used to enforce the network to capture correct disease features in the presence of label uncertainty and disentangle the uncertainty-associated features in the latent space. In the inference phase, the uncertainty-guided adaptive threshold is applied to unseen samples to improve the model's decision on hard cases.

## 2. METHODOLOGY

### 2.1. Modeling label uncertainty

Given an image $x_i$, we let $y_i^{(1)}, y_i^{(2)}, ..., y_i^{(M)}$ be the assigned labels from $M$ individual graders. In practice, the minimum of $M$ is 3 to gain a significant voting. Letting $c_1, ..., c_K$ be $K$ classes of disease severity, the empirical distribution of intra-observer variability is $\hat{\mathbf{p}}_\mathbf{i} = \left[ \hat{p}_i^{(1)}, ..., \hat{p}_i^{(K)} \right]$, with $\hat{p}_i^{(k)} = \frac{1}{M} \sum_m \mathcal{H}(y_i^{(m)} = c_k)$ where $\mathcal{H}(\cdot) = 1$ when the inside condition stands, otherwise 0 [8]. Then the label uncertainty of $x_i$ can be defined by entropy which indicates the stability of a system:

$$u_i = -\sum_k \hat{p}_i^{(k)} \log \hat{p}_i^{(k)}. \tag{1}$$

$u_i$ reaches its peak when human disagreements are equally distributed on $K$ classes, indicating the most unstable state.

Note that label uncertainty differs from model uncertainty, also known as epistemic uncertainty. The latter measures the model's confidence in making a decision. Those probabilistic approximation-based method for estimating model uncertainty, such as Monte Carlo Dropout (MC-Dropout) [9], cannot be directly applied to label uncertainty.

### 2.2. Multi-stream screening model

Fig. 2 shows an overview of the training phase (left) and the inference phase (right) of the multi-stream screening model. The three sub-streams (US-Net, SC-Net and HC-Net) employ the same deep neural network (DNN) structure as backbones. Notably, the proposed model is independent with DNN structures and compatible with various DNN backbones.

#### 2.2.1. US-Net stream

The US-Net is an auxiliary stream for disentangling uncertainty-associated features and predicting uncertainty scores for unseen samples (see Section 2.2.2). It is optimized with a mean squared error (MSE) loss (Fig.2: the purple learning flow). Given the training data pair $(x_i, u_i)$, the MSE loss is:

$$\mathcal{L}_{MSE} = \frac{1}{N} \sum_i (u_i - \tilde{u}_i)^2, \tag{2}$$

where $\tilde{u}_i$ is the predicted uncertainty score.

#### 2.2.2. SC-Net and HC-Net streams

In the training phase, the image samples are divided into simple cases and hard cases via a predefined threshold on the empirical uncertainty score $u$. Then the two cases are fed into the SC-Net and the HC-Net for learning disease severity representation, respectively. The learning of disease severity can be formulated as a multi-classification problem.

Regarding simple cases, the underlying disease severity features can be well captured by a vanilla DNN. Since human graders easily achieve a high consistency for simple cases, we follow previous studies [1, 2, 3, 4]using the majority-voting result $\bar{y}$ as the training ground truth. For the training pair $(x_i, \bar{y}_i)$, the optimization is performed with a cross-entropy loss (shown as the green learning flow in Fig. 2):

$$\mathcal{L}_{CE} = -\frac{1}{N} \sum_i \sum_k \bar{p}_i^{(k)} \log \left( \tilde{p}_i^{(k)} \right), \bar{p}_i^{(k)} \in \bar{\mathbf{p}}_\mathbf{i}, \tilde{p}_i^{(k)} \in \tilde{\mathbf{p}}_\mathbf{i} \tag{3}$$

where $\bar{\mathbf{p}}_\mathbf{i}$ denotes the one-hot encoding vector of $\bar{y}_i$ and $\tilde{\mathbf{p}}_\mathbf{i}$ is the class probability yielded by the softmax layer of DNN.

Compared to simple cases, hard cases are more difficult for both human and vanilla DNNs to recognize the disease severity correctly. Therefore, directly applying the majority-voting-based ground truth and the SC-Net for hard case learning is impractical. Taking advantage of label uncertainty information, we design the HC-Net with several specific strategies to address the problem. The red arrow in Fig. 2 illustrates the learning flow of HC-Net.

**Variability-based encoding.** Rethinking the above one-hot encoding method for a multi-classification problem, where the target class is assigned with a scalar 1 while other classes are 0, a strong assumption is that all samples belonging to the same class contain equal amount of ground truth information in the label space, irrespective of their intrinsic difficulties in disease screening being different. This could lead to a model bias in identifying the correct disease features for hard cases. Instead, we propose the variability-based encoding method which applies the empirical distribution of intro-observer variability $\hat{\mathbf{p}}_\mathbf{i} = \left[ \hat{p}_i^{(1)}, ..., \hat{p}_i^{(K)} \right]$ as the ground-truth. The label uncertainty information retained in the encoded vector can help ensuring the model to capture disease features properly in the presence of uncertainty. Particularly, when all graders arrive at the same decision (i.e., no uncertainty exists), $\hat{\mathbf{p}}_\mathbf{i}$ equals to the one-hot vector $\bar{\mathbf{p}}_\mathbf{i}$.

**Uncertainty-guided focal loss.** Inspired by the focal loss [10], we propose the uncertainty-guided focal loss to replace the cross-entropy loss, which can promote the HC-Net to pay more attentions to hard cases during training:

$$\mathcal{L}_{UGF} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{k=1}^{K} \left[ \left( 1 - \tilde{p}_i^{(k)} \right)^{g(u_i)} \cdot \hat{p}_i^{(k)} \log \left( \tilde{p}_i^{(k)} \right) \right], \tag{4}$$

where $\left( 1 - \tilde{p}_i^{(k)} \right)^{g(u_i)}$ is imposed to adjust the model's focus
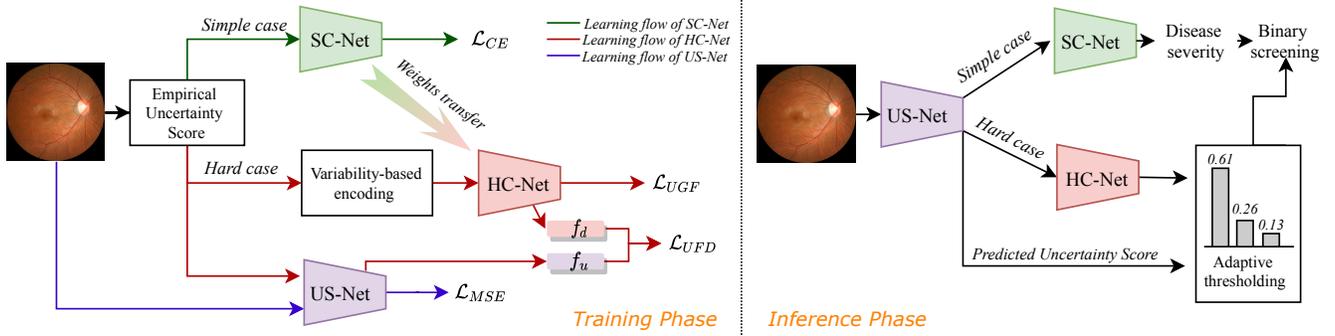
**Fig. 2**. Overview of the training phase and the inference phase of the uncertainty-guided multi-stream screening model.

according to the difficulty of given samples, i.e., attaching more importance to the samples with lower prediction confidences and higher uncertainty scores; $g(u_i) = \gamma \cdot u_i$ is an adjusting function relying on $u_i$ with a constant weight $\gamma$.

**Uncertainty feature decoupling loss.** According to the observations in Fig. 1, the potential correlation between uncertainty and disease severity may bias the screen model's decision. Therefore, we propose the uncertainty feature decoupling loss to disentangle disease features $f_{di}$ from uncertainty features $f_{un}$ in the latent space.

$$\mathcal{L}_{UFD} = -\frac{1}{N}\sum_i \max(0, h(u_i) - \mathcal{D}(f_{di}(x_i), f_{un}(x_i))),$$
(5)

where $f_{di}(x_i)$ and $f_{un}(x_i)$ are the flattened feature vectors of image $x_i$ output by the last convolutional layer of the HC-Net and the pre-trained US-Net, respectively; $h(u_i) = \min(\alpha \cdot u_i, 1)$ imposes a $u_i$-guided dynamic margin to ensure a lower bound of the two features' distance, where the case is harder, the margin is larger. $\mathcal{D}(\cdot)$ is the Pearson distance [11].

Consequently, the final loss for optimizing the HC-Net is:

$$\mathcal{L}_{joint} = \mathcal{L}_{UGF} + \mathcal{L}_{UFD}.$$
(6)

**Uncertainty-guided adaptive threshold.** The practical population-scale disease screening models are normally expected to make a binary decision in the inference phase to support clinical recommendations, e.g., "non-referable" versus "referable" cases [1, 2, 3, 4]. A common way is applying a fixed threshold (e.g. 0.5) to $\tilde{\mathbf{p}}_i$ to identify the best trade-off between sensitivity and specificity of a classification model [12, 13, 14]. The fixed threshold could degrade the flexibility facing cases with varying uncertainty. Instead, we design an uncertainty-guided adaptive threshold for this process.

Given an unseen sample, first its uncertainty score $\tilde{u}$ is predicted by the US-Net to decide whether the sample should be allocated to SC-Net or HC-Net. For HC-Net screening, an adaptive threshold $\tau$ is applied to the inference probability of the negative class:

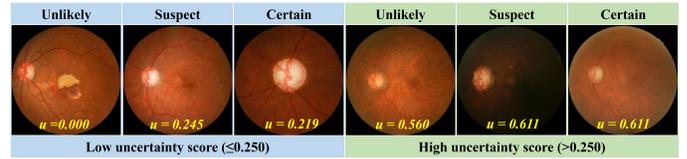$$\tau = 1 - \frac{K-1}{K}(1-\tilde{u})^\beta,$$
(7)



**Fig. 3**. Image examples with different GON severity levels and empirical uncertainty scores from the GON dataset.

where $K$ is the number of GON severity classes and $\beta$ is an adjustable weight; $\tau$ rises when $\tilde{u}$ increases, meaning that for samples with larger uncertainty, the model is allowed to be more inclined to classify them as "referable".

## 3. EXPERIMENTS

### 3.1. Experimental settings

**Datasets.** The evaluations were performed on the fundus image-based GON screening dataset *LabelMe* [3]. Images were collected from various clinical settings in China. Twenty-one ophthalmologists participated in the grading process. Each image was assigned to different graders sequentially until three consistent individual decisions were achieved. The five-class grading criteria was applied [3], including unlikely, suspect and certain GON, and poor quality/location. Images with poor quality/location were excluded and the final dataset consists of 47012 images. The dataset was randomly divided into training, validation and testing sets, as detailed in Table 1. Note that in all subsets, the mean uncertainty scores (Eq. 1) of suspect GON samples are significantly higher than the other two classes, in line with our previous observation *O2*. According to the observation in Fig.1.b, we empirically set the predefined threshold of $u$ as 0.25. The bottom two rows show the details of simple cases and hard cases divided by this threshold. Fig.3 shows some examples with different GON severity levels and uncertainty scores from this dataset.

**Settings.** All images were central-cropped and downsized to

**Table 1**. The details of the GON dataset. MUS: mean uncertainty score.

| | Training set | | Validation set | | Testing set | |
|---|---|---|---|---|---|---|
| | N | MUS | N | MUS | N | MUS |
| **Unlikely** | 27525 | 0.068 | 6147 | 0.068 | 1511 | 0.067 |
| **Suspect** | 2338 | 0.410 | 518 | 0.423 | 94 | 0.412 |
| **Certain** | 6947 | 0.162 | 1537 | 0.158 | 395 | 0.151 |
| **Simple cases** | 28050 | 0.0003 | 6243 | 0.0003 | 1550 | 0.0003 |
| **Hard cases** | 8760 | 0.452 | 1959 | 0.447 | 450 | 0.442 |

$299 * 299$ with pixel values rescaled to the range of $[0, 1]$. We employed Xception [15] as the backbone DNN for all sub-streams. The SC-Net and US-Net were initialized with random weights. The HC-Net were initialized with the pre-trained SC-Net weights in a transfer learning manner [4], due to that the shallow layers learn generic low-level domain features that can be shared in similar tasks to accelerate training. The batch size was 220 and the epoch number was 100. The Adam optimizer [16] with an initial learning rate of 0.01 plus an adaptively-decay scheduler was adopted. The hyper-parameters $\{\gamma, \alpha, \beta\}$ were empirically set as $\{4, 1.4, 2\}$ (the decision process is detailed in the supplementary).

**Evaluation metrics.** For the multi-classification of GON severity, we computed *F1*-score for each single severity class and the classification accuracy for the overall performance. For the final binary GON screening (i.e., referable vs. non-referable), sensitivity ($SE$), specificity ($SP$) and area under the ROC curve ($AUC$) were evaluated.

### 3.2. Results

**Ablation experiments.** The ablation experiments were performed in two groups - the whole testing set and only the hard testing cases, respectively. The candidate models include a base model without any strategies (equivalent to the SC-Net) and models with gradually-added strategies. The *uncertainty-guided adaptive threshold* is excluded here since it is only applicable to the final binary screening setting. Table 2 shows the $F_1$ scores of individual classes and the overall classification accuracy. For all classes in both groups, the $F_1$ scores and the overall accuracy generally increases with adding more strategies to the base model. The final model combining three strategies outperforms all other models in both groups, and the overall improvement is much more significant in the hard case group ($72.89\% \rightarrow 84.22\%$), indicating the effectiveness of the proposed strategies.

**Screening performance.** We tested the final binary GON screening performance of the multi-stream model with all strategies (also including the *uncertainty-guided adaptive threshold*). According to the screening criteria [3], unlikely GON is regarded as "non-referable" while the other two classes are "referable". To show our model's compatibility with different DNN structures, we employed three DNNs

**Table 2**. Results of the ablation experiments. *M1*: Base model; *M2*: *M1* + variability-based encoding; *M3*: *M2* + uncertainty-guided focal loss; *M4*: *M3* + uncertainty feature decoupling loss. The highest value is in **bold**.

| | The whole dataset | | | | Hard cases only | | | |
|---|---|---|---|---|---|---|---|---|
| | Unlikely | Suspect | Certain | *Overall* | Unlikely | Suspect | Certain | *Overall* |
| *M1* | 95.84 | 30.99 | 88.83 | 92.11 | 80.89 | 31.40 | 76.60 | 72.89 |
| *M2* | 95.83 | 43.04 | 88.70 | 92.28 | 83.13 | 47.41 | 77.42 | 76.01 |
| *M3* | 96.82 | 55.90 | 90.80 | 93.94 | 85.59 | 58.33 | 79.42 | 79.33 |
| *M4* | **97.08** | **62.35** | **91.56** | **94.49** | **88.70** | **67.97** | **85.61** | **84.22** |

**Table 3**. Binary screening performances of the proposed multi-stream model (with three different backbones as shown in brackets) and two baselines. The highest value is in **bold**.

| | The whole dataset | | | Hard cases only | | |
|---|---|---|---|---|---|---|
| | $SE$ | $SP$ | $AUC$ | $SE$ | $SP$ | $AUC$ |
| **BCNet [3]** | 87.59 | 92.96 | 95.54 | 80.23 | 85.83 | 87.80 |
| **DENet [17]** | 92.51 | 95.96 | 96.79 | 84.23 | 88.10 | 90.84 |
| **CaliNet [18]** | 90.11 | 91.27 | 93.74 | 81.01 | 86.22 | 89.04 |
| **Ours (Inception-V3)** | 90.60 | 94.76 | 96.13 | 85.21 | 85.22 | 93.56 |
| **Ours (ResNet50)** | **92.59** | 96.36 | 98.21 | **89.44** | 88.13 | 95.30 |
| **Ours (Xception)** | 91.65 | **97.68** | **98.90** | 89.12 | **89.57** | **95.79** |

as backbone, including Inceptivon-V3 [19], ResNet-50 [20] and Xception [15]. We also compared our method with two state-of-the-art models for large-scale fundus image-based GON screening: an Inception-V3-based binary classification network (BCNet) [3] and a disc-aware ensemble network (DENet) [17], as well as a label uncertainty-based model calibration method (CaliNet) [18]. Table 3 shows the results in two groups in terms of $SE, SP$ and $AUC$ scores. The models with three different backbones can all achieve $AUC$ scores larger than $96.00\%$ in the whole dataset group and $93.50\%$ in the hard case group, showcasing its satisfied DNN-compatibility. And their performances are generally better or comparable than the two baselines in both groups. Likewise, the improvement in the hard case group is much more significant, e.g., the $AUC$ score is raised from $90.84\%$ from DENet, which is the best result of baselines, to $95.79\%$ from our method with the Xception backbone.

### 3.3. Conclusion

In this paper, we investigated how to leverage the label uncertainty existing in medical image annotation as prior guidance to meliorate disease screening models' decisions. We developed a multi-stream model for the cases with different uncertainty levels, where multiple uncertainty-guided strategies were incorporated specifically for improvement on cases with high label uncertainty. The evaluations conducted in a GON screening case study showed the effectiveness of our method. Our method can benefit general DL models developed on medical image dataset annotated by multiple graders.

**Compliance with Ethical Standards.** This study was performed in line with the principles of the Declaration of Helsinki. Approval was granted by the he institutional review board of Zhongshan Ophthalmic Center (2017KYPJ049). The review board determined that informed consent was not necessary in this study due to the retrospective nature and fully anonymised use of images.

## 4. REFERENCES

[1] Zhixi Li, Stuart Keel, Chi Liu, et al., "An automated grading system for detection of vision-threatening referable diabetic retinopathy on the basis of color fundus photographs," *Diabetes care*, vol. 41, no. 12, pp. 2509–2516, 2018.

[2] Varun Gulshan, Lily Peng, Marc Coram, et al., "Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs," *Jama*, vol. 316, no. 22, pp. 2402–2410, 2016.

[3] Zhixi Li, Yifan He, Stuart Keel, et al., "Efficacy of a deep learning system for detecting glaucomatous optic neuropathy based on color fundus photographs," *Ophthalmology*, vol. 125, no. 8, pp. 1199–1206, 2018.

[4] Daniel S Kermany, Michael Goldbaum, Wenjia Cai, et al., "Identifying medical diagnoses and treatable diseases by image-based deep learning," *Cell*, vol. 172, no. 5, pp. 1122–1131, 2018.

[5] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, et al., "Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, pp. 590–597.

[6] Lisa S Abrams, Ingrid U Scott, George L Spaeth, et al., "Agreement among optometrists, ophthalmologists, and residents in evaluating the optic disc for glaucoma," *Ophthalmology*, vol. 101, no. 10, pp. 1662–1667, 1994.

[7] Zhibin Liao, Hany Girgis, Amir Abdi, et al., "On modelling label uncertainty in deep neural networks: automatic estimation of intra-observer variability in 2d echocardiography quality assessment," *IEEE transactions on medical imaging*, vol. 39, no. 6, pp. 1868–1883, 2019.

[8] Maithra Raghu, Katy Blumer, Rory Sayres, et al., "Direct uncertainty prediction for medical second opinions," in *International Conference on Machine Learning*. PMLR, 2019, pp. 5281–5290.

[9] Yarin Gal and Zoubin Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *international conference on machine learning*. PMLR, 2016, pp. 1050–1059.

[10] Tsung-Yi Lin, Priya Goyal, Ross Girshick, et al., "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.

[11] MH Fulekar, *Bioinformatics: applications in life and environmental sciences*, Springer Science & Business Media, 2009.

[12] Elizabeth A Freeman and Gretchen G Moisen, "A comparison of the performance of threshold criteria for binary classification in terms of predicted prevalence and kappa," *Ecological modelling*, vol. 217, no. 1-2, pp. 48–58, 2008.

[13] Seung Seog Han, Myoung Shin Kim, Woohyung Lim, et al., "Classification of the clinical images for benign and malignant cutaneous tumors using a deep learning algorithm," *Journal of Investigative Dermatology*, vol. 138, no. 7, pp. 1529–1538, 2018.

[14] Junghwan Cho, Ki-Su Park, Manohar Karki, et al., "Improving sensitivity on identification and delineation of intracranial hemorrhage lesion using cascaded deep learning models," *Journal of digital imaging*, vol. 32, no. 3, pp. 450–461, 2019.

[15] François Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.

[16] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[17] Huazhu Fu, Jun Cheng, Yanwu Xu, et al., "Disc-aware ensemble network for glaucoma screening from fundus image," *IEEE Transactions on Medical Imaging*, vol. 37, no. 11, pp. 2493–2501, 2018.

[18] Martin Holm Jensen, Dan Richter Jørgensen, Raluca Jalaboi, et al., "Improving uncertainty estimation in convolutional neural networks using inter-rater agreement," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 540–548.

[19] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, et al., "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.