

RetCAD version 2.2

Computer aided detection for Age-related Macular Degeneration, Diabetic Retinopathy and Glaucoma

About this white paper

This white paper applies to RetCAD (version 2.2). It describes the general principles of the RetCAD software and presents validations of the software compared to human experts, on various datasets.

Contents

Introduction	3
ROC analysis	4
RetCAD: How does it work?	6
RetCAD: Performance evaluation	7
Messidor	7
Operating points for RetCAD DR	8
Comparison with other systems and human experts	8
Messidor-2	9
Operating points for RetCAD DR	10
Comparison with other systems and human experts	10
Private1 dataset	11
Operating points for RetCAD AMD	12
Comparison with other systems and human experts	12
Mixed AMD-DR dataset	13
Operating points for RetCAD AMD/DR	14
Comparison with other systems and human experts	14
ORIGA dataset	15
Operating points for RetCAD GLC	16
Comparison with other systems and human experts	17
REFUGE dataset	17
Operating points for RetCAD GLC	18
Comparison with other systems and human experts	19

Introduction

RetCAD was developed by Thirona Retina. RetCAD is a class IIa CE-certified medical device software product that uses deep learning to analyze color fundus images for the presence of Age-related Macular Degeneration (AMD), Diabetic Retinopathy (DR) and Glaucoma (GLC).

RetCAD takes a color fundus (CF) image as input and produces several outputs. These outputs include a quality assessment of the input image, heatmaps indicating possibly abnormal areas, and a score for each of these retinal diseases. The scores for AMD and DR indicate the severity of the disease, whereas the vertical cup-to-disc (VCDR) and the GLC score give an indication of the presence of Glaucoma.

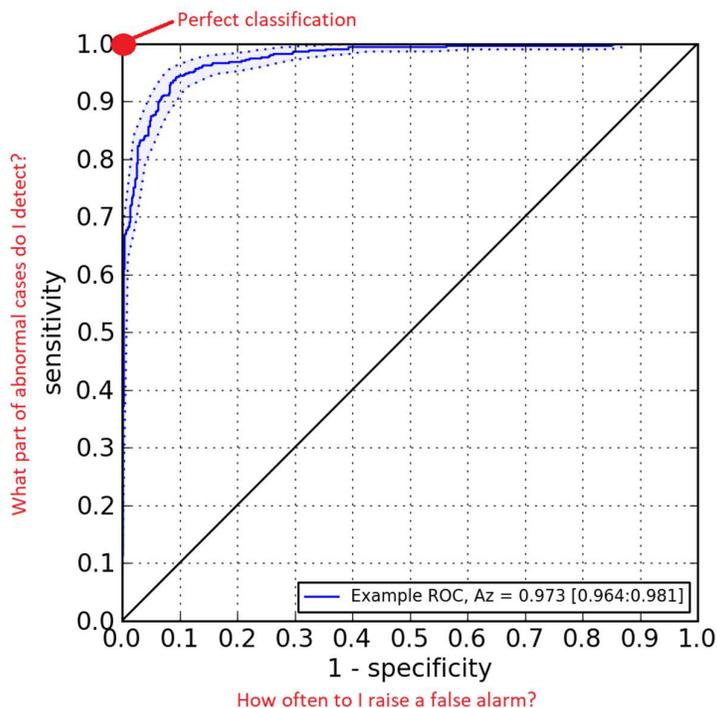
Users can take the output into account in their clinical work: they can decide if a new image should be acquired, in case the quality assessment indicates suboptimal image quality; they can decide to refer a patient for further testing for the presence of AMD, DR, GLC or other retinal abnormalities in case the heatmaps display suspicious regions that are verified by a human operator or when the scores are above certain thresholds.

ROC analysis

Definitions:

- **Sensitivity:** proportion of positive images (i.e. having an abnormality) that have been correctly labelled as positive.
- **Specificity:** proportion of negative images (i.e. not having an abnormality) that have been correctly labelled as negative.
- **ROC curve:** This curve is created by plotting the sensitivity (also called the True Positive Rate) against the False Positive Rate (1 - specificity) at various threshold settings. A shaded area is added to the curve that indicates the 95% confidence interval of the curve, computed using bootstrap analysis.
- **Az:** area under the ROC curve. This number is equivalent to the probability that a randomly picked positive case receives a higher score than a randomly picked negative case. It is bound between 0 and 1: at 0.5 the system is equivalent to guessing, at 1 the system shows perfect classification.
- **T:** Threshold value. Different threshold values correspond to different points on the ROC curve. The point on the ROC curve closest to perfect classification (the upper left corner) is often considered as the optimal threshold, but it is not necessarily the optimal threshold for the most cost-effective screening. **T** is the value set by the user to determine which images are labelled abnormal/normal.

An indication for the accuracy of a diagnostic test is the traditional academic point system:

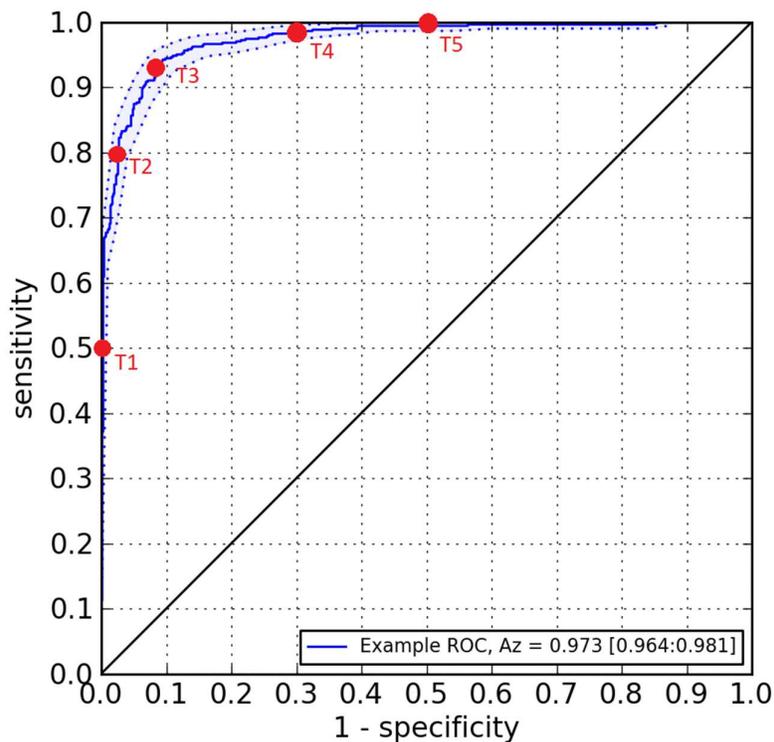


- 0.90-1 = excellent (A)
- 0.80-0.90 = good (B)
- 0.70-0.80 = fair (C)
- 0.60-0.70 = poor (D)
- 0.50-0.60 = fail (F)

In the curve shown on the left, the Az value is 0.973 which according to the above classification would be considered excellent (A).

Table 1: Different threshold values with corresponding sensitivity and specificity levels.		
Threshold	True Positive Rate (sensitivity)	True Negative Rate (specificity)
T1	50 %	100 %
T2	80 %	98 %
T3	93 %	90 %
T4	98 %	70 %
T5	100 %	50 %

This can be presented graphically into an ROC curve:



A relatively low threshold value of the software corresponds with a higher sensitivity, but at the cost of a lower specificity. A relatively high threshold value of the software corresponds with a higher specificity, but at the cost of a lower sensitivity. Hence, the threshold value is a trade-off between sensitivity and specificity. Shaded regions around the ROC curve (shown later in this report) indicate the 95% confidence intervals as computed using a statistical procedure called bootstrapping.

To summarize: An ROC curve demonstrates several things/characteristics of the test:

- ✓ It shows the trade-off between sensitivity and specificity (any increase in sensitivity will often be accompanied by a decrease in specificity).
- ✓ The closer the curve follows the left-hand border and then the top border of the ROC space, the more accurate the test.
- ✓ The closer the curve comes to the 45-degree diagonal of the ROC space, the less accurate the test.
- ✓ The area under the curve is a measure of the test's performance.

RetCAD: How does it work?

RetCAD is software based on convolutional neural networks, a state-of-the-art technique in machine learning. In the process of analyzing the input CF image, it compares regions in the image with regions extracted from normal and abnormal CF images. These images form the training data set of the software. The software is always tested on independent data: the test or validation set.

CF cameras from different manufacturers produce images of different quality because of hardware differences. In addition, image acquisition protocols can vary across acquisition sites, for example: the illumination, angular resolution (field of view) and the resolution of the image can vary. Furthermore, the patients may originate from different populations in which the appearance of the retina, such as color and pigmentation, may vary. In some patients the fluid in their eyeballs is not clear and this can make it difficult to make a good quality image. If a patient blinks during the acquisition of an image, the image may be substandard.

Specific algorithms to improve and normalize the input CF image prior to analysis are included in the RetCAD software. However, these algorithms are not perfect and cannot produce a high quality image if the quality of the input image is too low. Therefore, a quality measure for each image is also computed and the user of the software could decide to obtain a new image in case the quality is considered low by RetCAD.

The output scores of RetCAD are based on the AREDS (AMD) and ICDR (DR) grading protocols and can be roughly divided into the following classes:

AMD (according to the AREDS grading protocol)

- 0-0.5: No AMD
- 0.5-1.5: Early AMD
- 1.5-2.5: Intermediate AMD
- 2.5-3.0: Advanced AMD (both dry and wet-form)

DR (according to the ICDR grading protocol)

- 0-0.5: No DR
- 0.5-1.5: Mild DR
- 1.5-2.5: Moderate DR
- 2.5-3.5: Severe DR
- 3.5-4.0: Proliferative DR

For GLC, the vertical cup-to-disk-ratio (VCDR) and the GLC_score are computed and serve as an indication for the presence of GLC. These scores can roughly be divided as follows:

- <0.5: No suspicion of GLC
- ≥0.5: Suspicion of GLC

For the detection of referable DR, referable AMD and glaucoma, we generally recommend the following thresholds:

- AMD: 1.5
- DR: 1.5
- GLC (VCDR or GLC score): 0.5

RetCAD: Performance evaluation

The RetCAD software has been evaluated on several datasets. The images in these datasets were acquired using different types of CF cameras at different resolutions. The performance of the RetCAD software is directly compared with that of human experts. The following sections describe evaluations on various datasets.

Messidor

The Messidor database is a publicly available set of 1200 CF images which were acquired by three ophthalmologic departments using a color video 3CCD camera on a Topcon TRC NW6 non-mydratic retinography with a 45 degree field of view. The images were captured using 8 bits per color plane at 1440x960, 2240x1488, or 2304x1536 pixels. 800 images were acquired with pupil dilation (one drop of Tropicamide at 0.5%) and 400 without dilation. More information about the database can be found following the website link¹.

For each image in the database a reference DR severity grade, set by medical experts, was provided. Four severity grades were used: No DR, mild DR, moderate DR and severe DR.

The RetCAD software was applied to each of the 1200 images in the dataset and the RetCAD software was evaluated by comparing the RetCAD DR score with the DR severity grade as set by the medical experts.

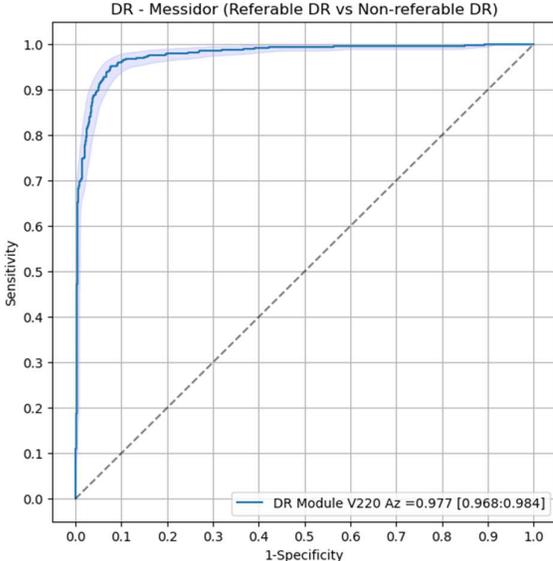
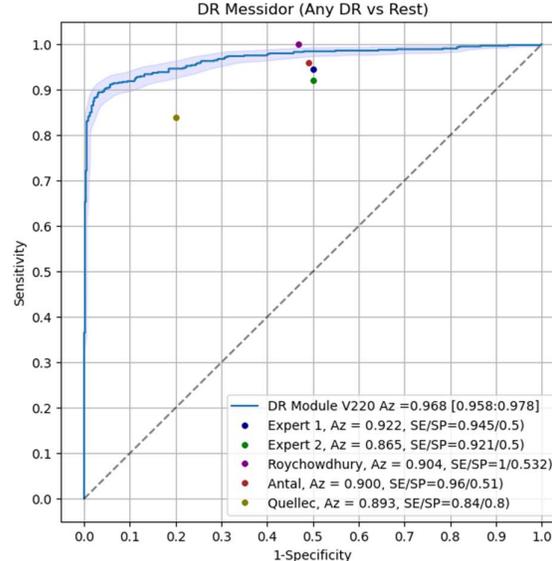
Two types of evaluations were performed:

1. The image was deemed positive if the reference grade was **referable DR**, i.e. severity level of moderate or severe DR.
2. The image was deemed positive if the reference grade was **any DR**, i.e. mild, moderate or severe DR.

The results of the evaluation are summarized in an ROC graph. In this ROC graph, the operating point of human experts is added for the second evaluation.

Note that the RetCAD software was not trained with any of the images that are part of the Messidor data set.

¹ Kindly provided by the Messidor program partners (see <http://www.adcis.net/en/DownloadThirdParty/Messidor.html>).

Reference: Referable DR	Reference: Any DR
Number of negative images: 699 Number of positive images: 501	Number of negative images: 546 Number of positive images: 654
	
Conclusion: The RetCAD software achieves an Az value of 0.977 for the identification of images with referable DR.	Conclusion: The RetCAD software achieves an Az value of 0.968 for the identification of images with any DR.

Operating points for RetCAD DR

In Table 2, sensitivity and specificity values of RetCAD referable DR detection are given at several threshold values for this specific dataset. Threshold with * is defined as the optimal threshold (best tradeoff sensitivity / specificity) for this dataset. DR Thresholds that can be used to retrieve the DR classification as presented in the section “RetCAD: How does it work?” are 0.5, 1.5, 2.5 and 3.5.

Threshold	True Positive Rate (sensitivity)	True Negative Rate (specificity)
0.25	98.2 %	74.0 %
0.50	97.6 %	82.4 %
0.94*	95.0 %	92.4 %
1.50	80.8 %	97.6 %
2.50	25.5 %	99.7 %
3.50	3.2%	100.0 %

Comparison with other systems and human experts

Several scientific publications have presented DR detection systems that were evaluated in the Messidor data set. One publication also reported the sensitivity/specificity for two human experts. All studies use the criteria of “any DR” for positive cases, i.e. mild or more severe are considered as the positive class. The table below reports the performances of the computer systems, including RetCAD DR, and the two human experts.

Author	Az value	Se/Sp	Year	Link
RetCAD DR	0.968	0.95/0.92	2023	-
Antal et al.	0.900	0.96/0.51	2012	http://arxiv.org/abs/1410.8577
Quellec et al.	0.893	0.84/0.80	2016	http://www.ncbi.nlm.nih.gov/pubmed/26774796

Roychowdhury et al.	0.904	1.00/0.53	2014	http://www.ncbi.nlm.nih.gov/pubmed/25192577
Sánchez et al.	0.876	0.92/0.50	2011	http://www.ncbi.nlm.nih.gov/pubmed/21527381
Expert 1	0.922	0.95/0.50	2011	http://www.ncbi.nlm.nih.gov/pubmed/21527381
Expert 2	0.865	0.91/0.50	2011	http://www.ncbi.nlm.nih.gov/pubmed/21527381

Messidor-2

The Messidor-2 dataset is a collection of Diabetic Retinopathy (DR) examinations, each consisting of two macula-centered eye fundus images (one per eye). Part of the dataset (*Messidor-Original*) was kindly provided by the Messidor program partners (see <http://messidor.crihan.fr>). The remainder (*Messidor-Extension*) consists of examinations obtained from the Brest University Hospital.

In the original Messidor dataset, some fundus images came in pairs (one image of both the left and right eye), some others were single (one image per patient). *Messidor-Original* consists of all image pairs from the original Messidor dataset, that is 529 examinations (1058 images).

In order to populate *Messidor-Extension*, diabetic patients were recruited in the Ophthalmology department of Brest University Hospital (France) between October 16, 2009 and September 6, 2010. Eye fundi were imaged, without pharmacological dilation, using a Topcon TRC NW6 non-mydratic fundus camera with a 45 degree field of view. Only macula-centered images were included in the dataset. *Messidor-Extension* contains 345 examinations (690 images).

Overall, Messidor-2 contains 874 examinations (1748 images). All patients in the database were graded for the presence of referable DR, i.e. moderate or more DR, by three medical experts and a consensus reference was made based on these gradings. In total, 190 patients had referable DR, and 684 patients did not have referable DR.

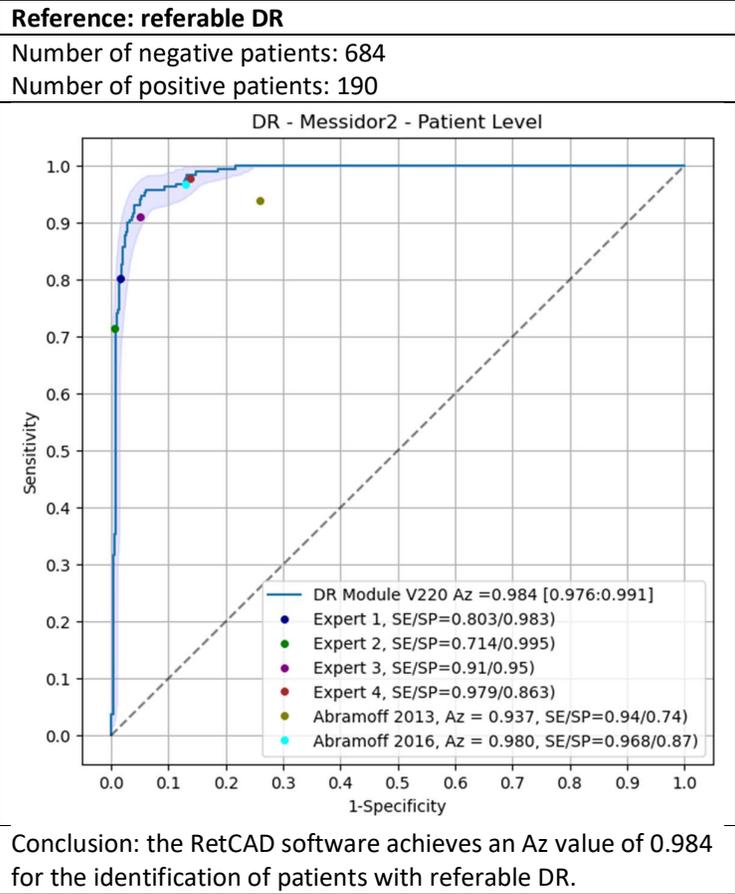
More information about the database can be found following the website link².

The RetCAD software was applied to each of the 1748 images in the dataset and the DR component of the RetCAD software was evaluated. In the evaluation, the highest score of the two images of a patient was set to be the patient-based score for DR. This score is compared with the provided reference score as constructed by a consensus of three medical experts (<https://www.ncbi.nlm.nih.gov/pubmed/27701631>).

The results of the evaluation are summarized in an ROC graph. Sensitivity and specificity of the three medical experts who scored the 874 examinations were measured by comparing the score to the consensus scoring of the other two human experts. The operating points of the human experts are added in the plot, but it thus has to be noted these were measured against a slightly different reference standard.

Note that the RetCAD software was not trained with any of the images that are part of the Messidor2 data set.

² Kindly provided by the LaTIM laboratory (see <http://latim.univ-brest.fr/>) and the Messidor program partners (see <http://messidor.crihan.fr/>)



Operating points for RetCAD DR

In Table 4, sensitivity and specificity values of RetCAD for DR detection are given at several threshold values for this specific dataset. Threshold with * is defined as the optimal threshold (best tradeoff sensitivity / specificity) for this dataset. DR Thresholds that can be used to retrieve the DR classification as presented in the section “RetCAD: How does it work?” are 0.5, 1.5, 2.5 and 3.5.

Table 4: Operating points of RetCAD for referable DR detection

Threshold	True Positive Rate (sensitivity)	True Negative Rate (specificity)
0.25	100.0 %	57.9 %
0.50	100.0 %	72.1 %
1.50	95.8 %	91.2 %
1.64*	95.8 %	94.0 %
2.50	41.1 %	99.3 %
3.50	7.4 %	99.6 %

Comparison with other systems and human experts

Performance of other state-of-the-art DR detection systems on the Messidor2 database have been reported. Additionally, the performance of human graders were reported in one of these publications (Abramoff et al, 2013) and were added.

Table 5: Performance of other software packages and human experts on the Messidor-2 dataset

Author	Az value	Se/Sp	Year	Link
RetCAD DR	0.984	0.96/0.94	2023	-

Abramoff et al.	0.937	0.97/0.59	2013	https://www.ncbi.nlm.nih.gov/pubmed/23494039
Abramoff et al.	0.980	0.97/0.87	2016	https://www.ncbi.nlm.nih.gov/pubmed/27701631
Expert 1	-	0.80/0.98	2013	https://www.ncbi.nlm.nih.gov/pubmed/23494039
Expert 2	-	0.71/1.00	2013	https://www.ncbi.nlm.nih.gov/pubmed/23494039
Expert 3	-	0.91/0.95	2013	https://www.ncbi.nlm.nih.gov/pubmed/23494039
Expert 4	-	0.98/0.86	?	https://www.ncbi.nlm.nih.gov/pubmed/23494039

Private1 dataset

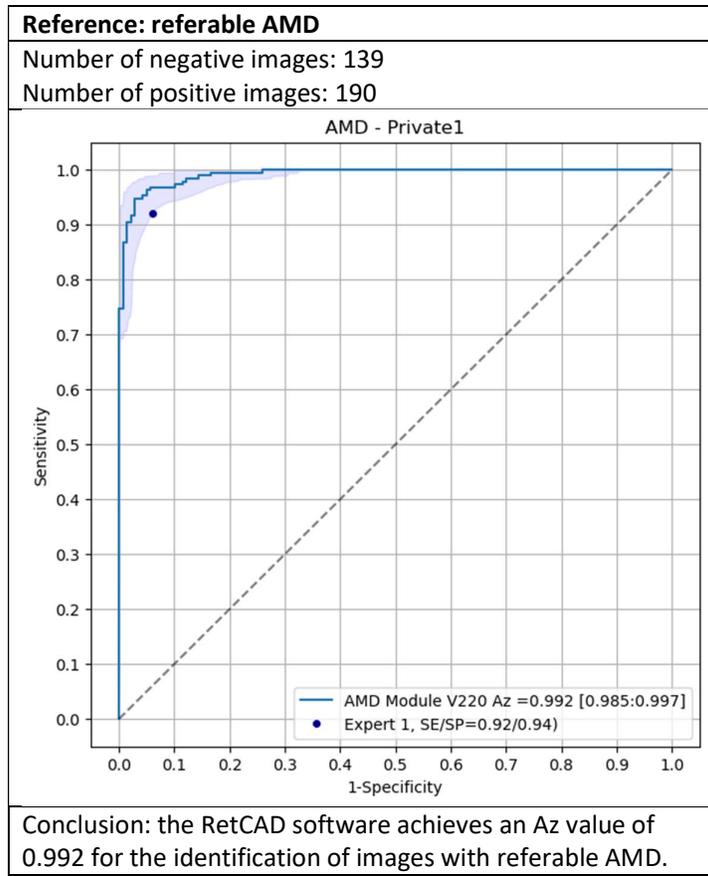
The private1 dataset is a dataset consisting of 329 macula centered images that were acquired at an ophthalmologic department in a hospital using either a Topcon TRC 501X model digital fundus camera at 50 degree field of view or a Canon CR-Dgi model non-mydratic retinal camera at 45 degree field of view. Pupil dilation was achieved with 1.0% tropicamide and 2.5% phenylephrine. All images were macula centered and image resolution varied between 1360x1024 to 3504x2336 pixels.

The images in the database were graded for presence of referable AMD by an expert with over 5 years of experience in grading fundus photographs. Referable AMD is defined as having at least 15 small drusen (>63µm) or more than one intermediate sized drusen (>126µm) or any sign of advanced AMD.

The RetCAD software was applied to each of the 329 images in the dataset and the AMD component of the RetCAD software was evaluated by comparing the RetCAD AMD score with the reference as set by the expert.

The results of the evaluation are summarized in an ROC graph. In this ROC graph, the operating point of a second human expert (over 5 years of experience in grading fundus images) is added for comparison.

Note that the RetCAD software was not trained with any of the images that are part of this data set.



Operating points for RetCAD AMD

In Table 6, sensitivity and specificity values of RetCAD for AMD detection are given at several threshold values for this specific dataset. Threshold with * is defined as the optimal threshold (best tradeoff sensitivity / specificity) for this dataset. AMD Thresholds that can be used to retrieve the AMD classification as presented in the section “RetCAD: How does it work?” are 0.5, 1.5 and 2.5

Table 6: Operating points of RetCAD for AMD detection

Threshold	True Positive Rate (sensitivity)	True Negative Rate (specificity)
0.25	99.5 %	81.3 %
0.50	97.4 %	88.5 %
0.91*	94.7 %	97.1 %
1.50	81.6 %	99.3 %
2.50	31.1 %	100.0 %

Comparison with other systems and human experts

Table 7: Performance of other software packages and human experts on the private1 dataset

Author	Az value	Se/Sp	Year	Link
RetCAD AMD	0.992	0.95/0.97	2023	-
Expert 1	-	0.92/0.94	2018	-

Mixed AMD-DR dataset

The Mixed AMD-DR dataset is a dataset consisting of 600 images that were acquired at an ophthalmologic department in a hospital using a Canon CR-2PlusAF digital fundus camera at 45 degree field of view. No pupil dilating eye-drops were administered. Image resolution varied between 2376x1584 to 3456x5184 pixels. The patients that were imaged had either signs of AMD, or DR, or both, or they were not affected by either disease.

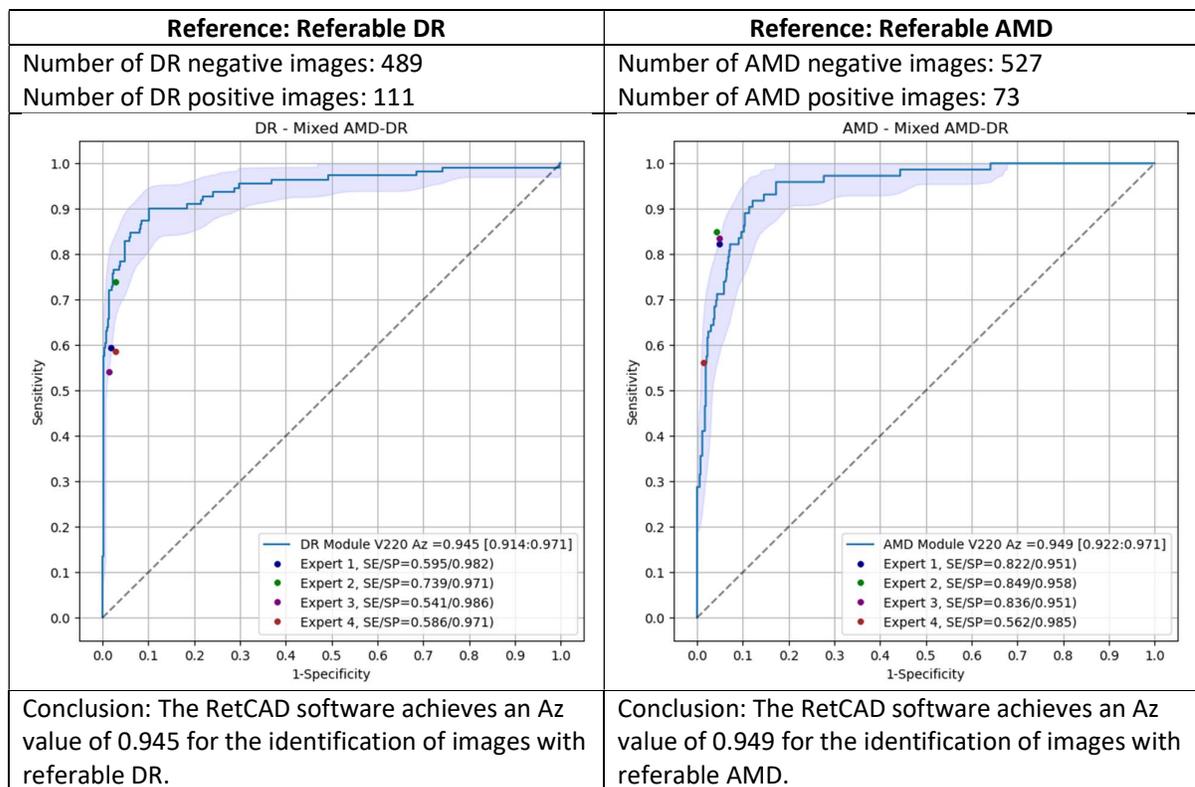
The reference for the images in this Mixed AMD-DR dataset was set by an experienced ophthalmologist. Grading criteria were based on the ICDR and AREDS classifications. In total 73 images were graded as having referable AMD (defined as intermediate or worse AMD), 111 were graded as referable DR (defined as moderate or worse DR), 3 were graded as both referable AMD and DR, and 408 were graded as neither referable AMD nor referable DR.

The RetCAD software was applied to each of the 600 images in the dataset and both the AMD and DR components of the RetCAD software were evaluated by comparing the RetCAD scores with the reference as set by the medical expert.

The results of the evaluation are summarized in two ROC graphs, one for AMD and one for DR.

Note that the RetCAD software was not trained with any of the images that are part of this data set.

The work has been published in a scientific journal publication in Acta Ophthalmologica:
<https://doi.org/10.1111/aos.14306>



Operating points for RetCAD AMD/DR

In Table 8, sensitivity and specificity values of RetCAD for AMD and DR detection are given at several threshold values for this specific dataset. Thresholds with * and ** are defined as the optimal threshold (best tradeoff sensitivity / specificity) for this dataset for DR and AMD respectively. DR Thresholds that can be used to retrieve the DR classification as presented in the section “RetCAD: How does it work?” are 0.5, 1.5, 2.5 and 3.5. AMD Thresholds that can be used to retrieve the AMD classification as presented in the section “RetCAD: How does it work?” are 0.5, 1.5 and 2.5.

Threshold	True Positive Rate DR (sensitivity)	True Negative Rate DR (specificity)	True Positive Rate AMD (sensitivity)	True Negative Rate AMD (specificity)
0.25	90.1 %	89.6 %	97.3 %	68.3 %
0.30*	90.1 %	89.8 %	N/A	N/A
0.50	87.4 %	90.8 %	95.9 %	72.3 %
1.50	76.6 %	96.5 %	91.8 %	85.6 %
1.74**	N/A	N/A	91.8 %	87.9 %
2.50	35.1 %	99.8 %	52.1 %	98.1 %
3.50	22.5 %	99.8 %	N/A	N/A

Comparison with other systems and human experts

Author	Az value	Se/Sp	Year	Link
RetCAD AMD	0.949	0.92/0.88	2023	-
RetCAD DR	0.945	0.90/0.90	2023	-
Expert 1	-	AMD: 0.82/0.95, DR: 0.60/0.98	2020	https://doi.org/10.1111/aos.14306

Expert 2	-	AMD: 0.85/0.96, DR: 0.74/0.97	2020	https://doi.org/10.1111/aos.14306
Expert 3	-	AMD: 0.84/0.95, DR: 0.54/0.99	2020	https://doi.org/10.1111/aos.14306
Expert 4	-	AMD: 0.56/0.99, DR: 0.59/0.97	2020	https://doi.org/10.1111/aos.14306

ORIGA dataset

The ORIGA dataset is a dataset consisting 650 images that were collected in a population based study, Singapore Malay Eye Study (SiMES)³. This study aims to assess the causes and risk factors of blindness and visual impairment in the Singapore Malay community. It was conducted over a 3 year period from 2004 to 2007 by Singapore Eye Research Institute and funded by the National Medical Research Council. SiMES examined 3,280 Malay adults aged 40 to 80, of which 149 are glaucoma patients. Retinal fundus images for both eyes were taken for each subject in the study. All retinal images have been de-identified by removing any individually identifiable information before being deposited to ORIGA.

The dataset consists of 650 annotated retinal images⁴. Each image is tagged with grading information such as cup-to-disk ratio, ISNT rule, disc haemorrhage, RNFL defects and glaucoma state. All images have image resolution of 3072x2048 pixels.

In total 168 images were graded as glaucoma suspicious whereas 482 were graded as non-glaucoma suspicious.

The RetCAD software was applied to each of the 650 images in the dataset and the GLC components of the RetCAD software, GLC score and VCDR, were evaluated by comparing the RetCAD scores with the reference provided with the dataset.

The results of the evaluation are summarized in a ROC graph. The performance measures of other systems have been added to the ROC. However, it has to be noted that these systems were designed by using images from the ORIGA dataset during training in a cross-validation setup. This results in a positive bias for those systems.

Note that the RetCAD software was not trained with any of the images that are part of this dataset.

³ T.Y. Wong, "Prediction of Diseases via Ocular Imaging: The Singapore Retinal Archival and Analysis Imaging Network", Inaugural Ocular Imaging Symposium, June 2008.

⁴ Z. Zhang, F.S. Yin, J. Liu, W.K. Wong, N.M. Tan, B.H. Lee, J. Cheng, T.Y. Wong: ORIGA(-light): an online retinal fundus image database for glaucoma analysis and research. <https://pubmed.ncbi.nlm.nih.gov/21095735/>

Reference: Glaucoma suspicious	Reference: Glaucoma suspicious
Number of negative images: 482 Number of positive images: 168	Number of negative images: 482 Number of positive images: 168
Conclusion: The RetCAD software with its VCDR Module achieves an Az value of 0. 812 for the identification of images with suspicious Glaucoma	Conclusion: The RetCAD software with its GLC_Score Module achieves an Az value of 0.838 for the identification of images with suspicious Glaucoma

Operating points for RetCAD GLC

In Table 10, sensitivity and specificity values of RetCAD for GLC are given at several threshold values for this specific dataset. Thresholds with * and ** are defined as the optimal threshold (best tradeoff sensitivity / specificity) for this dataset and the GLC (VCDR) and GLC (GLC score) Module respectively. GLC Threshold that can be used to retrieve the GLC classification as presented in the section “RetCAD: How does it work?” is 0.5.

Threshold	True Positive Rate VCDR (sensitivity)	True Negative Rate VCDR (specificity)	True Positive Rate GLC_Score (sensitivity)	True Negative Rate GLC_Score (specificity)
0.05**	100.0 %	0.00 %	75.0%	80.5%
0.30	100.0 %	0.00 %	41.1%	94.2%
0.40	98.8 %	14.1 %	33.9%	96.1%
0.50	84.5 %	64.7 %	24.4%	97.7%
0.53*	73.2 %	74.1 %	23.2%	98.1%
0.60	30.4 %	94.2 %	19.0%	98.3%
0.70	4.8 %	99.8 %	15.5%	98.8%

Comparison with other systems and human experts

Author	Az value	Se/Sp	Year	Link
RetCAD GLC (VCDR)	0.812	0.85/0.65	2023	-
RetCAD GLC (GLCscore)	0.838	0.24/0.98	2023	-
Chen et al.	0.831	-	2015	https://pubmed.ncbi.nlm.nih.gov/26736362/
Cheng et al.	0.838	-	2016	https://pubmed.ncbi.nlm.nih.gov/28268570/
Xu et al.	0.823	0.58/0.85	2013	https://link.springer.com/chapter/10.1007/978-3-642-40760-4_56
Fu et al.	0.851	-	2018	https://ieeexplore.ieee.org/document/8252743
Bajwa et al.	0.874	0.71/0.85	2019	https://bmcmedinformdecismak.biomedcentral.com/articles/10.1186/s12911-019-0842-8
Guo et al.	0.831	-	2018	https://pubmed.ncbi.nlm.nih.gov/32646472/
Sreng et al.	0.889	-	2020	https://www.mdpi.com/2076-3417/10/14/4916

REFUGE dataset

The REFUGE challenge database⁵ consists of 1200 retinal color fundus images stored in JPEG format, with 8 bits per color channel, acquired by ophthalmologists or technicians from patients sitting upright and using one of two devices: a Zeiss Visucam 500 fundus camera with a resolution of 2124 × 2056 pixels (400 images) and a Canon CR-2 device with a resolution of 1634 × 1634 pixels (800 images). The images are centered at the posterior pole, with both the macula and the optic disc visible, to allow the assessment of the ONH and potential retinal nerve fiber layer (RNFL) defects. These pictures correspond to Chinese patients (52% and 55% female in offline and online test sets, respectively) visiting eye clinics, and were retrieved retrospectively from multiple sources, including several hospitals and clinical studies. Only high-quality images were selected to ensure a proper labelling, and any personal and/or device information was removed for anonymization.

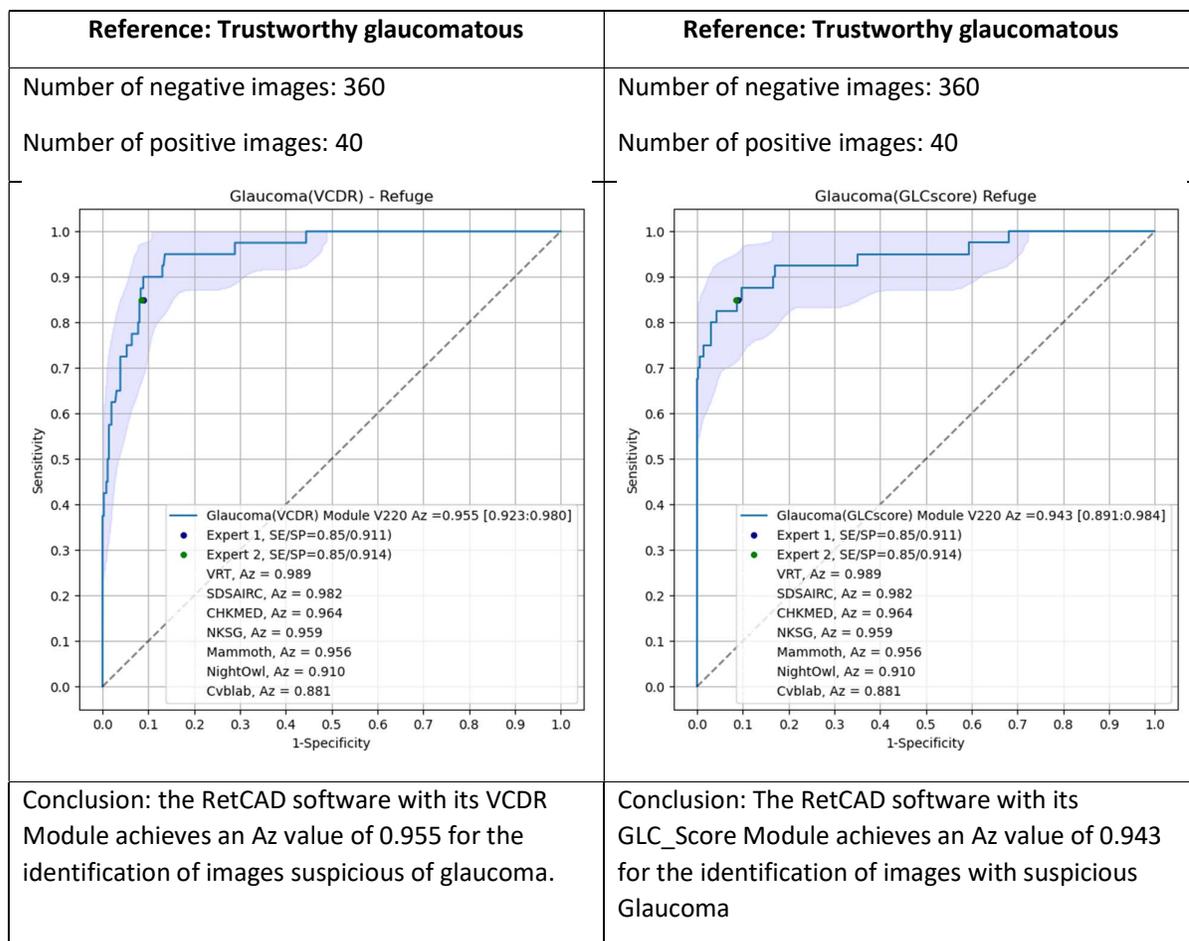
The test set of the REFUGE dataset has been used to evaluate RetCAD. This set consists of 400 images, acquired with the Canon CR-2 device and were provided with a reference glaucoma grading: trustworthy glaucomatous or non-glaucomatous. In total, 40 images were labeled as glaucomatous and 360 images were labeled as non-glaucomatous. These diagnostics were assigned based on the comprehensive evaluation of the subjects' clinical records, including follow-up fundus images, IOP measurements, optical coherence tomography images and visual fields (VF). The glaucomatous cases correspond to subjects with glaucomatous damage in the ONH area and reproducible glaucomatous VF defects.

The RetCAD software was applied to each of the 400 images in the dataset and the GLC component of the RetCAD software, VCDR and GLCscore, were evaluated by comparing the RetCAD scores with the reference provided with the dataset.

The results of the evaluation are summarized in a ROC graph. The performance measures of other systems have been added to the ROC.

Note that the RetCAD software was not trained with any of the images that are part of this dataset.

⁵ J.I. Orlando et al, "REFUGE Challenge: A unified framework for evaluating automated methods for glaucoma assessment from fundus photographs", Medical Image Analysis, vol 59, 2020



Operating points for RetCAD GLC

In Table 12, sensitivity and specificity values of RetCAD for GLC are given at several threshold values for this specific dataset. Thresholds with * and ** are defined as the optimal threshold (best tradeoff sensitivity / specificity) for this dataset and the GLC (VCDR) and GLC (GLCscore) Module respectively. GLC Threshold that can be used to retrieve the GLC classification as presented in the section “RetCAD: How does it work?” is 0.5.

Threshold	True Positive Rate VCDR (sensitivity)	True Negative Rate VCDR (specificity)	True Positive Rate GLC_Score (sensitivity)	True Negative Rate GLC_Score (specificity)
0.005**	100.0 %	0.0 %	85.0%	91.1%
0.30	100.0 %	0.0 %	52.5%	100%
0.40	100.0 %	21.1 %	45.0%	100%
0.50	95.0 %	75.8 %	42.5%	100%
0.55*	90.0 %	91.1 %	32.5%	100%

0.60	57.5 %	98.3 %	32.5%	100%
0.70	22.5 %	100.0 %	27.5%	100%

Comparison with other systems and human experts

Table 13: Performance of other software packages and human experts on the Refuge dataset				
Author	Az value	Se/Sp	Year	Link
RetCAD GLC (VCDR)	0.955	0.95/0.86	2023	-
RetCAD GLC (GLC_Score)	0.943	0.43/1	2023	-
Expert 1	-	0.85/0.911	2020	https://doi.org/10.1016/j.media.2019.101570
Expert 2	-	0.85/0.914	2020	https://doi.org/10.1016/j.media.2019.101570
VRT	0.989	-	2020	https://doi.org/10.1016/j.media.2019.101570
SDSAIRC	0.982	-	2020	https://doi.org/10.1016/j.media.2019.101570
CHKMED	0.964	-	2020	https://doi.org/10.1016/j.media.2019.101570
NKSG	0.959	-	2020	https://doi.org/10.1016/j.media.2019.101570
Mammoth	0.956	-	2020	https://doi.org/10.1016/j.media.2019.101570
NightOwl	0.910	-	2020	https://doi.org/10.1016/j.media.2019.101570
Cvblab	0.881	-	2020	https://doi.org/10.1016/j.media.2019.101570