# GLAUCOMA DETECTION FROM RAW CIRCUMPAPILLARY OCT IMAGES USING FULLY CONVOLUTIONAL NEURAL NETWORKS

*Gabriel García*[1], *Rocío del Amor*[1], *Adrián Colomer*[1], *Valery Naranjo*[1]

[1]Instituto de Investigación e Innovación en Bioingeniería (I3B),
Universitat Politècnica de València, Camino de Vera s/n, 46022, Valencia, Spain.

## ABSTRACT

Nowadays, glaucoma is the leading cause of blindness worldwide. We propose in this paper two different deep-learning-based approaches to address glaucoma detection just from raw circumpapillary OCT images. The first one is based on the development of convolutional neural networks (CNNs) trained from scratch. The second one lies in fine-tuning some of the most common state-of-the-art CNNs architectures. The experiments were performed on a private database composed of 93 glaucomatous and 156 normal B-scans around the optic nerve head of the retina, which were diagnosed by expert ophthalmologists. The validation results evidence that fine-tuned CNNs outperform the networks trained from scratch when small databases are addressed. Additionally, the VGG family of networks reports the most promising results, with an area under the ROC curve of 0.96 and an accuracy of 0.92, during the prediction of the independent test set.

***Index Terms***— Glaucoma detection, deep learning, circumpapillary OCT, fine tuning, class activation maps.

## 1. INTRODUCTION

Glaucoma has become the leading cause of blindness worldwide, according to [1]. It is characterized by causing progressive structural and functional damage to the retinal optic nerve head (ONH). Recent studies advocate that roughly 50% of people suffering from glaucoma in the world are undiagnosed and ageing populations suggest that the impact of glaucoma will continue to rise, affecting 111.8 million people in 2040 [2]. Therefore, early treatment of this chronic disease could be essential to prevent irreversible vision loss.

Currently, a complete glaucoma study usually includes medical history, fundus photography, visual field (VF) analysis, tonometry and optic nerve imaging tests such as optical coherence tomography (OCT). Most of the state-of-the-art studies addressed the glaucoma detection via fundus image

analysis, making use of visual field tests and relevant parameters like the intraocular pressure (IOP) [3, 4]. Specifically, J. Gmez-Valverde et al. [5] performed a comparison between convolutional neural networks (CNNs) trained from scratch and using fine-tuning techniques. Also, the authors in [6, 7] considered the use of transfer learning and fine-tuning methods applied to very popular state-of-the-art network architectures to identify glaucoma on fundus images. Other studies such as [8,9] carried out a combination between OCT B-scans and fundus images to obtain an RNFL thickness probability map which was used as an input to the CNNs. In this paper, contrary to the studies of the literature, we propose an end-to-end system for glaucoma detection based only on raw circumpapillary OCT images, without using another kind of images or external expensive tests related to the VF and IOP parameters. It is important to highlight that circumpapillary OCT images as shown in Fig. 1 correspond to circular scans located around the ONH, where rich information about different retinal layers structures can be found. Additionally, several studies claimed that circumpapillary retinal nerve fiber layer (RNFL) is essential to detect early glaucomatous damage [10–12]. For that reason, one of the main novelties of this paper is focused on demonstrating that a single circumpapillary OCT image may be of great interest when carrying out an accurate glaucoma detection.

We propose two different data-driven learning strategies to develop computer-aided diagnosis systems capable of discerning between glaucomatous and healthy eyes just from B-scans around the ONH. Several CNNs trained from scratch and different fine-tuned state-of-the-art architectures were
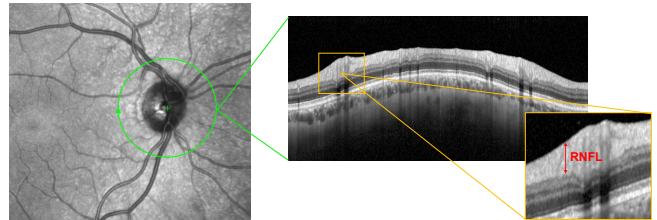
**Fig. 1**. B-scan around the retinal ONH corresponding to a circumpapillary OCT image. RNFL is highlighted in red.

considered. Furthermore, we propose, for the first time in this kind of images, the class activation maps computation in order to compare the location information reported by the clinicians with the heat maps generated by the developed models. Heat maps allow highlighting the regions in which the networks pay attention to determine the class of each specific sample.

## 2. MATERIAL

The experiments detailed in this paper were performed on a private database composed of 249 OCT images of dimensions $M \times N = 496 \times 768$ pixels. In particular, 156 normal and 93 glaucomatous circumpapillary samples were analysed from 89 and 59 patients, respectively. Each B-scan was diagnosed by experts ophthalmologists from Oftalvist Ophthalmic Clinic. Note that *Heidelberg Spectrallis* OCT system was employed to acquire the circumpapillary OCT images with an axial resolution of 4-5 $\mu$m.

## 3. METHODOLOGY

### 3.1. Data Partitioning

A data partitioning stage was carried out to divide the database into different training and test sets. Specifically, $\frac{4}{5}$ of the circumpapillary images, which corresponds to 73 glaucomatous and 124 normal samples, from 12 and 18 patients respectively, composed the training set, whereas the test set was defined by $\frac{1}{5}$ of the data (20 with glaucoma and 32 normal B-scans from 12 and 18 patients). In addition, for the training set, we also performed an internal cross-validation (ICV) stage to control the overfitting, as well as to select the best neural network hyper-parameters. Finally, the independent test set was used to evaluate the definitive predictive models, which were created using the entire training set.

### 3.2. Learning from scratch

Similarly to the methodology exposed in [5], we propose in this paper the use of shallow CNNs from scratch to address the glaucoma detection, taking into account the significant differences between our grey-scale circumpapillary OCT images and other large databases containing natural images, which are widely used for transfer-learning techniques.

During the internal cross-validation (ICV) stage, an empirical exploration was carried out to determine the best hyper-parameter combination in terms of minimisation of the binary cross-entropy loss function. Different network architectures composed of diverse learning blocks were developed. In particular, convolutional, pooling, batch normalisation and dropout layers were considered to address the feature extraction stage. The variable components of each layer, such as the convolutional filters, pooling size, dropout coefficients, as

well as the number of convolutional layers in each block were optimised during the experimental phase. Regarding the top model, the use of flatten, dropout and fully-connected layers with a different number of neurons was studied. Also, global max and global average pooling layers were analysed in order to reduce the number of trainable parameters. Moreover, we implemented an optimal weighting factor of [1.35, 0.79] during the training of the models to alleviate the unbalanced problem between classes.

After the ICV stage, the best CNN architecture was found using four convolutional blocks, as it is detailed in Table 1. It is remarkable the use of the global max-pooling (GMP) layer applied in the last block, which allows extracting the maximum activation of each convolutional filter before the classification layer. Also, note that batch normalization and dropout layers were not used because no improvement was reported during the validation phase. Only a dense layer with a *softmax* activation and 2 neurons, corresponding to glaucoma and healthy classes, was defined.

**Table 1**. Proposed CNN architecture trained from scratch.

| Layer name | Output shape | Filter size |
|---|---|---|
| Input layer | 496 x 768 x 1 | N/A |
| Conv1_1 | 496 x 768 x 32 | 3 x 3 x 32 |
| MaxPooling | 248 x 384 x 32 | 2 x 2 x 32 |
| Conv2_1 | 248 x 384 x 64 | 3 x 3 x 64 |
| MaxPooling | 124 x 192 x 64 | 2 x 2 x 64 |
| Conv3_1 | 124 x 192 x 128 | 3 x 3 x 128 |
| MaxPooling | 62 x 96 x 128 | 2 x 2 x 128 |
| Conv4_1 | 62 x 96 x 256 | 3 x 3 x 256 |
| MaxGlobalPool | 256 | N/A |
| Dense (softmax) | 2 | N/A |

The optimal hyper-parameters combination was achieved by training the CNNs during 150 epochs, using Adadelta optimizer with a learning rate of 0.05 and a batch size of 16. It should be noticed that we also proposed the use of data augmentation (DA) techniques [13] to elucidate how important is the creation of artificial samples when addressing small databases. Specifically, a factor ratio of 0.2 was applied here to perform random geometric and dense elastic transformations from the original images.

### 3.3. Learning by fine tuning

Deeper architectures networks could improve the models' performance, but a large number of images annotated by experts would be necessary for training a deep CNN from scratch. For this reason, we propose in this section the use of fine-tuning techniques [14], which allows training CNNs with greater depth using the weights pre-trained on large databases, without the need to train from scratch. In particular, we applied a deep fine-tuning [15] strategy to transfer

the wide knowledge acquired by several state-of-the-art networks, such as VGG16, VGG19, InceptionV3, Xception and ResNet, when they were trained on the large *ImageNet* data set. Attending to the small database used in this work, only the coefficients of the last convolutional blocks (4 and 5) were retrained with the specific knowledge corresponding to the circumpapillary OCT images. The rest of coefficients were frozen with the values of the weights pre-trained with 14 million of natural images contained in *Imagenet* database.

Additionally, similarly to the proposed learning from scratch strategy, an empirical exploration of different hyper-parameters and top-model architectures was considered for all networks. It is important to notice that InceptionV3, Xception and ResNet architectures reported a poor performance due to their extensive depth (42, 36 and 53 convolutional layers, respectively). However, the family of VGG architectures achieved the best performance, in line with the findings in the literature [5]. Specifically, VGG16 base model is composed of five convolutional blocks according to Fig. 2, where blue boxes correspond to convolutional layers activated with *ReLu* functions and red-grey boxes represent max-pooling layers. VGG19 base model is composed of the same architecture, but including an extra convolutional layer in the last three blocks.

A top model composed of global max pooling and dropout layers with a coefficient of 0.4, followed by a softmax layer with two neurons, provided the best model performance when VGG architectures were fine-tuned (see Fig. 2). Regarding the selection of hyper-parameters combination, Adadelta optimizer with a learning rate of 0.001 reported the best learning curves when the model was forward, and backward, propagated during 125 epochs with a batch size of 16, trying to minimise the binary cross-entropy loss function.

Note that an initial down-sampling ×0.5 of the original images was necessary to alleviate the GPU memory problems during the training phase. Besides, replicating ×3 the channels of the grey-scale was necessary to adapt the input images in order to fine tune the CNNs. Data augmentation (DA) techniques with a factor of 0.2 were also considered.
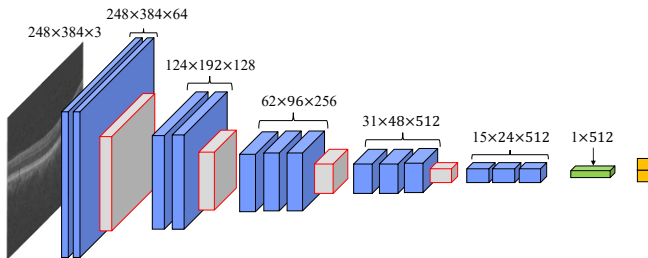


**Fig. 2**. Network architecture used to discern between glaucomatous and healthy OCT samples by fine-tuning the VGG16 base model. Note that numeric values of the filters are correctly defined in the image, although they do not correspond to the representation size of the boxes due to space problems.

# 4. RESULTS AND DISCUSSION

## 4.1. Validation results

In this stage, we present the results achieved during the ICV stage for each of the proposed CNNs. We expose in Table 2 a comparison of the CNNs trained from scratch, in terms of mean ± standard deviation. Several figures of merit are calculated to evidence the differences between using or not data augmentation (DA) techniques. In particular, sensitivity (SN), specificity (SPC), positive predictive value (PPV), negative predictive value (NPV), F-score (FS), accuracy (ACC) and area under the ROC curve (AUC) are employed.

**Table 2**. Classification results reached during the ICV stage from the proposed CNNs trained from scratch.

|  | Without DA | With DA |
|---|---|---|
| **SN** | $0.7657 \pm 0.2032$ | $\mathbf{0.8771 \pm 0.1281}$ |
| **SPC** | $\mathbf{0.9270 \pm 0.1302}$ | $0.8047 \pm 0.1514$ |
| **PPV** | $\mathbf{0.8721 \pm 0.0662}$ | $0.7477 \pm 0.14061$ |
| **NPV** | $0.8808 \pm 0.0971$ | $\mathbf{0.9224 \pm 0.0678}$ |
| **FS** | $\mathbf{0.8016 \pm 0.1309}$ | $0.7980 \pm 0.10745$ |
| **ACC** | $\mathbf{0.8679 \pm 0.0781}$ | $0.8315 \pm 0.0985$ |
| **AUC** | $0.9152 \pm 0.0490$ | $\mathbf{0.9319 \pm 0.0386}$ |

Significant differences between CNNs trained with and without data augmentation techniques can be appreciated in Table 2, especially related to the sensitivity and specificity metrics. Worth noting that the learning curves relative to the CNN trained without implementing DA algorithms reported slight overfitting during the validation phase. This fact is evidenced in the high sensitivity standard deviation of the model.

Additionally, we also detail in Table 3 the validation results achieved from the fine-tuned VGG networks, since they provided a considerable outperforming with respect to the rest of state-of-the-art architectures during the ICV stage. Specifically, VGG16 reaches better results for all figures of merit, although both architectures report similar behaviour. In comparison to the CNNs trained from scratch, VGG16 provides the best model performance too.

**Table 3**. Results comparison between the best fine-tuned CNNs proposed during the validation phase.

|  | VGG16 | VGG19 |
|---|---|---|
| **SN** | $\mathbf{0.7800 \pm 0.1302}$ | $0.7400 \pm 0.1462$ |
| **SPC** | $\mathbf{0.9677 \pm 0.0334}$ | $0.9597 \pm 0.0283$ |
| **PPV** | $\mathbf{0.9401 \pm 0.0643}$ | $0.9180 \pm 0.0602$ |
| **NPV** | $\mathbf{0.8864 \pm 0.0662}$ | $0.8670 \pm 0.0692$ |
| **FS** | $\mathbf{0.8466 \pm 0.0720}$ | $0.8131 \pm 0.0936$ |
| **ACC** | $\mathbf{0.8984 \pm 0.0468}$ | $0.8786 \pm 0.0563$ |
| **AUC** | $\mathbf{0.9463 \pm 0.0339}$ | $0.9416 \pm 0.0501$ |

## 4.2. Test results

In order to provide reliable results, an independent test set was used to carry out the prediction stage. Table 4 shows a comparison between all proposed deep-learning models to evaluate their prediction ability by means of different figures of merit. Additionally, we expose in Fig. 3 the ROC curve relative to each proposed CNN to visualise the differences.

**Table 4**. Classification results achieved during the prediction stage from the proposed CNNs trained from scratch (FS) and fine-tuning the VGGs network architectures.

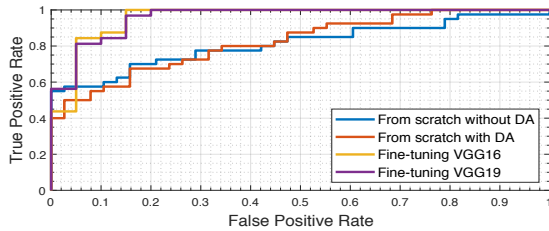|       | FS without DA | FS with DA | VGG16  | VGG19  |
|-------|---------------|------------|--------|--------|
| **SN**  | 0.7632 | 0.7895 | **0.8510** | **0.8510** |
| **SPC** | 0.7250 | 0.6750 | 0.9064 | **0.9688** |
| **PPV** | 0.7250 | 0.6977 | 0.8490 | **0.9444** |
| **NPV** | 0.7632 | 0.7714 | 0.9063 | **0.9118** |
| **FS**  | 0.7436 | 0.7407 | 0.8500 | **0.8947** |
| **ACC** | 0.7436 | 0.7308 | 0.8846 | **0.9230** |
| **AUC** | 0.8132 | 0.8230 | 0.9578 | **0.9594** |



**Fig. 3**. ROC curves corresponding to the prediction results reached from the different proposed CNNs.

Test results exposed in Fig. 4 are in line with those achieved during the validation phase. However, due to the randomness effect of the data partitioning (which is accentuated in small databases), significant differences may exist in the prediction of each subset. This fact mainly affects to the CNNs trained from scratch because all the weights of the network were trained with the images of a specific subset, whereas the proposed fine-tuned architectures keep most of the weights frozen. Regarding the ROC curves comparison, Fig. 3 shows that fine-tuned CNNs report a significant improvement in relation to the networks trained from scratch.

It is important to remark that an objective comparison with other state-of-the-art studies is difficult because there are no public databases of circumpapillary OCT images. Additionally, each group of researchers addresses glaucoma detection using a different kind of images. Notwithstanding, we detail a subjective comparison with other works based on similar methodologies applied to fundus images. In particular, [5] fine-tuned the VGG19 architecture and achieved an AUC of

0.94 predicting glaucoma. Also, [7] reached an AUC of 0.91 applying transfer learning techniques to the ResNet architecture. Otherwise, authors in [16] proposed a CNN from scratch obtaining AUC values of 0.83 and 0.89 from two independent databases. Basing on this, the proposed learning methodology exceeds the state-of-the-art results, achieving an AUC of 0.96 during the prediction of the test set.

### Class Activation Maps (CAMs)

We compute the class activation maps to generate heat maps highlighting the interesting regions in which the proposed model is paying attention to determine the class of each specific circumpapillary OCT image. In Fig. 4, we expose the CAMs relative to random specific glaucomatous and normal samples in order to elucidate what is VGG19 taking into account to discern between classes.
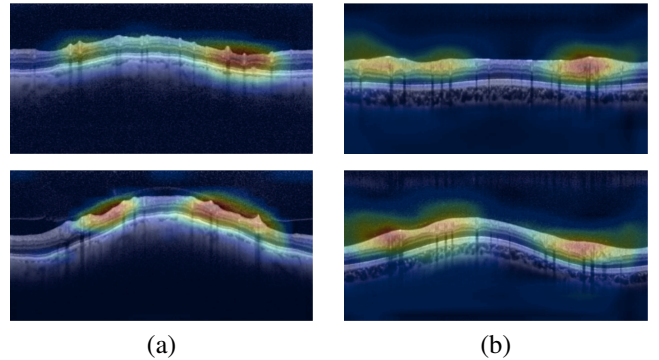


(a)                                    (b)

**Fig. 4**. Heat maps extracted from the CAMs computation for (a) glaucomatous and (b) healthy circumpapillary images.

The findings from the CAMs are directly in line with the reported by expert clinicians, who claim that a thickening of the RNFL is intimately linked with healthy patients, whereas a thinning of the RNFL evidence a glaucomatous case. That is just what heat maps in Fig. 4 reveal. Therefore, the results suggest that the proposed circumpapillary OCT-based methodology can provide a great added value for glaucoma diagnosis taking into account that information similar to that of specialists is reported by the model without including any previous clinician knowledge.

## 5. CONCLUSION

In this paper, two different deep-learning methodologies have been performed to elucidate the added value enclosed in the circumpapillary OCT images for glaucoma detection. The reported results suggest the fine-tuned VGG family of architectures as the most promising networks. The extracted CAMs evidence the great potential of the proposed model since it is able to highlight areas such as the RNFL, in line with the clinical interpretation. In future research lines, external validation of the proposed strategy with large databases is considered.

# 6. REFERENCES

[1] Jost B Jonas, Tin Aung, Rupert R Bourne, Alain M Bron, Robert Ritch, and Songhomitra Panda-Jonas, "Glaucoma–authors' reply," *The Lancet*, vol. 391, no. 10122, pp. 740, 2018.

[2] Wei Wang, Miao He, Zihua Li, and Wenyong Huang, "Epidemiological variations and trends in health burden of glaucoma worldwide," *Acta ophthalmologica*, vol. 97, no. 3, pp. e349–e355, 2019.

[3] Seong Jae Kim, Kyong Jin Cho, and Sejong Oh, "Development of machine learning models for diagnosis of glaucoma," *PLoS One*, vol. 12, no. 5, pp. e0177726, 2017.

[4] Peiyu Wang, Jian Shen, Ryuna Chang, Maemae Moloney, Mina Torres, and Bruce et al Burkemper, "Machine learning models for diagnosing glaucoma from retinal nerve fiber layer thickness maps," *Ophthalmology Glaucoma*, vol. 2, no. 6, pp. 422–428, 2019.

[5] Juan J Gómez-Valverde, Alfonso Antón, Gianluca Fatti, Bart Liefers, Alejandra Herranz, Andrés Santos, Clara I Sánchez, and María J Ledesma-Carbayo, "Automatic glaucoma classification using color fundus images based on convolutional neural networks and transfer learning," *Biomedical optics express*, vol. 10, no. 2, pp. 892–913, 2019.

[6] Naoto Shibata, Masaki Tanito, Keita Mitsuhashi, Yuri Fujino, Masato Matsuura, Hiroshi Murata, and Ryo Asaoka, "Development of a deep residual learning algorithm to screen for glaucoma from fundus photography," *Scientific reports*, vol. 8, no. 1, pp. 1–9, 2018.

[7] Mark Christopher, Akram Belghith, Christopher Bowd, James A Proudfoot, Michael H Goldbaum, Robert N Weinreb, Christopher A Girkin, Jeffrey M Liebmann, and Linda M Zangwill, "Performance of deep learning architectures and transfer learning for detecting glaucomatous optic neuropathy in fundus photographs," *Scientific reports*, vol. 8, no. 1, pp. 1–13, 2018.

[8] Hassan Muhammad, Thomas J Fuchs, Nicole De Cuir, Carlos G De Moraes, Dana M Blumberg, Jeffrey M Liebmann, Robert Ritch, and Donald C Hood, "Hybrid deep learning on single wide-field optical coherence tomography scans accurately classifies glaucoma suspects," *Journal of glaucoma*, vol. 26, no. 12, pp. 1086, 2017.

[9] Kaveri A Thakoor, Xinhui Li, Emmanouil Tsamis, Paul Sajda, and Donald C Hood, "Enhancing the accuracy of glaucoma detection from oct probability maps using convolutional neural networks," in *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2019, pp. 2036–2040.

[10] Yoshiyuki Kita, Ritsuko Kita, Ai Nitta, Chiaki Nishimura, and Goji Tomita, "Glaucomatous eye macular ganglion cell complex thickness and its relation to temporal circumpapillary retinal nerve fiber layer thickness," *Japanese journal of ophthalmology*, vol. 55, no. 3, pp. 228–234, 2011.

[11] Donald C Hood and Ali S Raza, "On improving the use of oct imaging for detecting glaucomatous damage," *British Journal of Ophthalmology*, vol. 98, no. Suppl 2, pp. ii1–ii9, 2014.

[12] Christopher KS Leung, Wai-Man Chan, Wing-Ho Yung, Alan CK Ng, Jackson Woo, Moon-Kong Tsang, and KK Raymond, "Comparison of macular and peripapillary measurements for the detection of glaucoma: an optical coherence tomography study," *Ophthalmology*, vol. 112, no. 3, pp. 391–400, 2005.

[13] Sebastien C Wong, Adam Gatt, Victor Stamatescu, and Mark D McDonnell, "Understanding data augmentation for classification: when to warp?," in *2016 international conference on digital image computing: techniques and applications (DICTA)*. IEEE, 2016, pp. 1–6.

[14] Shin Hoo-Chang, Holger R Roth, Mingchen Gao, Le Lu, Ziyue Xu, Isabella Nogues, Jianhua Yao, Daniel Mollura, and Ronald M Summers, "Deep convolutional neural networks for computer-aided detection: Cnn architectures, dataset characteristics and transfer learning," *IEEE transactions on medical imaging*, vol. 35, no. 5, pp. 1285, 2016.

[15] Nima Tajbakhsh, Jae Y Shin, Suryakanth R Gurudu, R Todd Hurst, Christopher B Kendall, Michael B Gotway, and Jianming Liang, "Convolutional neural networks for medical image analysis: Full training or fine tuning?," *IEEE transactions on medical imaging*, vol. 35, no. 5, pp. 1299–1312, 2016.

[16] Xiangyu Chen, Yanwu Xu, Damon Wing Kee Wong, Tien Yin Wong, and Jiang 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.