

1 Establishing a Generalized Deep Learning System for Detection of  
2 Glaucomatous Optic Neuropathy using Fundus Photographs

3  
4 Hanruo Liu<sup>1</sup>, MD, PhD; Liu Li<sup>2</sup>, BEng; I. Michael Wormstone<sup>3</sup>, PhD; Chunyan Qiao<sup>1</sup>, MD,  
5 PhD; Chun Zhang<sup>4</sup>, MD, PhD; Ping Liu<sup>5</sup>, MD, PhD; Shuning Li<sup>1</sup>, MD, PhD; Huaizhou  
6 Wang<sup>1</sup>, MD, PhD; Dapeng Mou<sup>1</sup>, MD, PhD; Ruiqi Pang<sup>1</sup>, MD; Diya Yang<sup>1</sup>, MD, PhD; Lai  
7 Jiang<sup>2</sup>, BEng; Yihan Chen<sup>1</sup>, MD; Man Hu<sup>6</sup>, MD, PhD; Yongli Xu<sup>7</sup>, PhD; Hong Kang<sup>8</sup>, PhD;  
8 Xin Ji<sup>9</sup>, BEng; Robert Chang<sup>10</sup>, MD, PhD; Clement Tham<sup>11</sup>, MD, PhD; Carol Cheung<sup>11</sup>, PhD;  
9 Daniel Shu Wei Ting<sup>12</sup>, MD, PhD; Tien Yin Wong<sup>12</sup>, MD, PhD; Zulin Wang<sup>2</sup>, PhD; Robert  
10 N. Weinreb<sup>13</sup>, MD, PhD; Mai Xu<sup>2\*</sup>, PhD; Ningli Wang<sup>1\*</sup>, MD, PhD

- 11 1. Beijing Institute of Ophthalmology, Beijing Tongren Hospital, Capital Medical University,  
12 Beijing Ophthalmology & Visual Science Key Lab, Beijing, China.
- 13 2. School of Electronic and Information Engineering, Beihang University, Beijing 100191,  
14 China.
- 15 3. School of Biological Sciences, University of East Anglia, Norwich Research Park, Norwich  
16 NR4 7TJ, UK
- 17 4. Department of Ophthalmology, Peking University Third Hospital, Beijing, China
- 18 5. Ophthalmology Hospital, First Hospital of Harbin Medical University, Harbin 150001,  
19 Heilongjiang, China
- 20 6. National Key Discipline of Pediatrics, Ministry of Education, Department of Ophthalmology,  
21 Beijing Children's Hospital, Capital Medical University, Beijing, China
- 22 7. Department of Mathematics, Beijing University of Chemical Technology, Beijing, China
- 23 8. Nankai University, the College of Computer Science, Tianjin, China
- 24 9. Beijing Shanggong Medical Technology co., Ltd, Beijing, China

- 25 10. Department of Ophthalmology, Byers Eye Institute at Stanford University, Palo Alto, CA, US
- 26 11. Department of Ophthalmology & Visual Sciences, Faculty of Medicine, The Chinese
- 27 University of Hong Kong, 147K Argyle Street, Kowloon, Hong Kong, China
- 28 12. Singapore National Eye Center, Singapore Eye Research Institute, Duke-National University
- 29 of Singapore Graduate School of Medicine, National University of Singapore and National
- 30 University, Singapore 168751, Singapore
- 31 13. Shiley Eye Institute, University of California San Diego, 9500 Gilman Drive, MC 0946, La
- 32 Jolla, CA 92093, US

33 \*These authors contributed equally to this work.

34 Correspondence should be addressed to Ningli Wang (email: wningli@vip.163.com), or Mai Xu.

35 (maixu@buaa.edu.cn).

36 Running head: Deep Learning System for Detection of Glaucomatous Optic Neuropathy

37 Number of Tables: 3

38 Number of Figures: 2

39 Word count (total): Abstract: 338 Text: 2996 words (excluding abstract, references and tables), Key

40 Points 100 words

41 **Key Points**

42

43 Question:

44 How does a deep learning system (DLS) compare with professional human graders in  
45 detecting glaucomatous optic neuropathy (GON)?

46

47 Findings:

48 The DLS showed a sensitivity of 96.2% and specificity of 97.7% for detecting GON in a  
49 local validation dataset; 93.6-96.1% sensitivity and 95.6-97.1% specificity in three clinical-  
50 based datasets; 91.0% sensitivity and 92.6% specificity in a real-world distribution dataset;  
51 87.7% sensitivity and 80.8% specificity in a multi-ethnic dataset; 82.2% sensitivity and 70.4%  
52 specificity in a website-based dataset.

53

54 Meaning:

55 This assessment of fundus images suggests DLS can provide a tool with high sensitivity,  
56 specificity and might expedite screening for GON.

57 **Abstract:**

58 **IMPORTANCE** A deep learning system (DLS) that could automatically detect  
59 glaucomatous optic neuropathy (GON) with high sensitivity and specificity might expedite  
60 screening for GON.

61 **OBJECTIVE** To establish a DLS for detection of GON using retinal fundus images and  
62 convoluted neural networks (GD-CNN) that has the ability to be generalized across  
63 populations.

64 **DESIGN, SETTING, AND PARTICIPANTS** A DLS for the classification of GON was  
65 developed for automated classification of GON using retinal fundus images. To build and  
66 validate GD-CNN, a total of 355 339 fundus images were included. Of those, 241 032 images  
67 and 114 307 images were selected as the training and validation dataset, respectively. The  
68 generalization of the DLS was tested in several validation datasets, which allowed assessment  
69 of the DLS in a clinical setting without exclusions, testing against variable image quality  
70 based on fundus photographs obtained from websites, evaluation in a population-based study  
71 that reflects a natural distribution of glaucoma patients within the cohort and an additive  
72 dataset that has a diverse ethnic distribution. An online learning system was established to  
73 transfer the trained and validated DLS to generalize the results with fundus images from new  
74 sources. To better understand the DLS decision making process, a prediction visualization  
75 test was performed that identified regions of the fundus images utilized by the DLS for  
76 diagnosis.

77 **EXPOSURES** Use of a deep learning system.

78 **MAIN OUTCOMES AND MEASURES** Area under the receiver operating characteristics  
79 curve (AUC), sensitivity and specificity for DLS with reference to professional graders.

80

81 **RESULTS** The AUC of the GD-CNN model in primary local validation datasets was 0.996  
82 (95% CI, 0.995-0.998), with sensitivity of 0.962, and specificity of 0.977. The most common  
83 reasons for both false-negative and false-positive grading by GD-CNN (46.3% and 32.3%)  
84 and manual grading (44.2% and 34.0%) was pathologic or high myopia.

85

86 **CONCLUSIONS AND RELEVANCE** Application of GD-CNN to fundus images from  
87 different settings and varying image quality demonstrated a high sensitivity, specificity and  
88 generalization for detecting GON. These findings suggest automated DLS might enhance  
89 current screening programs in a cost-effective and time-efficient manner.

90 Glaucoma is the leading cause of irreversible blindness.<sup>1</sup> It is predicted to affect 80 million  
91 people worldwide by 2020 and 111.8 million by 2040.<sup>2</sup> Glaucoma is a chronic  
92 neurodegenerative disease of the eye.<sup>3</sup> The majority of glaucoma patients are unaware of  
93 their condition until late in the course of their disease, when central visual acuity is affected.<sup>4</sup>  
94 Screening and early detection of glaucoma, along with timely referral and treatment, is a  
95 generally accepted strategy for preventing vision loss.<sup>5</sup> Digital fundus image evaluation has  
96 emerged as a modality for large-scale glaucoma screening due to convenience and relative  
97 affordability.<sup>6,7</sup> Nevertheless, this process of manual image assessment is labor intensive and  
98 time-consuming.<sup>8</sup> In addition, glaucoma diagnosis from fundus images is subjective, and  
99 efficiency is likely linked to the experience and skill of the observer.

100

101 Artificial intelligence has been successfully applied in image-based medical diagnoses, such  
102 as skin cancer, breast cancer, brain tumors and diabetic retinopathy.<sup>9-13</sup> The deep learning  
103 system (DLS) approach also has recently been adopted to provide high sensitivity and  
104 specificity (>90%) for detecting glaucomatous optic neuropathy (GON) from high-quality  
105 retinal fundus images.<sup>14</sup> However, the use of DLS for medical diagnosis has inferior  
106 performance when applied to data obtained from different sources.<sup>12,14</sup> This is an important  
107 consideration, as ideally a DLS would need to be generally utilized in different settings in  
108 which the images will be of varying quality, ethnicity and population sources if maximum  
109 reach and clinical benefit is to be achieved.<sup>15-17</sup>

110

111 In this study, we established a large-scale database of fundus images for glaucoma diagnosis  
112 ('FIGD' database) and developed from fundus images Glaucoma Diagnosis with Convolved  
113 Neural Networks (GD-CNN), as an advanced DLS approach for automatically detecting  
114 GON that has the ability to be generalized across populations.

115

## 116 **Methods**

### 117 **Training datasets**

118 The study was conducted according to the tenets of the Declaration of Helsinki and it was  
119 approved by the institutional review board (IRB) of Beijing Tongren Hospital (identifier,  
120 TRECKY2018-034). As the study was a retrospective review and analysis of fully  
121 anonymized colour retinal fundus images, the medical ethics committee exempted the need  
122 for the patients' informed consent.

123

124 To establish an automatic diagnosis system for GON, a total of 274 413 fundus images were  
125 obtained from the Chinese Glaucoma Study Alliance (CGSA, Appendix 1.1 available online  
126 at [www.aojournal.org](http://www.aojournal.org)) between 2009 and 2017 (Table 1). The CGSA uses a tele-  
127 ophthalmology platform and a cloud-based online dataset (<http://www.funduspace.com>  
128 Accessed May 2017), which has established its own electronic data capture system to achieve  
129 effective data quality control. For each patient, two fundus images of each eye were recorded.  
130 For this study, each image in the training dataset was subjected to a tiered grading system  
131 consisting of multiple layers of trained graders of increasing expertise. Each image imported  
132 into the database started with a label matching the most recent diagnosis of the patient. The  
133 first tier of graders consisted of five trained medical students and non-medical  
134 undergraduates. They conducted initial quality control according to the following rules: 1) the  
135 image did not contain severe resolution reductions or significant artifacts; 2) the image field  
136 included the entire optic nerve head and macula; 3) the illumination was acceptable i.e. not  
137 too dark or too light; 4) the image was focused sufficiently for grading the optic nerve head  
138 and retinal nerve fiber layer (RNFL). The second tier of graders consisted of twenty-two  
139 Chinese board-certified ophthalmologists or postgraduate ophthalmology trainees (>2 years'

140 experience) who had passed a pre-training test. In the process of grading, each image was  
141 assigned randomly to two ophthalmologists for grading. Each grader independently graded  
142 and recorded each image according to the criteria of GON (Table 2). The third tier of graders  
143 consisted of two senior independent glaucoma specialists (>10 years of experience with  
144 glaucoma diagnosis); they were consulted to adjudicate disagreement in tier 2 grading  
145 (Appendix 1.2, available online at [www.aaojournal.org](http://www.aaojournal.org)). Following this process images  
146 were classified as unlikely, probable, and definite GON. Referable GON was defined as  
147 probable or definite GON.

148

#### 149 **GD-CNN Model**

150 The training images with assigned labels were utilized to establish a state-of-the-art DLS,  
151 GD-CNN, based on the Residual Net (ResNet) platform.<sup>18</sup> (eFigures 1 & 2, Appendix 2.0 -  
152 [www.aaojournal.org](http://www.aaojournal.org)). In the current study, we restricted the analysis to the binary  
153 classification problem of glaucoma in fundus images. The basic operation of ResNet is to  
154 apply convolution repeatedly, which is computationally quite expensive for high-resolution  
155 images. Therefore, we pre-process images by down-sampling them to 224×224 pixel  
156 resolution. In addition, these images were centered on the optic cup and contained part of the  
157 surrounding vessels, as glaucoma is highly correlated with alteration in these regions.<sup>19</sup> To  
158 achieve this, the optic cups were automatically detected by recognition of the area with the  
159 highest intensity on the grayscale map of each fundus image; this was found to consistently  
160 be associated with the optic cup. Next, we calculate the mean values of red, green and blue  
161 (RGB) channels, respectively, among all the fundus images in the training dataset. Then, for  
162 each sample, we remove the three mean values on RGB channels, such that the input to GD-  
163 CNN is around 0 for relieving the over-fitting issue.<sup>20</sup> As such, the redundancy of the fundus  
164 image can be removed for the binary classification of glaucoma in GD-CNN. Since the GON



165 diagnosis was formulated as a binary classification problem, predicting whether GON was  
166 positive or negative, a cross-entropy function was applied in GD-CNN as the loss function.  
167 For each parameter assessed, GD-CNN was trained to minimize the cross-entropy loss over  
168 the large-scale training samples of positive and negative GON. The minimization was  
169 achieved through the back-propagation algorithm with the stochastic gradient descent  
170 optimizer. Once training of GD-CNN was established, the system was applied to validation  
171 sets.

172

### 173 **Validation datasets**

174 Details of all validation datasets are described in Table 1 and eTable 1. The initial local  
175 validation dataset did not overlap with the image data used in training. Images previously not  
176 seen by the network were presented to GD-CNN for assessment and automated diagnosis.  
177 The images were also independently assessed by three experienced professional graders (>2  
178 years' experience) in detecting referable GON.

179

### 180 **Online deep learning (ODL) system**

181 The central challenge of applying DLSs in medicine is the ability to guarantee  
182 generalizability in prediction. Generalization refers to the ability of DLSs to successfully  
183 perform when assessing previously unseen samples from different data sources. An ODL  
184 system was developed to improve the generalization ability of the GD-CNN model, making  
185 automatic GON diagnosis practical. In the ODL system, the GD-CNN model is used to  
186 sequentially predict GON with a Human-Computer Interaction (HCI) loop (eFigure 2 A). The  
187 HCI loop consisted of three iterative steps: (1) The computer used GD-CNN to initially  
188 diagnose glaucoma of fundus images with a high sensitivity rate; (2) the ophthalmologists  
189 manually confirmed the positive samples predicted by the computer; (3) the confirmed

190 samples fine-tuned the GD-CNN model, which was used for initial GON diagnosis of the  
191 subsequent fundus images (i.e., go to step 1).

192

### 193 **Visualization of Prediction**

194 Following Zeiler and Fergus,<sup>21</sup> we visualized the contributions of different regions to GD-  
195 CNN prediction of GON on fundus images. The visualization is represented by heat maps,  
196 which highlight strong prognostic regions of the fundus images. The experiment of occlusion  
197 testing was conducted to obtain the visualization results. First, original fundus image was  
198 resized into a 360x360 RGB image. Then, a 60x60 gray block was used to slice through the  
199 fundus image (with a stride of 10 pixels), alongside both horizontal and vertical axes.  
200 Consequently, the fundus image generates 961 (=31x31) visualization testing images, each of  
201 which has a 60x60 gray block at different position, respectively. Second, the visualization  
202 testing images were predicted using the GD-CNN model. For each visualization test image,  
203 the prediction probability output refers to the value of the visualization heat map at the  
204 corresponding position. Hence, the visualization heat map was 31x31. Finally, the heat map  
205 was mapped to the original fundus image to visualize the importance of each region in GON  
206 prediction.

207

208 The deep features refer to the output of the final max pooling layer, which is in 512  
209 dimensions. In order to visualize the distribution of the deep features from different  
210 categories, the dimensionality of deep features was reduced by t-distributed stochastic  
211 neighbor embedding visualization (t-SNE) from 512 to 3. Note that t-SNE is a state-of-the-art  
212 nonlinear dimensionality reduction method. The deep features from glaucoma and negative  
213 glaucoma are clustered into two groups once the training loss converges. The groups of two

214 clusters can be clearly separated, verifying the effectiveness of the deep features learned in  
215 GD-CNN.

216

## 217 **Statistical analysis**

218 The performance of our algorithm was evaluated in terms of area under the curve (AUC) of  
219 receiver operating characteristic (ROC) curves. 95% confidence intervals for AUC were  
220 calculated non-parametrically through logit-transformation-based confidence intervals, which  
221 was found to have good coverage accuracy over unbiased samples. In addition to AUC,  
222 sensitivity and specificity of each operating point in ROC curves were also measured with 2-  
223 sided 95% confidence intervals. These confidence intervals were calculated as Clopper-  
224 Pearson intervals, which are “exact” intervals based on cumulative probabilities.

225

226 Furthermore, to determine if the ODL system has an effect on diagnosing glaucoma,  
227 McNemar tests were conducted between the original GD-CNN model and the fine-tuned GD-  
228 CNN models. Specifically, two 2x2 contingency tables were applied to count the diagnosis  
229 changes after ODL, for positive and negative samples, respectively. Then a Chi-squared  
230 based P value was calculated along with the sensitivity/specificity over each validation  
231 dataset.

232

233 All statistical analyses were computed using the Stats Models (version 0.6.1) python package  
234 and Matlab AUC (version 1.1) package.

235

## 236 **Results**

### 237 **Training, validation and evaluation of the GD-CNN model**

238 From a total of 274 413 fundus images initially obtained from CGSA, 269 601 images passed  
239 initial image quality review and were graded for GON by the second-tier graders of Chinese  
240 board-certified ophthalmologists. The median quantity of images per ophthalmologist graded  
241 was 14 756 (range, 8 762-55 389) and ten ophthalmologists graded more than 15 000 images.  
242 13 254 images of disagreement in tier 2 grading were adjudicate by senior glaucoma  
243 specialists. 241 032 images (definite GON 29 865 (12.4%), probable GON 11 046 (4.6%),  
244 unlikely GON 200 121 (83%) from 68 013 patients were selected, using random sampling, to  
245 train the GD-CNN model. Validation and evaluation of the GD-CNN model was assessed  
246 using the remaining 28 569 images from CGSA. Distribution of the three diagnostic  
247 categories was 15.8% definite GON, 2% probable GON and 82.2% unlikely GON (eTable 1).  
248 In local validation dataset, the AUC of the GD-CNN model was 0.996 (95%CI, 0.995-0.998),  
249 and sensitivity and specificity in detecting referable GON was comparable with that of  
250 trained professional graders (96.2%vs 96.0%; P = 0.76; 97.7% vs 97.9%; P = 0.81  
251 respectively) (eFigure 3). To evaluate the ability of the GD-CNN to work across different  
252 populations, three clinical based studies were performed to reflect the routine functioning of  
253 an ophthalmic center. When images from these cohorts from different hospitals were  
254 diagnosed through GD-CNN and compared to clinical evaluation, performance remained  
255 high (Table 3), such that the AUC for referable GON ranged from 0.995 to 0.987 , with both  
256 sensitivity and specificity of greater than 90% (range: 93.6-96.1% and 95.6-97.1%  
257 respectively). Further evaluation was undertaken using the Handan Eye Study dataset to  
258 provide a real-world distribution of glaucoma patients. In this case AUC was 0.964 with a  
259 sensitivity of 91.0% and specificity 92.6% (Table 3). To test GD-CNN across a range of  
260 ethnic backgrounds, a multi-ethnic dataset (73.0% White, 19.3% Black/African American,  
261 5.4% Asian, 0.3% Middle Eastern) from the Hamilton Glaucoma Center was utilized, with  
262 AUC of 0.923, sensitivity of 87.7% and specificity 80.8%. GD-CNN showed an AUC of

263 0.823 with 82.2% sensitivity and 70.4% specificity in a varied range of image quality dataset  
264 from worldwide web (Table 3).

265

### 266 **Understanding the basis for incorrect diagnosis**

267 Among the local validation datasets, an additional analysis was conducted to further evaluate  
268 GD-CNN's performance, to better establish the basis for false positive and negative diagnosis  
269 (eTable 2). The most common reason for undetected GON from fundus images was  
270 pathological or high myopia for both GD-CNN and manual grading (n = 51 [46.3%] and n =  
271 50 [44.2%] respectively). Interestingly, the most likely cause for a false-positive  
272 classification by DLS or manual grading was also pathological or high myopia (n = 191  
273 [32.3%] and n = 183 [34%] respectively). Physiologically large cupping was also a common  
274 cause of false positives with manual diagnosis (n = 138 [25.6%]), and to a lesser degree with  
275 GD-CNN (n = 94 [16.0%]).

276

### 277 **Implementation of the ODL system**

278 The ODL system was implemented in the tele-ophthalmic image reading platform of Beijing  
279 Tongren Hospital (Appendix 1.4), which collected a group of fundus images every week  
280 (around 600 images). It was found that the ODL system both sensitivity and specificity  
281 improve with each group of samples collected sequentially over a five-week period (eFigure  
282 2).

283

### 284 **Visualization of prediction**

285 To visualize the learning procedure and represent the areas contributing most to the DLS, we  
286 created a heatmap which superimposed a convolutional visualization layer at the end of our  
287 network; performed on 1000 images (Figure 1 and eFigure 4). The regions of interest

288 identified to have made the greatest contribution to the neural network's diagnosis were also  
289 shared with 91.8% of ophthalmologists (Figure 2A). All areas containing optic nerve head  
290 variance and neuroretinal rim loss were located correctly on all the images used for testing,  
291 while RNFL defects and peripapillary atrophy (PPA) on occasions did not present a clear  
292 point of interest with an accuracy of 90.0% and 87.0% respectively. Figure 2B represents a t-  
293 distributed stochastic neighbor embedding visualization of this data set by our automated  
294 method, clearly showing 2 clusters of fundus images and indicating the ability of our model  
295 to separate normal from those with glaucoma.

296

## 297 **Discussion**

298 In this study, we focused on automating the diagnosis of glaucoma from fundus images by  
299 establishing a DLS (GD-CNN) with an ability to work across numerous populations.  
300 Previous studies have reported automated methods for the evaluation of glaucoma with most  
301 employing technology on feature extraction<sup>22-26</sup>. Recently, the DLS approach also has been  
302 adopted to provide high sensitivity and specificity for detecting GON from high-quality  
303 retinal fundus images.<sup>14,27,28</sup> The ambition of deep learning is to create a “fully-automated”  
304 screening model, which can automatically learn the features for glaucoma diagnoses without  
305 any human effort, avoiding misalignment or/and misclassification caused by introduced  
306 errors in the localization and segmentation. Compared with previous work, the GD-CNN  
307 model differs from conventional learning-based algorithms in a number of aspects.

308

309 The GD-CNN model was trained using a larger dataset than previous studies<sup>13,14,27-32</sup>. It is  
310 reasonable to assume that access to a greater pool of training images is likely to increase the  
311 accuracy of the DLS to detect glaucoma. A major issue with deep learning algorithms is their  
312 general applicability to systems and settings beyond the site of development. To address this

313 issue, additional data sets were employed. Datasets resulting from ophthalmic settings are  
314 likely to provide a higher incidence of glaucoma patients than is present in the general  
315 population. Therefore, to provide a realistic disease-screening test for GD-CNN, a population  
316 dataset obtained from the Handan Eye Study was employed, which provided a real-world  
317 ratio of individuals with and without diagnosed glaucoma<sup>33,34</sup>. Ethnicity can also present  
318 different anatomical/clinical features and incidence of glaucoma<sup>35</sup>. A number of the cohorts  
319 derived from Chinese centers have limited ethnic diversity. Therefore, to test GD-CNN  
320 across a range of ethnic backgrounds a multi-ethnic dataset, which includes White, African  
321 American, Asian, and Middle Eastern, from the Hamilton Glaucoma Center was utilized.  
322 Despite the different challenges imposed by these different data sets, GD-CNN consistently  
323 performed with high sensitivity and specificity. Another major factor that can impact on the  
324 generalization of DLSs is the image quality provided on which the DLS is making decisions  
325 and diagnosis. To address this important concern, GD-CNN was externally evaluated using a  
326 multi-quality image dataset of retinal fundus photographs established from website sources.  
327 Examination of 884 images available on the worldwide web using GD-CNN as expected  
328 proved a greater challenge, but analysis showed acceptable performance with AUC of 0.823  
329 with 82.2% sensitivity and 70.4% specificity.

330

331 The current study addressed the issue of false positive and negative diagnosis by the DLS and  
332 manual grading. The main reason for both false-negative and false-positive diagnosis by GD-  
333 CNN and manual grading was high or pathologic myopia, which are characterized by  
334 peripapillary atrophy (beta-zone), shallow cups, tilting and/or torsion optic disc. More studies  
335 assessing textural based properties are planned to allow more accurate classification by the  
336 algorithm which can distinguish among the optic disc region, central  $\beta$ -zone and peripheral  $\alpha$ -  
337 zone of peripapillary atrophy and other retinal areas.

338

339 To further evaluate the ability of the GD-CNN model across multiple populations, an ODL  
340 system was proposed in which the GD-CNN model iteratively updated with an HCI loop.  
341 Consequently, in the ODL system, the generalization ability of GD-CNN can be improved  
342 through human-computer interaction, such that each can educate and inform the other. An  
343 ODL system using a pre-trained GD-CNN model to reinforce training on limited local images  
344 would likely generate a more accurate model requiring less time for local dataset  
345 classifications. In principle, the ODL system we have described here could potentially be  
346 employed on a wide range of medical images across multiple disciplines. Further benefit  
347 may come from the use of AI with digital images like a combination of structural and  
348 functional testing, and even multiple other orthogonal datasets, for example, cardiovascular  
349 data and genomic data, to further enhance the value of data utilization for the health care  
350 system.

351

## 352 **Conclusions**

353 The GD-CNN model, which was driven by a large-scale database of fundus images, has high  
354 sensitivity and specificity for detecting glaucoma. The experimental results show the  
355 potential of automated DLSs in enhancing current screening programs in a cost-effective and  
356 time efficient manner. The generalization of this approach might be facilitated by training the  
357 GD-CNN model on large-scale data and implementing GD-CNN in an ODL system, which  
358 may be further refined through a human computer interface.

359

- 360 1. Tham YC, Li X, Wong TY, et al. Global prevalence of glaucoma and projections of  
361 glaucoma burden through 2040: a systematic review and meta-analysis.  
362 *Ophthalmology*. 2014;121(11):2081-2090. doi: 10.1016/j.opthta.2014.05.013



- 363 2. Quigley HA, Broman AT. The number of people with glaucoma worldwide in 2010  
364 and 2020. *Br J Ophthalmol*. 2006;90(3):262-267. doi:10.1136/bjo.2005.081224
- 365 3. Hood DC, Raza AS, de Moraes CG, et al. Glaucomatous damage of the macula. *Prog*  
366 *Retin Eye Res*. 2013; 32:1-21. doi: 10.1016/j.preteyeres.2012.08.003
- 367 4. Tatham AJ, Weinreb RN, Medeiros FA. Strategies for improving early detection of  
368 glaucoma: the combined structure-function index. *Clin Ophthalmol*. 2014; 8:611-621.  
369 doi: 10.2147/OPHTH.S44586
- 370 5. Weinreb RN, Aung T, Medeiros FA. The pathophysiology and treatment of glaucoma:  
371 a review. *JAMA*. 2014; 311(18):1901-1911. doi: 10.1001/jama.2014.3192
- 372 6. Pizzi LT, Waisbourd M, Hark L, et al. Costs of a community-based glaucoma  
373 detection programme: analysis of the Philadelphia Glaucoma Detection and  
374 Treatment Project. *Br J Ophthalmol*. 2018; 102(2):225-232. doi:  
375 10.1136/bjophthalmol-2016-310078
- 376 7. Zhao D, Guallar E, Gajwani P, et al. Optimizing glaucoma screening in high-risk  
377 population: design and 1-year findings of the screening to prevent (SToP) glaucoma  
378 study. *Am J Ophthalmol*. 2017; 180:18-28. doi: 10.1016/j.ajo.2017.05.017
- 379 8. Fleming C, Whitlock EP, Beil T, et al. Screening for primary open-angle glaucoma in  
380 the primary care setting: an update for the US preventive services task force. *Ann Fam*  
381 *Med*. 2005; 3(2):167-170. doi: 10.1370/afm.293
- 382 9. Esteva A, Kuprel B, Novoa RA, et al. Dermatologist-level classification of skin  
383 cancer with deep neural networks. *Nature*. 2017; 542(7639):115-118. doi:  
384 10.1038/nature21056
- 385 10. Ehteshami Bejnordi B, Veta M, Johannes van Diest P, et al. Diagnostic Assessment  
386 of Deep Learning Algorithms for Detection of Lymph Node Metastases in Women  
387 with Breast Cancer. *JAMA*. 2017; 318(22):2199-2210. doi: 10.1001/jama.2017.14585

- 388 11. Korfiatis P, Kline TL, Coufalova L, et al. MRI texture features as biomarkers to  
389 predict MGMT methylation status in glioblastomas. *Med Phys*. 2016; 43(6):2835-  
390 2844. doi: 10.1118/1.4948668
- 391 12. Gulshan V, Peng L, Coram M, et al. Development and Validation of a Deep Learning  
392 Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs.  
393 *JAMA*. 2016; 316(22):2402-2410. doi: 10.1001/jama.2016.17216
- 394 13. Ting DSW, Cheung CY, Lim G, et al. Development and validation of a deep learning  
395 system for diabetic retinopathy and related eye diseases using retinal images from  
396 multiethnic populations with diabetes. *JAMA*. 2017; 318(22):2211-2223. doi:  
397 10.1001/jama.2017.18152
- 398 14. Li Z, He Y, Keel S, et al. Efficacy of a deep learning system for detecting  
399 glaucomatous optic neuropathy based on color fundus photographs. *Ophthalmology*.  
400 2018; 125(8):1199-1206. doi: 10.1016/j.ophttha.2018.01.023
- 401 15. Wong TY, Bressler NM. Artificial Intelligence with Deep Learning Technology  
402 Looks into Diabetic Retinopathy Screening. *JAMA*. 2016; 316(22):2366-2367. doi:  
403 10.1001/jama.2016.17563
- 404 16. Castelvechi D. Can we open the black box of AI? *Nature*. 2016; 538(7623):20-23.  
405 doi: 10.1038/538020a
- 406 17. Verghese A, Shah NH, Harrington RA. What this computer needs is a physician:  
407 humanism and artificial intelligence. *JAMA*. 2018; 319(1):19-20. doi:  
408 10.1001/jama.2017.19198
- 409 18. He K, Zhang X, Ren S. et al. Deep residual learning for image recognition.  
410 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.  
411 2016:770-778. doi: 10.1109/CVPR.2016.90
- 412 19. Haleem MS, Han L, van Hemert J. et al. Automatic extraction of retinal features from

- 413 colour retinal images for glaucoma diagnosis: A review. *Computerized Medical*  
414 *Imaging and Graphics*. 2013; 37(7-8):581-596. doi:  
415 10.1016/j.compmedimag.2013.09.005
- 416 20. Szegedy, C, Liu W, Jia Y, et al. Going deeper with convolutions. Proceedings of the  
417 *IEEE Conference on Computer Vision and Pattern Recognition*. 2015:1-9.  
418 doi:10.1109/CVPR.2015.7298594
- 419 21. Zeiler MD, Fergus R. Visualizing and understanding convolutional networks.  
420 Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.  
421 2014:818-833. doi: 10.1007/978-3-319-10590-1\_53
- 422 22. Singh A, Dutta MK, Partha SM, et al. Image processing based automatic diagnosis of  
423 glaucoma using wavelet features of segmented optic disc from fundus image. *Comput*  
424 *Methods Programs Biomed*. 2016; 124:108-120. doi: 10.1016/j.cmpb.2015.10.010
- 425 23. Issac A, Partha SM, Dutta MK. An adaptive threshold-based image processing  
426 technique for improved glaucoma detection and classification. *Comput Methods*  
427 *Programs Biomed*. 2015; 122(2):229-244. doi: 10.1016/j.cmpb.2015.08.002
- 428 24. Chakrabarty L, Joshi GD, Chakravarty A, et al. Automated Detection of Glaucoma  
429 from Topographic Features of the Optic Nerve Head in Color Fundus Photographs. *J*  
430 *Glaucoma*. 2016; 25(7):590-597. doi: 10.1097/IJG.0000000000000354
- 431 25. Chen X, Xu Y, Wong DWK, et al. Glaucoma detection based on deep convolutional  
432 neural network. *Conf Proc IEEE Eng Med Biol Soc*. 2015:715-718. doi:  
433 10.1109/EMBC.2015.7318462
- 434 26. Annan L, Jun C, Damon WKW, et al. Integrating holistic and local deep features for  
435 glaucoma classification. *Conf Proc IEEE Eng Med Biol Soc*. 2016:1328-1331. doi:  
436 10.1109/EMBC.2016.7590952
- 437 27. Christopher M, Belghith A, Bowd C, et al. Performance of Deep Learning

- 438 Architectures and Transfer Learning for Detecting Glaucomatous Optic Neuropathy  
439 in Fundus Photographs. *Sci Rep.* 2018;8(1):16685. doi: 10.1038/s41598-018-35044-9
- 440 28. Shibata N, Tanito M, Mitsuhashi K, et al. Development of a deep residual learning  
441 algorithm to screen for glaucoma from fundus photography. *Sci Rep.* 2018;8(1):14665.  
442 doi: 10.1038/s41598-018-33013-w
- 443 29. Meier J, Bock R, Michelson G, et al. Effects of Preprocessing Eye Fundus Images on  
444 Appearance Based Glaucoma Classification. Proceedings of the International  
445 Conference on Computer Analysis of Images and Patterns. 2007:165-172. doi:  
446 10.1007/978-3-540-74272-2\_21
- 447 30. Bock R, Meier J, Michelson G, et al. Classifying Glaucoma with Image-Based  
448 Features from Fundus Photographs. Proceedings of the 29th DAGM Symposium,  
449 Heidelberg, Germany, September 12-14, 2007. doi:10.1007/978-3-540-74936-3\_36
- 450 31. Bock R, Meier J, Nyúl L, et al. Glaucoma risk index: automated glaucoma detection  
451 from color fundus images. *Med Image Anal.* 2010; 14(3): 471-481. doi:  
452 10.1016/j.media.2009.12.006
- 453 32. Keerthi SS, Shevade SK, Bhattacharyya C, et al. Improvements to Platt's SMO  
454 Algorithm for SVM Classifier Design. *Neural Computation.* 2014; 13(3):637-649.  
455 doi:10.1162/089976601300014493
- 456 33. Wang NL, Hao J, Zhen Y, et al. A Population-based Investigation of Circadian  
457 Rhythm of Intraocular Pressure in Habitual Position Among Healthy Subjects: The  
458 Handan Eye Study. *J. Glaucoma.* 2016; 25(7):584-589. doi:  
459 10.1097/IJG.0000000000000351
- 460 34. Zhang Y, Li SZ, Li L, et al. The Handan Eye Study: comparison of screening methods  
461 for primary angle closure suspects in a rural Chinese population. *Ophthalmic*  
462 *Epidemiol.* 2014; 21(4):268-275. doi: 10.3109/09286586.2014.929707

463 35. Cho HK, Kee C. Population-based glaucoma prevalence studies in Asians. *Surv*  
464 *Ophthalmol.* 2014; 59(4):434-447. doi: 10.1016/j.survophthal.2013.09.003

465

466 **Legends:**

467 **Figure 1. Visualization of deep features of the GD-CNN deep learning system.**

468 Visualization maps generated from deep features, which can be superimposed on the input  
469 image to highlight the areas the model considered important in making its diagnosis.

470 **Figure 2. Training loss and visualization of deep features at different training iterations.**

471 (A) Training loss with accuracy with training iterations. (B) Feature clustering with the  
472 progress of training. The dimensionality of deep features was nonlinearly reduced by t-  
473 distributed stochastic neighbor embedding (t-SNE) method for visualization.

474

475

476 **Acknowledgements**

477 **Funding/Support:** The research has received funding from National Natural Science Fund  
478 Projects of China (81271005), Beijing municipal administration of hospitals Qingmiao  
479 projects (QMS20180210), The priming scientific research foundation for the junior  
480 researcher in Beijing Tongren Hospital (2016-YJJ-ZZL-021), Beijing Tongren Hospital Top  
481 Talent Training Program. Medical Synergy Science and Technology Innovation Research  
482 (Z181100001918035)

483 **Role of the Funder/Sponsor:** The funding organizations had no role in design and conduct  
484 of the study; collection, management, analysis, and interpretation of the data; preparation,  
485 review, or approval of the manuscript; and decision to submit the manuscript for publication.

486 **Conflict of Interest Disclosures:** All authors have completed and submitted the authorship  
487 forms.

488 **Author Contributions:** Drs H. Liu, M. Xu and N. Wang had full access to all of the data in  
489 the study and take responsibility for the integrity of the data and the accuracy of the data  
490 analysis.

491 Drs M. Xu and N. Wang contributed equally.

492 Concept and design: H. Liu, M. Xu, N. Wang

493 Acquisition, analysis, or interpretation of data: H. Liu, L. Li, C. Qiao, C. Zhang, P. Liu, S. Li,  
494 H. Wang, D. Mou, R. Pang, D. Yang, L. Jiang, Y. Chen, M. Hu, Y. Xu, H. Kang, X. Ji, C.  
495 Tham, C. Cheung

496 Critical revision of the manuscript for important intellectual content: H. Liu, I. M.  
497 Wormstone, D. Yang, R. Chang, D. S. Ting, Z. Wang, T. Y. Wong, M. Xu, R. N. Weinreb

498 Statistical analysis: H. Liu, D. Yang, L. Li, L. Jiang

499 Obtained funding: H. Liu, M. Xu, N. Wang

500 Administrative, technical, or material support: H. Liu, C. Qiao, C. Zhang, P. Liu, S. Li, H.  
501 Wang, Z. Wang, M. Xu, N. Wang

502 Supervision: H. Liu, Z Wang, M. Xu, N. Wang  
503

Table 1. Summary of Source Datasets

Source Datasets	Images No.	Eye s <sup>a</sup> No.	Individuals No.	Age, Mean <sup>b</sup> (SD)	Female <sup>b</sup> /Total (%)	Cohort	Ethnicity/Race	Camera	Assessor
CGSA	274 413	138 210	69 105	54.1 (14.5)	20 167 (55.8%)	Clinic-based	Han Chinese (78.3%)	Topcon, Canon, Carl Zeiss	Professional grader team
Beijing Tongren Hospital	20 466	10 308	5 154	52.8 (16.7)	1 068 (49.7%)	Clinic-based	Han Chinese (81.7%)	Topcon, Canon	2 Ophthalmologists; arbitration by 1 glaucoma specialist
Peking University Third Hospital	12 718	64 60	3 230	57.2 (10.9)	327 (43.1%)	Clinic-based	Han Chinese (79.5%)	Topcon	2 Ophthalmologists; arbitration by 1 glaucoma specialist
Harbin Medical University First Hospital	9 305	4 732	2 366	59.9 (11.2)	771 (57.3%)	Clinic-based	Han Chinese (82.9%)	Topcon	2 Professional senior graders; arbitration by 1 glaucoma specialist
Handan Eye Study	29 676	13 404	6 702	55.2 (10.9)	2 589 (42.2%)	Population-based	Han Chinese (80.1%)	Topcon, Canon	3 Glaucoma specialists
Hamilton Glaucoma Center	7 877	3 938	1 969	58.2 (19.2)	1041 (52.9%)	Clinic-based	White (73.0%), Black/African American (19.3%), Asian (5.4%), Middle Eastern (0.3%)	Topcon, Canon	3 Glaucoma specialists
Website	884	884	884	N/A	N/A	Website-based	N/A	N/A	2 Professional senior graders; arbitration by 1 glaucoma specialist

<sup>a</sup> For each patient, 2 fundus images were taken and recorded of each eye. <sup>b</sup> Individual data including age sex and ethnicity/race were available for CGSA (52.3%), Beijing Tongren Hospital (41.7%), Peking University Third Hospital (23.5%), Harbin Medical University First Hospital (56.9%), Handan Eye Study (99.6%), Hamilton Glaucoma Center (100%), Website (N/A).

507 Table 2. The Classification for Glaucomatous Optic Neuropathy

Classification	Clinical Features
Unlikely glaucomatous optic neuropathy	With no sign of the following
Probable glaucomatous optic neuropathy	At least two conditions positive: $0.7 \leq \text{VCDR} < 0.85$ ; Rim Width $\leq 0.1 \text{ DD}$ ; General Rim Thinning $\geq 60^\circ$ or localized Rim Thinning $< 60^\circ$ (11-1 o'clock or 5-7 o'clock); RNFL defects; Splinter Hemorrhages, Peripapillary Atrophy (Beta zone)
Definite glaucomatous optic neuropathy	Any of the following conditions: $\text{VCDR} \geq 0.85$ ; RNFL defects corresponds with thinning area of rim or notches.
VCDR: vertical cup-to-disc ratio. DD: disc diameter RNFL: retinal nerve fiber layer	

508  
509



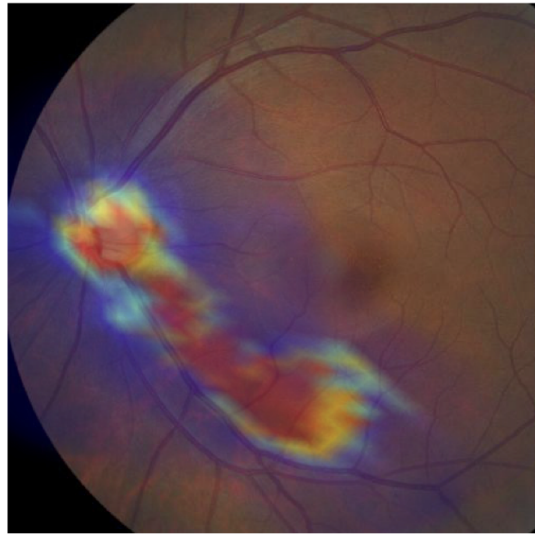
Table 3. The Performance of the GD-CNN in Validation Datasets

Datasets (No. of Images)	AUC (95% CI)	% (95% CI)		Confusion Result No. (%)				Total Concordant Images
		Sensitivity	Specificity	True- Positive	False- Positive	False- Negative	True- Negative	
Local Validation								
Chinese Glaucoma Study Alliance (N =28 569)	0.996 (0.995- 0.998)	96.2 (95.4 – 96.9)	97.7 (97.5- 97.9)	2 786 (9.8)	588 (2.1)	110 (0.4)	25 085 (87.8)	27 871 (97.6)
Clinical Validation								
Beijing Tongren Hospital (N =20 466)	0.995 (0.996- 0.996)	96.1 (95.2- 96.9)	97.1 (96.8- 97.3)	2 226 (10.9)	534 (2.6)	90 (0.4)	17 616 (86.1)	19 842 (97.0)
Peking University Third Hospital (N = 12 718)	0.994 (0.991- 0.996)	96.0 (93.9- 97.2)	96.1 (95.8- 96.5)	593 (4.7)	468 (3.7)	26 (0.2)	11 631 (91.5)	12 224 (96.1)
Harbin Medical University First Hospital (N = 9 305)	0.987 (0.982- 0.991)	93.6 (90.9- 95.6)	95.6 (95.1- 96.0)	435(4.7)	392 (4.2)	30 (0.3)	8 448 (90.8)	8 883 (95.5)
Population Screening Validation								
Handan Eye Study (N = 29 676)	0.964 (0.952- 0.972)	91.0 (88.4- 93.1)	92.6 (92.2- 92.8)	543 (1.8)	2 175 (7.3)	54 (0.2)	26 904 (90.7)	27 447 (92.5)
Multi-ethnic Validation								
Hamilton Glaucoma Center (N=7 877)	0.923 (0.916- 0.930)	87.7 (86.8- 88.5)	80.8 (78.9- 82.5)	5224 (66.3)	369 (4.7)	733 (9.3)	1551 (19.7)	6 775 (86.0)
Multi-quality Validation								
Website (N = 884)	0.823 (0.787- 0.855)	82.2 (76.9- 86.6)	70.4 (65.8- 74.7)	212 (31.0)	126 (18.4)	46 (6.7)	300 (43.9)	512 (74.9)



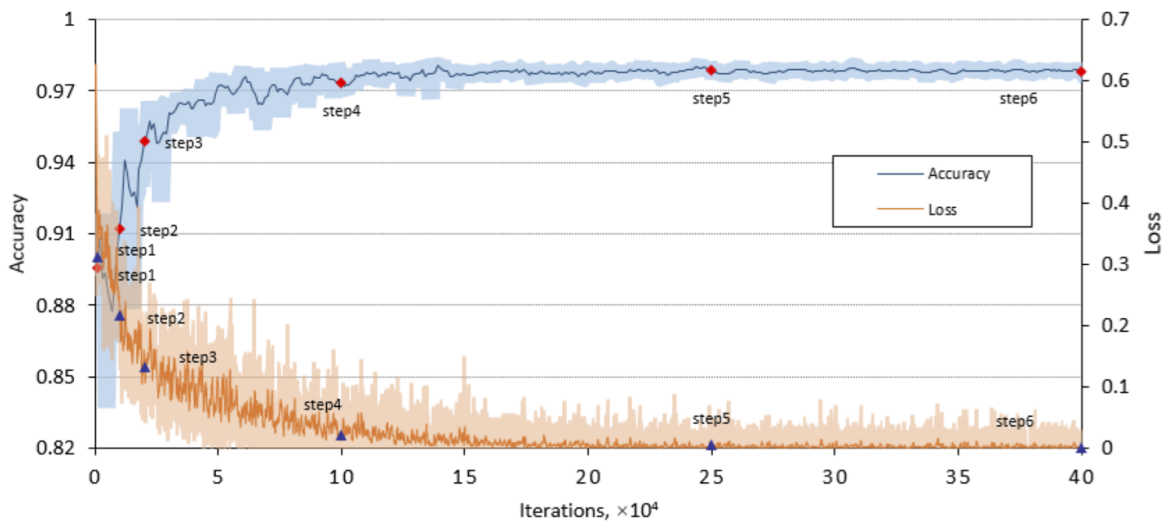


Original



Heatmap

A



B

