1 **APPENDIX**

2

3 **Image Masking**

4     Image masking was attempted to remove or minimize eyelash and eyelid artifacts present

5 in many of the ultra-widefield (UWF) images. The U-Net convolutional network for biomedical

6 image segmentation[1] was applied to the UWF color and UWF fundus autofluorescence (FAF)

7 images. The segmentation network uses a U-Net architecture to translate the raw UWF color and

8 UWF FAF images into a mask map with two classes including the image for analysis and a

9 portion discarded as noise. The U-Net model was trained with 20 UWF color images from the

10 cognitively healthy control cohort that had undergone manual image masking to remove eyelash

11 and eyelid artifacts. Manually masked images were binary images, with noise pixels (eyelid and

12 eyelash artifacts) set to 0, and others pixels set to 1. The area under the curve (AUC) of the

13 automated segmentation task was 0.970 when compared to a sample of 20 manual segmentations

14 based on the leave-one-out cross-validation (LOOCV).[2] The results of our convolutional neural

15 network (CNN) for prediction of a symptomatic Alzheimer's disease (AD) diagnosis were

16 compared with and without the image masking. There was no improvement in performance in

17 the validation or test sets when image masking was used. Thus, image masking was not used in

18 the results reported in this study.

19

20 **Image Cropping**

21     Results of the combination model incorporating all retinal imaging modalities, OCT and

22 OCTA quantitative data, and patient data were compared with and without cropping of UWF

23 color images. With the entire UWF image incorporated, the model using all images and all data

24 inputs achieved an AUC of 0.854 [95% CI: 0.717, 0.990] on the validation set and 0.826 [95%

25 CI: 0.715, 0.937] on the test set, while use of cropped images allowed the model to achieve an

26 AUC of 0.854 [95% CI: 0.718, 0.991] on the validation set and 0.836 [95% CI: 0.729, 0.943] on

27 the test set. Thus, image cropping slightly improved the predictive value of the model. The

28 model used cropped UWF color and UWF FAF images.

29

30 **Image Resizing**

31     All images were resized prior to use in the model. Images were resized to 128x128 pixels

32 from 1800x1800 for UWF color and FAF, 409x409 – 450x450 for OCTA, and 281x281 for

33  ganglion cell-inner plexiform layer (GC-IPL) maps. This resizing resulted in 7x7x64

34  convolutional features that, when flattened out, are represented as 3136-dimensional vectors.

35  While such resizing prevents fine details from being preserved, sample size limitations did not

36  allow us to preserve all of the fine details of each imaging modality. While model architectures

37  with larger images were trialed, these architectures were unable to obtain significant

38  performance improvements during the development phase, quantified on the validation set. Thus,

39  while finer details have the potential of greatly improving model performance, we would need a

40  significantly larger sample size to seize the benefits of such resolution gain.

41

42  **Attention Maps**

43      To interpret the classification result made by the model, attention maps were generated to

44  visualize the discriminative image regions used by the model to distinguish symptomatic AD

45  subjects from controls using class activation mappings (CAMs).[3] After the model was fully

46  trained, we fed the images, including OCTA images and UWF color and UWF FAF images, into

47  the image feature extractor and created the feature maps for each imaging modality, $f_{OCTA}$, $f_{Color}$

48  $_{and\ FAF}$, and $f_{GC-IPL}$ (Figure 1). We then projected back the weights of the corresponding fully

49  connected layer onto the feature maps to produce the CAMs. The final attention maps were

50  generated by scaling the CAMs to a heatmap with the same size as the input images. Figures 2

51  and 3 include examples of the OCTA attention maps that were generated. Supplemental Figure 2

52  contains attention map examples of the UWF images.

53
54
55  **Model Structure**

56      The image feature extractor used in this study has a similar structure to the first five layers of

57  the ResNet18[4] neural network. The model included a total of 166,936 parameters. More

58  specifically, the CNN-based feature extractor consists of 157,504 parameters (1 layer of 64

59  7x7x3 filters, 4 layers of 64 3x3x64 filters, and 5 layers of batch normalizations of 64

60  dimensions). The fully connected layers that output the modality-wise pre-classification

61  probabilities consist of 9,432 parameters (3,137 weights each for UWF color and FAF, OCTA,

62  and GC-IPL maps, and 21 weights for the quantitative features).

63      Two primary differences between our model and the original ResNet18[4] are: 1) Stride size of

64  the first convolutional layer was changed from 2 to 1. As input images are compressed from

65    large images, we found that a stride size of 2 in the first layer may lose some detailed

66    information in the compressed images; 2) Pooling size was changed from 2x2 to 4x4 to reduce

67    the dimension of output feature maps, which leads to fewer parameters in the following fully-

68    connected (FC) layers and mitigates over-fitting. As our feature extractor has the same structure

69    with ResNet18, it can still be initialized by ResNet18 pre-trained on ImageNet.[5] The detailed

70    structure of the image feature extractor can be found in Supplemental Figure 1.

71        A significant concern with training a CNN with a dataset of this size is over-fitting. To

72    mitigate over-fitting, we employed the following technologies.

73        1）Model Architecture

74            Given the limitation of our relatively smaller training dataset, we used only the first five

75    layers of the ResNet18[4] as the feature extractor and discarded other layers. The number of layers

76    is typically chosen based on performance of the training and validation sets after assessing

77    different layer configurations. At the same time, we shared the feature extractor across different

78    modalities, which increases the amount of data available to train the feature extraction

79    component of the model (i.e. the CNN model). Provided that different imaging modalities have

80    their own idiosyncrasies, feeding them through the same feature extractor seems counterintuitive.

81    However, noting that modality differences are predominantly high-level image characteristics,

82    we customized the feature extractor for each modality by appending a modality-specific FC layer

83    at the output of the CNN-based feature extractor, as shown in Figure 1. As a result, these FC

84    layers learned the unique characteristics of each modality. Given that the feature extractor

85    contained only five layers, it is likely to be extracting general image features that are then

86    adapted for each modality via FC layers. After passing through the FC layers, pre-classification

87    signals were generated separately for each imaging modality describing how likely an eye

88    carried an AD diagnosis. We also increased the average pooling from 2x2 to 4x4 (the "avg pool,

89    ¼" by 4 in Supplemental Figure 1) between different convolutional layers to help reduce the

90    dimensions of the image feature extractor's outputs, which led to fewer parameters for the

91    corresponding FC layers. Finally, given the dataset size, we were aware that overloading the

92    model with excess input data from each subject could also cause over-fitting. Thus, providing the

93    model with the entire set of volumetric OCT images would not likely improve performance. In

94    light of this, quantitative OCT data that summarized the volumetric OCT findings with fewer

95    data points was used in our study. In future studies, using volumetric OCT images may be further

96    explored; however, given that our model achieved similar AUC values with inputs of all images

97    only and inputs including quantitative OCT data, it remains unclear if volumetric data would

98    improve performance.

99        2) L2 and L1 Regularization

100       L2 regularization is a technique where the sum of squared parameters of a model

101    (multiplied by a coefficient) is added into the loss function as a penalty term to be minimized. It

102    tends to cause the learning algorithm to perceive the input as having higher variance, which

103    makes it shrink the weights on features with low covariance with the target label.[6,7] L1

104    regularization adds the sum of the absolute values of the individual parameters to the loss

105    function. In comparison to L2 regularization, L1 regularization results in a solution that is

106    sparser, which is consistent with the fact that some parameters may have an optimal value of

107    zero. Thus, L1 regularization has also been used for feature selection.[6] Both L1 and L2

108    regularization make it difficult for the regularized network to learn local noise in the dataset and

109    force the network to learn only those features which are often seen across the training set.

110    According to the properties of L2 and L1 regularization, we added L2 regularization to the entire

111    model, and L1 regularization on the FC layers to the final loss function to limit the capacity of

112    the model and prevent over-fitting.[8,9] We tested different weights from 0.001 to 10 for the

113    regularization loss and found 0.01 worked best. We also tried dropout with different dropout

114    rates without benefit.

115        3) Data Augmentation

116       Data augmentation can be used to reduce over-fitting and increase the amount of training

117    data. It creates new images by transforming (rotating, translating, scaling, flipping, distorting)

118    and adding some noise (e.g. Gaussian noise) to the images in the training dataset. Both the

119    original image and the created images are fed into the neural network. In our model, we

120    increased the diversity of images for training models through rotating, shifting, cropping, and

121    zooming images.

122        4) Transfer Learning and Fine-tuning

123       We chose to initialize parameters of our feature extractor with weights of the pre-trained

124    ResNet[4] models to aid in the classification task. This model was then fine-tuned with our labeled

125    imaging data. This method transferred the knowledge from general large-scale image datasets

126    (ImageNet)[5] to our task. Gulshan et al.[10] also used this method to speed up the training of a deep

127 learning model, which was designed to detect diabetic retinopathy in retinal fundus photographs,

128 although they had a relatively larger dataset.

129

130 **Inter-eye Correlations**

131 We acknowledge that, when evaluating each eye separately from each patient in the test

132 set as we have done herein, inter-eye correlation can result in inflated performance metrics. As a

133 result, we recalculated the performance metrics by repeatedly sampling a single eye for each of

134 the 34 subjects in the test set. We reviewed AUC summaries for 300, 500, 1000, and 10,000

135 samples, each of which demonstrated stable means and standard deviations. For 10,000 samples,

136 calculated considering data from both eyes of each patient, the "patient-wise" AUC for the best

137 performing model was 0.84 (SD 0.034). This result is very similar to the results in Table 1,

138 which consider only a single eye for each patient and demonstrate the "eye-wise" AUC. Thus, in

139 our case, inter-eye correlation does not seem to affect performance estimates.

140

141 **Performance Reporting**

142 We presented AUC as the main performance metric because it is widely used and

143 accepted. However, we recognize that AUC can sometimes be misleading in situations where

144 there is an imbalance between cases and controls. In Appendix Figure 1, we include the

145 performance recall curve (AUPRC) for the test set for the best performing model, including GC-

146 IPL maps, quantitative data, and patient data. Precision in classification tasks is also referred to

147 as positive predictive value (PPV). Recall in classification tasks is also referred to as true

148 positive rate (TPR) or sensitivity. We chose the optimal point according to F1-score, which

149 combines precision and recall into one metric by calculating the harmonic mean between the

150 two. The optimal point is achieved when the threshold for probability of AD, P(AD), is set to

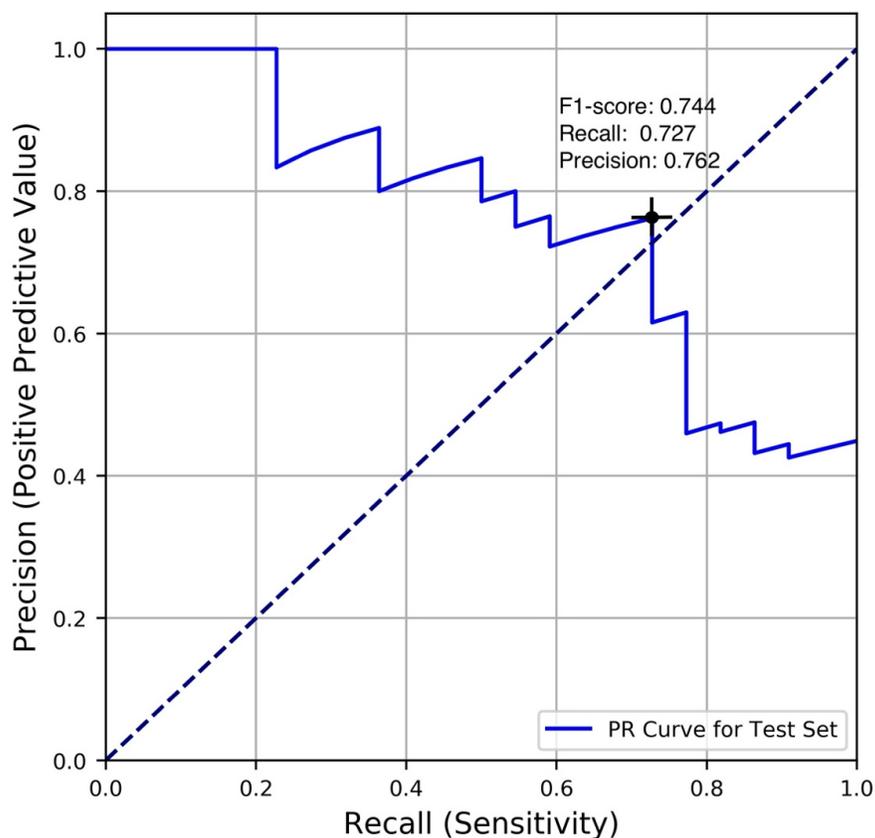151 0.224. The model achieves F1-score of 0.744, recall 0.727, and precision 0.762.

152
153
154

155

156

157

158

Supplemental material

BMJ Publishing Group Limited (BMJ) disclaims all liability and responsibility arising from any reliance placed on this supplemental material which has been supplied by the author(s)

*Br J Ophthalmol*

159

160     **Appendix Figure 1**



161
162

163     **Age Differences in AD and Control Subjects**

164         Supplemental Table 1 demonstrates that AD subjects in our study were significantly older

165     (2.73 years on average) than control subjects (p=0.015). To address any concern that the model

166     might simply be detecting the age difference between the two groups, we performed experiments

167     using age directly to predict AD using a decision tree model (to account for nonlinear effects)

168     using the same data partitions we used for the model (training, validation, and test sets). Results

169     for tree models of varying depths (1-10) shown in Appendix Table 1 below indicate that the

170     performance of age only (depth=5), 0.583 validation AUC, and 0.528 test AUC, is considerably

171     lower than that of the image-based model.

172

173     **Appendix Table 1**

| Max Depth | Validation Set Accuracy | Validation Set AUC | Test Set Accuracy | Test Set AUC |
|:---:|:---:|:---:|:---:|:---:|
| 1 | 0.733 | 0.458 | 0.618 | 0.480 |
| 2 | 0.800 | 0.500 | 0.647 | 0.502 |
| 3 | 0.733 | 0.458 | 0.618 | 0.480 |
| 4 | 0.733 | 0.458 | 0.647 | 0.549 |
| 5 | 0.733 | 0.583 | 0.618 | 0.528 |
| 6 | 0.667 | 0.542 | 0.588 | 0.482 |
| 7 | 0.667 | 0.542 | 0.588 | 0.482 |
| 8 | 0.667 | 0.542 | 0.588 | 0.482 |
| 9 | 0.667 | 0.542 | 0.559 | 0.460 |
| 10 | 0.667 | 0.542 | 0.559 | 0.460 |

## REFERENCES

1. Ronneberger O, Fischer P, Brox T. *U-Net: Convolutional Networks for Biomedical Image Segmentation.* Vol 9351: Springer, Cham; 2015.
2. Molinaro A, Simon R, Pfeiffer R. Prediction error estimation: a comparison of resampling methods. *Bioinformatics.* 2005;21(15):3301-3307.
3. Zhou B, Khosla A, Lapedriza A, Oliva A, Torralba A. Learning deep features for discriminative localization. *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2016:2921-2929.
4. He K, Zhang, X., Ren, S., Sun, J. Deep residual learning for image recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2016:770-778.
5. Deng J, Dong W, Socher R, Li-Jia L, Li K, Fei-Fei L. ImageNet: A large-scale hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition.* 2009.
6. Goodfellow I, Bengio Y, Courville A. *Deep Learning.* Cambridge, MA: The MIT Press; 2016.
7. Zhang C, Bengio S, Hardt M, Recht B, Vinyals O. Understanding deep learning requires rethinking generalization. *International Conference on Learning Representations (ICLR).* 2017;arXiv:1611.03530.
8. Bühlmann P, van de Geer S. *Statistics for High-Dimensional Data.* Vol 1. Heidelberg: Springer-Verlag Berlin Heidelberg; 2011.

Supplemental material

BMJ Publishing Group Limited (BMJ) disclaims all liability and responsibility arising from any reliance placed on this supplemental material which has been supplied by the author(s)

*Br J Ophthalmol*

198   9.   Loshchilov I, Hutter F. Decoupled weight decay regularization. *International Conference*
199        *on Learning Representations (ICLR).* 2019;arXiv:1711.05101.
200   10.  Gulshan V, Peng L, Coram M, et al. Development and validation of a deep learning
201        algorithm for detection of diabetic retinopathy in retinal fundus photographs. *JAMA*.
202        2016;316(22):2402-2410.
203