- 1 Accurate prediction of glaucoma from color fundus images with a
- 2 convolutional neural network that relies on active and transfer
- 3 learning
- 4
- 5 Authors:
- 6 Ruben Hemelings^{ae*}, MS
- 7 Bart Elene, MS
- 8 João B. Breda^a, MD
- 9 Sophie Lemmens^a, MD
- 10 Maarten Meire^h, MS
- 11 Sayeh Pourjavan^f, MD
- 12 Evelien Vandewalle^{ab}, MD PhD professor
- 13 Sara Van de Veire^g, MD
- 14 Matthew B. Blaschko^c, PhD professor
- 15 Patrick De Boever^{de}, PhD professor
- 16 Ingeborg Stalmans^{ab}, MD PhD professor
- 17
- 18 Affiliations:
- ¹⁹ ^a Research Group Ophthalmology, KU Leuven, Herestraat 49, 3000 Leuven, Belgium
- ^b Ophthalmology Department, UZ Leuven, Herestraat 49, 3000 Leuven, Belgium
- ^c ESAT-PSI, KU Leuven, Kasteelpark Arenberg 10, 3001 Leuven, Belgium
- ^d Hasselt University, Agoralaan building D, 3590 Diepenbeek, Belgium
- ^e VITO NV, Boeretang 200, 2400 Mol, Belgium
- ¹ ^f Chirec Hospitals, Delta-site: Triomflaan 201, 1160 Brussel, Belgium
- ^g AZ Sint-Jan, Ruddershove 10, 8000 Brugge, Belgium
- ²⁶ ^h TC CS-ADVISE, KU Leuven, Kleinhoefstraat 4, 2440 Geel, Belgium
- 27
- 28 *corresponding author
- 29 Contact details
- 30 Affiliation: KU Leuven, VITO
- 31 Postal address: Vito Biologie, Industriezone Vlasmeer 7, 2400 Mol, Belgium
- 32 E-mail: <u>ruben.hemelings@kuleuven.be</u>
- 33 Phone: +32472748707
- 34

35 Abstract

Purpose: To assess the use of deep learning (DL) for computer-assisted glaucoma identification, and the
 impact of training using images selected by an active learning strategy, which minimizes labeling cost.

38 Additionally, this study focuses on the explainability of the glaucoma classifier.

39 Methods: This original investigation pooled 8433 retrospectively collected and anonymized color optic

40 disc-centered fundus images, in order to develop a deep learning-based classifier for glaucoma

41 diagnosis. The labels of the various deep learning models were compared with the clinical assessment by

42 glaucoma experts. Data were analyzed between March and October 2018. Sensitivity, specificity, area

43 under the receiver operating characteristic curve (AUC), and amount of data used for discriminating

44 between glaucomatous and non-glaucomatous fundus images, on both image and patient level.

45 **Results:** Trained using 2072 color fundus images, representing 42% of the original training data, the

trained DL model achieved an AUC of 0.995, sensitivity and specificity of respectively 98.0% (CI 95.5%-

47 99.4%) and 91% (Cl 84.0%–96.0%), for glaucoma versus non-glaucoma patient referral.

48 **Conclusions:** These results demonstrate the benefits of deep learning for automated glaucoma detection

49 based on optic disc-centered fundus images. The combined use of transfer and active learning in the

50 medical community can optimize performance of DL models, while minimizing the labeling cost of

51 domain-specific mavens. Glaucoma experts are able to make use of heat maps generated by the deep

52 learning classifier to assess its decision, which seems to be related to inferior and superior neuroretinal

rim (within ONH), and RNFL in superotemporal and inferotemporal zones (outside ONH).

54 **Key words:** glaucoma detection, deep learning, fundus image, artificial intelligence

55 Introduction

56 Glaucoma is currently responsible for approximately 12% of all cases of irreversible vision loss

57 (Kapetanakis et al. 2006). The number of patients is expected to increase in our ageing society.

58 Predictions indicate that over 110 million people worldwide may be diagnosed with the disease by 2040

59 (Tham et al. 2014). Glaucoma is a neurodegenerative disease characterized by retinal ganglion cell loss

as a result of multiple factors, including high intraocular pressure, optic nerve ocular blood flow

61 dysregulation and neurotoxicity. Progressive optic nerve fiber damage leads to visual field (VF) loss, 62 which often remains unnoticed by the patient because the initial VF loss is peripheral and is compensated 63 by the overlapping VF of the contralateral eye as well as by a compensatory 'filling-in' of these zones by 64 the brain. The resulting lack of early symptoms implies that a significant number of individuals remain 65 undiagnosed, even in high-income countries. Besides VF testing, structural assessment of the optic nerve 66 head (ONH) and retinal nerve fiber layer (RNFL) is crucial in the diagnosis and follow-up of glaucoma. 67 Optical coherence tomography (OCT) and fundus photography are two complementary imaging 68 modalities, with the latter allowing qualitative analysis like disc hemorrhages and color changes. General 69 population screening for glaucoma is currently not common practice (Ervin et al. 2012), as there is no 70 sufficient evidence of its cost-effectiveness to date (Tuulonen 2010; Burr et al. 2014). With the prospect of 71 a growing population affected by glaucoma, a thorough reassessment of glaucoma care is warranted.

72 Ophthalmology is pioneering with future possible application of artificial intelligence (AI) (Ting et al. 2018). 73 Gulshan et al. (2016) developed a convolutional neural network (CNN) for the detection of diabetic 74 retinopathy (DR) from fundus images, scoring areas under the receiver operating characteristic curves 75 (AUCs) of 0.991 and 0.990 on two validation sets. More recently, van der Heijden et al. (2018) reported 76 an AUC of 0.94 on a referral task for DR in a prospective study with nearly 900 patients. This pivotal 77 study led to FDA clearance of the first commercial automated grading tool for referable DR using deep 78 learning. Automated detection of age-related macular degeneration from color fundus images using a 79 pretrained deep learning encoder on the large public AREDS data set was independently described by 80 Burlina et al. (2018) and Grassman et al. (2018).

Automated glaucoma detection from fundus imaging has been actively studied prior to deep learning, with the majority of techniques relying on hand-crafted features, such as the vertical cup-to-disc ratio, extracted from fundus images. Deep learning architectures for glaucoma have been reported on topics including optic disc and cup segmentation (Fu et al. 2018), VF prediction (Wen et al. 2018), and automated glaucoma detection using small data sets (Shibata et al. 2018; Ahn et al. 2018; Muhammad et al. 2017; Asaoka et al. 2016; Maheshwari et al. 2017; Matsopoulos et al. 2008). In 2015, the first results on glaucoma classification with deep learning were published, using two data sets (<2000 images) (Chen
et al. (2015). More recently, Li et al. (2018) described automated glaucoma detection using 48116 fundus
images from an Asian population, reporting high sensitivity (95.6%), specificity (92.0%), and AUC (0.986)
on a validation set of more than 8000 images using pretrained deep learning encoders. The main strength
of their work is the recruitment of a large number of trained ophthalmologists, who graded the entire set of
fundus images for signs of glaucoma.

The current paper reports on the development of a glaucoma prediction model. Optic disc changes are
initially subtle and can be challenging to detect by a human grader. Our access to glaucomatous fundus
images – labeled based on a complete ophthalmologic examination (tonometry, OCT or confocal
scanning laser ophthalmoscopy) – allows the deep learning encoder to learn subtle features in fundus
images of early/moderate stage glaucoma patients. Hence, the first objective of this study was to develop
and validate a deep learned glaucoma classifier using color fundus images from a patient population,
measured against clinical diagnosis.

The second objective was to explore the added-value of active learning (Settles, 2009) on top of deep learning for automated glaucoma detection. Active learning is a special case of semi-supervised learning that aims to leverage uncertainty information from an unlabeled set in order to predict from which unlabeled images the classifier would benefit the most if they would become labeled. True labels, especially in the medical community, can be difficult to obtain. By employing an active learning system that maximizes classification performance, while minimizing the number of required labels, data sets and labeling efforts can be used more efficiently.

107 The third and final objective was to inspect the trained model's decision process using interpretable heat 108 maps. DL models learn concepts from the data itself, omitting the need for manual feature extraction, and 109 leading to state-of-the-art results, but lower transparency in understanding the classifier's decision 110 process. Heat maps that visualize the image areas that contributed the most towards glaucoma 111 classification might assist in opening the black box of the trained deep learning system.

112 Methods

113 Image and Label Acquisition

All 30° optic disc-centered color fundus images of 1620x1444 resolution were captured with a Zeiss 114 115 VISUCAM (Carl Zeiss Meditec, Jena, Germany) and used retrospectively in the current study. The 116 glaucomatous fundus images (6651) originate from 1353 unique patients (±4.9 images per patient) 117 imaged at the glaucoma clinic of the University Hospitals Leuven (Belgium) during several consultations 118 between 2009 and 2017. Over 60% of patients went to follow-up consultations, leading to images taken at 119 different points in time, which can differ due disease progression, hence useful for the model. The vast 120 majority of fundus images (1614) from 403 non-glaucoma (normal) individuals (±4 images per individual) 121 stem from a data set collected at three different locations in the context of an awareness campaign during 122 the World Glaucoma Week 2018 that took place between March 11th and 17th. Screening sessions were 123 organized at different Belgian hospitals (Brussels, Leuven and Bruges), and were aimed at raising public 124 awareness on the disease, with eligible participants restricted to age 40 or above. The images of the 125 healthy subjects at the screening sessions were taken at the same time, and show no signs of retinal 126 changes. However, they do hold additional information because of small changes due to focus, lighting, 127 eve movement etc. Additionally, a set of normal fundus images (168) from 88 individuals (±1.9 images per 128 individual, both eyes when applicable) were sourced from a 2016 glaucoma screening program at the 129 University Hospitals Leuven. This resulted in a total set of 1782 images of 491 non-glaucoma individuals. 130 For all images, information provided to data processor was limited to an anonymized patient identifier and 131 glaucoma type. The glaucoma diagnoses (following ICD standards) linked with the fundus images were 132 obtained through a full ophthalmologic examination. Patients were subjected to neuroretinal rim and nerve fiber layer analysis using either OCT (Spectralis OCT; Heidelberg Engineering, Heidelberg, 133 134 Germany) or confocal scanning laser ophthalmoscopy (Heidelberg retinal tomography; HRT; Heidelberg 135 Engineering, Heidelberg, Germany), tonometry (Goldmann Applanation Tonometry; Haag-Streit AT900; 136 Köniz, Switzerland) and visual field testing (Humphrey Visual Field Analyzer; Carl Zeiss Meditec, Jena, 137 Germany). The glaucoma experts are aware of the so-called red and green disease surrounding OCT 138 results and do verify the actual images to look for any artifacts or other sources of misinterpretation and 139 check the reliability of the analysis. The transition HRT to Spectralis OCT device rolled out in 2015.

140 Patients that were followed up prior to the switch are still imaged with HRT to ensure consistent 141 progression analysis. Visual field testing at the glaucoma clinic of UZ Leuven is achieved through 142 Humphrey or Octopus standard automated perimetry. Clinicians look for typical glaucomatous visual field 143 defects such as wedge shape defects, steps or nasal breakthrough. The glaucoma experts at UZ Leuven 144 incorporated progression analysis when available, to ensure accurate glaucoma diagnosis. The images 145 sourced from the screening program were evaluated by two glaucoma experts without the aid of OCT and 146 VF tests, but did include a slit lamp biomicroscopic examination including fundoscopy by a glaucoma 147 expert.

148 Image Preprocessing

149 All 8433 images were manually inspected by two independent retinal image experts to control for quality, 150 omitting images without visible optic disc. Because of the high-quality glaucoma labels based on a full 151 ophthalmological exam, even poor images can be used during training, to increase the robustness of the 152 deep learning model. This guality control does not match the guality that human experts require for 153 diagnosis, hence the task being carried out by retinal image experts with experience in deep learning in 154 ophthalmology. Quality assessed images deemed fit for analysis were initially center cropped to a square 155 of 1016x1016, removing any risk of influence caused by the image border, and subsequently resized to 156 224x224 to match the input layer of the ResNet-50 (He et al. 2016) neural network architecture.

157 Color fundus images are characterized by a large intra-image variance in intensity levels mainly due to

the curvature of the retina. Therefore, the images were convolved with a Gaussian kernel (30x30) to

159 estimate its background, and this was deducted from the original image. The result is a data set of

standardized fundus images, as illustrated in the top left of Figure 1.

161 In this study, data augmentation was implemented to artificially increase the number of original images

used to train the CNN. Augmentation techniques included in the training process of the final model were:

163 horizontal flip, brightness shift, and minor elastic deformation. All image augmentations were randomly

generated at the start of each mini-batch, as can be seen in the top right of Figure 1.

165 Transfer Learning

This study used the publicly available Keras (v2.2.0, TensorFlow v1.4.1 backend) ResNet-50 encoder pretrained on ImageNet (Deng et al. 2009), followed by additional layers to increase regularization. The complete deep neural network counted 182 layers of mathematical operations including convolutions and batch normalization (see supplementary material for full network details). During training, all pretrained encoder layers were frozen, except for the last 12 layers, to allow the model to learn features relevant for glaucoma detection. Standard binary cross entropy was used as cost function, and the Adam (Kingma & Ba, 2015) optimizer was used with a constant learning rate of 0.0001.

173 Active Learning

The employed ResNet-50 encoder features over 25 million parameters, requiring a high amount of unique training data to reach its full potential. This study opted for uncertainty sampling as the active learning criterion because of its widespread application in image classification (Joshi et al. 2009).

177 Uncertainty sampling refers to selecting new samples based on their close distance to the decision

boundary set by the classification system, which corresponds to a higher uncertainty. By querying these

179 labels first, the classifier is expected to reduce its uncertainty on these data, more quickly converging to a

180 stable solution. To benchmark the performance of this heuristic, this study also conducted an experiment

181 in which data to be labeled are sampled at random (see Figure 1 and supplementary material for

182 sampling details).

183 Saliency maps

184 The Keras Visualization Toolkit (Kotikalapudi, 2017) was used to generate saliency maps. Saliency maps

185 for deep learning accentuate the pixels (colored reddish) that contribute the most to the classification

186 output, i.e. if that pixel were to change, the classification output would be likely to change as well

187 (Simonyan, Vedaldi & Zisserman, 2014). The generated saliency maps of (1) randomly selected images

188 classified correctly by the trained model, and (2) the false positives (FP) and false negatives (FN) were

189 subsequently examined by two blinded glaucoma experts.

- In order to reveal a pattern, saliency maps of thirty oculus dextrus fundus images were manually aligned
 and averaged. The average saliency map was divided into six zones commonly used in ONH analysis,
- 192 with differences in saliency intensity quantified.

193 Evaluation metrics

All predictions by the deep learning models were evaluated against the ground truth label provided by the
University Hospitals Leuven. AUC was selected as main performance metric, with specificity and
sensitivity also reported. The evaluation phase was conducted using the SciPy Python library (Jones,
Oliphant & Peterson, 2001).

198 Results

A total number of 7038 images (83.5% of originally pooled number of images) of 1775 patients passed the manual image quality assessment and were further used in this study. Selected images of 1775 patients were allocated to training (70%; 1244 patients; 4935 images), validation (10%; 177 patients, 679 images), and test set (20%; 354 patients, 1424 images), based on anonymized patient identifier, ensuring that all images from the same patient were to be found in the same class. All glaucoma detection experiments were evaluated on the validation set of 679 images as proxy to select the optimal state of trainable parameters.

206 Final results are reported on the independent test set of 1424 unique images, corresponding to 354 207 individuals. For patient level prediction, all glaucoma predictions of images belonging to the same patient 208 are averaged, and then classified based on the 0.5 cut-off. Results on the patient level were considered 209 more appropriate for interpretation of the results, as referral decisions would be made on the patient level. 210 Table 1 outlines classification results for glaucoma detection (glaucoma vs non-glaucoma, abbreviated by 211 GLC and NO) for the active learning experiments and baseline model with all training images and labels 212 included at start of the training process. Confusion matrix and performance metrics are given, computed 213 over the original image with test-time augmentation (TTA). The latter corresponds to randomly 214 augmenting the image tenfold, using the same techniques as in training, followed by averaging the 215 prediction probabilities in order to decrease prediction uncertainty. The use of TTA led to reductions in 216 AUC error up to 14%.

Final models for the two active learning experiments were selected at 2072 training images, due to the marginal improvements when using additional data (Figure 1, graph bottom right). After seeing 42% (2072 images) of the training data, the model following the active learning strategy achieved an AUC of 0.995, with sensitivity at 98% and specificity at 91% on the test set, clearly benefitting from the employed heuristic that leveraged uncertainty information (Table 1). The performance gap is the most prominent when comparing the specificity of both models, with the random sampling technique yielding a modest 84% on patient level.

224 The baseline model trained with all original 4935 training images (accompanied by a large set of artificial images following data augmentation) obtained an AUC of 0.996 on patient referral level. Sensitivity and 225 226 specificity reach 99.2% and 93%, respectively, corresponding to a low number of false negatives (2) and 227 false positives (7). The grouping of images at the patient level led to a reduction in misclassification. 228 Images of misclassified patients were reviewed by two ophthalmologists specialized in glaucoma. False 229 positives could be grouped into (1) subpar image quality due to blurriness or artefacts like eye lashes (n = 230 4) and (2) signs of other ocular diseases like macular drusen (n = 1) and (3) peripapillarly atrophy (n = 2). 231 The fundus image of one false negative patient did not display any clear signs of glaucoma onset, while 232 the other one was a true FN.

The saliency analysis, aimed at explaining the classifier's decision process, is given in Figure 2. Careful analysis of over 500 saliency maps by two glaucoma experts revealed a recurrent pattern of elevated saliency in inferotemporal and superotemporal zones, either within (early/moderate stage, remaining neuroretinal rim) or outside (late stage, complete thinning) the ONH. This recurrent pattern was subsequently confirmed through the averaging of thirty optic disc-aligned saliency maps.

238 Discussion

This study resulted in an accurate deep learning based glaucoma classifier, achieving patient referral AUC of 0.995 on 1424 test images from 354 individuals, with only 42% (2072 images) of the complete training set (4935 images) used. The joint forces of transfer and active learning foster potential in the 242 domain of glaucoma classification from fundus images, allowing model training with a 58% reduction in
243 labeling requirements.

244 The development of a baseline model trained with all available training data (4935 images) and transfer 245 learning yielded an AUC of 0.996, sensitivity and specificity of 99.2% and 93%, on the test set. The merits 246 of transfer learning in the field of automated glaucoma detection using fundus images have been 247 illustrated using both small (Shibata et al. 2018; Ahn et al. 2018) and large (Li et al. 2018; Christopher et 248 al. 2018) (>5000 images) data sets. Li et al. (2009) trained a CNN for glaucoma classification using a data 249 set of 48116 images, reporting AUC, sensitivity and specificity of 0.986, 95.6% and 92%, respectively. 250 While the efforts to reach a labeled data set of this size are to be commended, one could question 251 whether the same performance can be reached in a more cost-effective manner, with significantly less 252 labeled fundus images used during training. Annotated medical image data are hard to gather, with images and associated glaucoma diagnosis employed in this study generated over several years. The 253 254 field of active learning encompasses a set of techniques that accelerate training by querying experts for 255 labels that would benefit the classification system the most. In this study, the addition of an active learning 256 component resulted in a model with 42% of the data used, while still attaining an AUC of 0.995 on patient 257 referral.

258 Two trained glaucoma clinicians analyzed more than 500 saliency maps, accompanied by the original 259 glaucomatous images, and indicate a recurrent pattern of salient regions in the inferotemporal and 260 superotemporal zones neighboring the ONH. These regions likely correspond to the RNFL areas that are 261 affected as a result of glaucoma. The hypothesis of a recurrent pattern of elevated saliency in 262 inferotemporal and superotemporal regions was supported by a statistical analysis using the manually 263 aligned average of 30 randomly selected saliency maps (Figure 2). The center part (disc area) of the 264 average saliency map provides additional evidence on a significant concentration of salient regions in the 265 inferior, temporal and superotemporal region of the ONH. The latter partly matches the findings described 266 by Christopher et al. (2018), who used an occlusion-based strategy to reveal salient regions in inferior 267 and superior zones within the disc. This study is the first to indicate that regions outside the ONH could

be valuable in glaucoma classification using deep learning. We aim to further investigate the importance
 of RNFL defects in glaucoma classification from fundus images in future work.

Manual image quality assessment led to 83.5% of available fundus images being actually of sufficient quality for analysis in this original investigation. Two retinal image experts graded each image, omitting those with an excessive presence of camera artefacts or missing optic nerve head. Image quality is essential to ensure proper functioning of the convolutional neural network. In this study, the latter is backed up by the analysis of false positives and false negatives by two ophthalmologists (performed in a blind manner), who indicated subpar image quality to be the culprit in several cases.

276 This study has several limitations. The class distribution, with over 70% glaucomatous images, is far from 277 the real-life prevalence one would encounter at screening sessions. The selected data imbalance is due 278 to the small availability of non-glaucomatous images, which are often not stored in hospitals. In addition, a 279 large set of the glaucoma images are intermediate or late stage (based on neuroretinal rim assessment), 280 while an important application of glaucoma classification with deep learning could be early detection. 281 Finally, the models trained and validated in this study used images of mainly Caucasian patients that 282 were captured with a fundus camera device from one vendor. To overcome this limitation, we are extending our work by validating and refining our current model using heterogenous data sets obtained 283 284 through international collaborations, with the goal to develop a model suitable for global screening.

285 Conclusions

286 This study achieves state-of-the-art results for automated glaucoma referral with a 60% decrease in 287 labeling cost through the combination of transfer learning, careful data augmentation, and uncertainty 288 sampling, a heuristic commonly used in the domain of active learning. Our iterative sampling process 289 provides novel evidence that deep learning can achieve excellent performance in glaucoma classification, 290 even when using a limited amount of labeled training data. These findings should motivate research 291 groups that have access to less data to help to advance the field of artificial intelligence applied to 292 ophthalmology. Finally, this study provides novel insights into the decision-making process of the trained 293 deep learning glaucoma classifier through the averaging of saliency maps, which seem to be highlighting

294	inferior and superior neuroretinal rim thinning (within ONH) as well as RNFL defects in superotemporal
295	and inferotemporal zones (outside ONH).
296	
297	
298	
299	
300	
301	
302	
303	
304	
305	
306	
307	
308	
309	
310	
311	
312	
313	

314 Tables

Table 1. Glaucoma detection with transfer and active learning – quantitative results

Model		Confusion matrix			Metrics			
			predicted					
		true		NO	GLC	Sensitivity	96% 87% 0.983	
			NO	324	47	Specificity		
Uncertainty sampling			GLC	42	1011	100	0.000	
(2072 training images; 42% of training set)			predicted					
		true		NO	GLC	Sensitivity	98%	
			NO	91	9	AUC	91% 0.995	
			GLC	5	249			
			predicted					
Random sampling		true		NO	GLC	Sensitivity	96%	
			NO	305	66	AUC	81% 0.972	
			GLC	46	1007			
training set)			predicted				0.001	
		true		NO	GLC	Sensitivity	98% 84% 0.986	
			NO	84	16	AUC		
			GLC	5	249	,	0.000	
	Image		predicted		Consitivity	069/		
		true		NO	GLC	Specificity	90%	
			NO	346	25	AUC	0.986	
Baseline ResNet-50 CNN			GLC	42	1011			
training set)	L 1		predicted				000/	
		true	ē		NO	GLC	Sensitivity	99%
			NO	93	7	ALIC	0 996	
			GLC	2	252	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	0.000	

324 Figures

- **Figure 1.** Top: Overview of data used, effect of image quality control, and subdivision in training,
- 326 validation and test set.



- Figure 2. Top: Overview of image preprocessing (ROI extraction, background subtraction) and data
- augmentation (horizontal flip, elastic deformation and brightness shift). Bottom: Overview of active
- learning process. 1: 14 preprocessed and augmented fundus images were used to finetune a CNN with
- pretrained ImageNet weights. 2: After convergence (no improvement of validation accuracy for two
- epochs), the model was validated on 679 images, with the results of each active learning iteration
- visualized in 3. 4: The model was also evaluated on an unlabeled set of (non-)glaucomatous images, with
- the 14 most uncertain samples (or random samples) transferred to the training set (5). This process was
- repeated until the unlabeled set was depleted.



Figure 3. A: Average saliency map of thirty aligned oculus dextrus images split into six sectors, with reddish color corresponding to high saliency (important area for glaucoma classification); optic disc is contained inside the inner circle. B: Quantification of average saliency in the six zones outside the disc area (complementary to A).



Acknowledgements 379

380 The first author is jointly supported by the Research Group Ophthalmology, KU Leuven and VITO NV. No 381 outside entities have been involved in the study design, in the collection, analysis and interpretation of 382 data, in the writing of the manuscript, nor in the decision to submit the manuscript for publication. Thus, the authors declare that there are no conflicts of interest in this work.

- 383
- 384

References 385

Ahn JM, Kim S, Ahn KS, Cho SH, Lee KB, Kim US (2018): A deep learning model for the detection of 386 387 both advanced and early glaucoma using fundus photography. PloS one **13**(11):e0207982. 388 doi:10.1371/journal.pone.0207982

389

390 Asaoka R, Murata H, Iwase A, Araie M (2016): Detecting Preperimetric Glaucoma with Standard 391 Automated Perimetry Using a Deep Learning Classifier. Ophthalmology 123(9):1974–1980. 392 doi: 10.1016/j.ophtha.2016.05.029.

393

394 Burlina PM, Joshi N, Pacheco KD, Freund DE, Kong J, Bressler NM (2018): Use of Deep Learning for Detailed Severity Characterization and Estimation of 5-Year Risk Among Patients With Age-Related 395

396 Macular Degeneration. JAMA Ophthalmol, Published online September 14, 2018.

- 397 doi:10.1001/jamaophthalmol.2018.4118
- 398 399 Burr J, Hernández R, Ramsay C, et al. (2014): Is it worthwhile to conduct a randomized controlled trial of 400 glaucoma screening in the United Kingdom? J Health Serv Res Policy, 19(1):42-51.
- 401 https://doi.org/10.1177/1355819613499748
- 402

403 Chen X, Xu Y, Kee Wong DW, Wong TY, Liu J (2015): Glaucoma detection based on deep convolutional 404 neural network. 37th Annual International Conference of the IEEE Engineering in Medicine and Biology 405 Society (EMBC), Milan 715-718. doi: 10.1109/EMBC.2015.7318462

406

407 Christopher M, Belghith A, Bowd C, et al. (2018): Performance of Deep Learning Architectures and 408 Transfer Learning for Detecting Glaucomatous Optic Neuropathy in Fundus Photographs. Sci Rep 8(1): 409 16685. doi:10.1038/s41598-018-35044-9 410

411 Deng J, Dong W, Socher R, Li L-J, Li K, Fei-Fei L (2009): ImageNet: A Large-Scale Hierarchical Image Database. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Miami, FL 248-255. 412 413 doi: 10.1109/CVPR.2009.5206848

414

415 Ervin AM, Boland MV, Myrowitz EH, et al. (2012): Screening for Glaucoma: Comparative Effectiveness. 416 Rockville (MD): Agency for Healthcare Research and Quality (US). (Comparative Effectiveness Reviews. 417 No. 59.)

418

419 Fu H, Cheng J, Xu Y, Wong DWK, Liu J and Cao X (2018): Joint Optic Disc and Cup Segmentation 420 Based on Multi-Label Deep Network and Polar Transformation. IEEE Transactions on Medical Imaging

- 421 37(7): 1597-1605. doi: 10.1109/TMI.2018.2791488
- 422

Grassmann F, Mengelkamp J, Brandl C, et al. (2018): A Deep Learning Algorithm for Prediction of Age-423 424 Related Eye Disease Study Severity Scale for Age-Related Macular Degeneration from Color Fundus 425 Photography. Ophthalmology 125(9): 1410-1420. doi:10.1016/j.ophtha.2018.02.037.

- 426 Gulshan V. Peng L. Coram M. et al. (2016); Development and Validation of a Deep Learning Algorithm for 427
- 428 Detection of Diabetic Retinopathy in Retinal Fundus Photographs. JAMA 316(22): 2402-2410.

429 doi:10.1001/jama.2016.17216

431 He K, Zhang X, Ren S, Sun J (2016): Deep residual learning for image recognition. Proceedings of the 432 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Piscataway, NJ 771–778. 433 434 Jones E, Oliphant E, Peterson P (2001): SciPy: Open Source Scientific Tools for Python. 435 http://www.scipy.org/ 436 Joshi AJ, Porikli F, Papanikolopoulos N (2009): Multi-class active learning for image classification. IEEE 437 Conference on Computer Vision and Pattern Recognition (CVPR), Miami, FL, 2372-2379. 438 doi: 10.1109/CVPR.2009.5206627 439 440 Kapetanakis VV, Chan MPY, Foster PJ, et al. (2006): Global variations and time trends in the 441 prevalence of primary open angle glaucoma (POAG): a systematic review and meta-analysis. Br 442 J Ophthalmol 100: 86-93. 443 444 Kingma DP, Ba J. Adam (2015): A method for stochastic optimization. International Conference on 445 Learning Representations (ICLR). 446 447 Kotikalapudi R (2017): keras-vis. https://github.com/raghakot/keras-vis. 448 449 Li Z, He Y, Keel S, Meng W, Chang R, He M (2018): Efficacy of a Deep Learning System for Detecting 450 Glaucomatous Optic Neuropathy Based on Color Fundus Photographs. Ophthalmology 125(8): 1199-451 1206. doi:10.1016/j.ophtha.2018.01.023 452 453 Li Z, Keel S, Liu C, He M. (2018): Can Artificial Intelligence Make Screening Faster, More Accurate, and 454 More Accessible? Asia Pac J Ophthalmol (Phila) 7: 436-441. 455 456 Maheshwari S, Pachori RB, Acharya UR (2017): Automated Diagnosis of Glaucoma Using Empirical 457 Wavelet Transform and Correntropy Features Extracted From Fundus Images. IEEE Journal of 458 Biomedical and Health Informatics 21(3): 803-813. doi: 10.1109/JBHI.2016.2544961 459 460 Matsopoulos GK, Asvestas PA, Delibasis KK, Mouravliansky NA, Zeyen TG (2008): Detection of 461 glaucomatous change based on vessel shape analysis. Comput Med Imaging Graph 32(3): 183-192. doi: 462 10.1016/j.compmedimag.2007.11.003. 463 464 Muhammad H, Fuchs TJ, De Cuir N, et al. (2017): Hybrid Deep Learning on Single Wide-field Optical 465 Coherence tomography Scans Accurately Classifies Glaucoma Suspects. J Glaucoma 26(12): 1086-1094. doi: 10.1097/IJG.000000000000765 466 467 468 Settles, B (2009): Active Learning Literature Survey. Computer Sciences Technical Report 1648, 469 University of Wisconsin-Madison. 470 471 Shibata N, Tanito M, Mitsuhashi K, et al. (2018): Development of a deep residual learning algorithm to 472 screen for glaucoma from fundus photography. Sci Rep 8(1): 14665. doi:10.1038/s41598-018-33013-w 473 474 Simonyan K, Vedaldi A, Zisserman A (2014): Deep Inside Convolutional Networks: Visualising Image 475 Classification Models and Saliency Maps. Proceedings of the 2014 International Conference on Learning 476 Representations (ICLR). 477 478 Tham Y, Li X, Wong T, Quigley H, Aung T, Cheng C (2014): Global Prevalence of Glaucoma and 479 Projections of Glaucoma Burden through 2040. Ophthalmology 121(11): 2081-2090. 480 doi:10.1016/j.ophtha.2014.05.013 481 482 Ting DSW, Pasquale LR, Peng L, et al. (2018): Artificial intelligence and deep learning in ophthalmology. 483 Br J Ophthalmol, Published Online First: 25 October 2018. doi: 10.1136/bjophthalmol-2018-313173

484

Tuulonen A. Cost-effectiveness of screening for open angle glaucoma in developed countries. Indian J
 Ophthalmol. 2011;59 Suppl(Suppl1):S24–S30. doi:10.4103/0301-4738.73684

487

van der Heijden AA, Abramoff MD, Verbraak F, Hecke MV, Liem A and Nijpels G (2018): Validation of

automated screening for referable diabetic retinopathy with the IDx-DR device in the Hoorn Diabetes

490 Care System. Acta Ophthalmol **96**: 63–68. doi:10.1111/aos.13613 491

492 Wen JC, Lee CS, Keane PA, et al. (2018). Forecasting Future Humphrey Visual Fields Using Deep

493 Learning. arXiv e-prints.