# Interpretable-by-design Deep Survival Analysis for Disease Progression Modeling

Julius Gervelmeyer[1(✉)], Sarah Müller[1], Kerol Djoumessi[1], David Merle[2], Simon J. Clark[2], Lisa Koch[1,3], and Philipp Berens[1(✉)]

[1] Hertie Institute for AI in Brain Health, University of Tübingen, Tübingen, Germany
{julius.gervelmeyer,philipp.berens}@uni-tuebingen.de
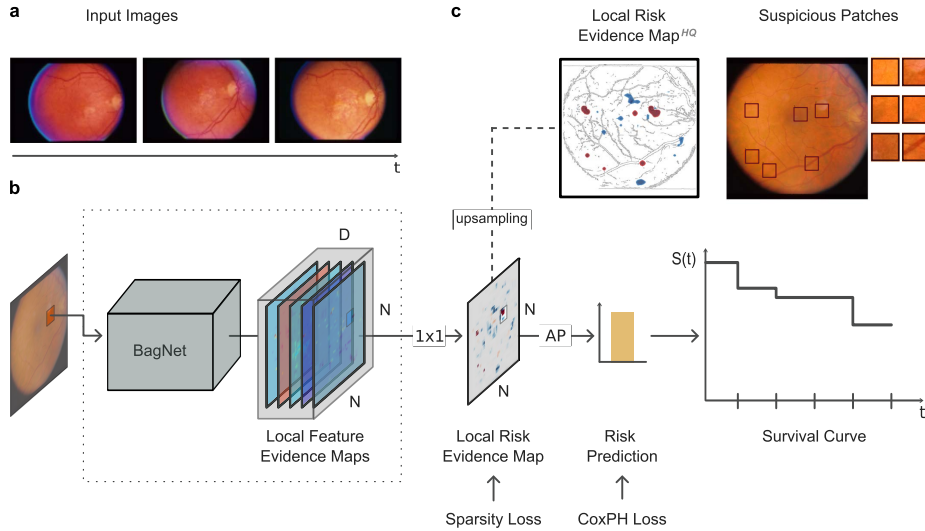[2] Department of Ophthalmology, University Eye Clinic, University of Tübingen, Tübingen, Germany
[3] Department of Diabetes, Endocrinology, Nutritional Medicine and Metabolism UDEM, Inselspital, Bern University Hospital, University of Bern, Bern, Switzerland

**Abstract.** In the elderly, degenerative diseases often develop differently over time for individual patients. For optimal treatment, physicians and patients would like to know how much time is left for them until symptoms reach a certain stage. However, compared to simple disease detection tasks, disease progression modeling has received much less attention. In addition, most existing models are black-box models which provide little insight into the mechanisms driving the prediction. Here, we introduce an interpretable-by-design survival model to predict the progression of age-related macular degeneration (AMD) from fundus images. Our model not only achieves state-of-the-art prediction performance compared to black-box models but also provides a sparse map of local evidence of AMD progression for individual patients. Our evidence map faithfully reflects the decision-making process of the model in contrast to widely used post-hoc saliency methods. Furthermore, we show that the identified regions mostly align with established clinical AMD progression markers. We believe that our method may help to inform treatment decisions and may lead to better insights into imaging biomarkers indicative of disease progression. The project's code is available at github.com/berenslab/interpretable-deep-survival-analysis.

**Keywords:** Interpretability · Deep survival analysis · Disease prognosis.

## 1 Introduction

Age-related macular degeneration (AMD) is the main cause of legal blindness in developed countries, caused by cumulative damage to the central retina leading to loss of central vision. This progressive disease severely impacts patients' quality of life by impairing tasks requiring sharp vision, motivating the need for early detection and intervention. AMD is characterised by retinal changes, such as drusen and pigment abnormalities, and is typically classified into early, intermediate and late disease stages, where vision is mostly endangered in the

**Fig. 1. Interpretable-by-design Deep Survival Model for AMD Progression.**
**a.** Fundus images of an eye at various screenings. **b.** An image is fed into the BagNet
survival model, which first yields local feature evidence maps (D the features dimension,
N its width and height) as in [6]. After a $1 \times 1$ convolution, we receive a local risk
evidence map, on which we apply a sparsity loss. Average pooling (AP) yields the
final risk prediction as in [10] on which we apply a CoxPH loss. In a second step,
the prediction is translated into a survival curve $S(t)$, the individual's probability of
"surviving" the event of interest – conversion to late AMD. **c.** The local risk evidence
map allows intuitive interpretation of predictions (left). Clinicians could be provided
with suspicious patches, i.e., regions that the model found to indicate a conversion risk
(right).

latter [12]. Retinal fundus images can reveal such indicators of AMD and con-
sequently, multiple studies have developed deep learning solutions to accurately
detect the disease stage [13,19] or predict the conversion to late AMD by a
specified time [4,25]. However, these black-box models do not provide inherent
interpretability and could only be explained by post-hoc saliency maps. These,
however, lack reliability, which is particularly problematic in high-risk medical
applications [3,22,21]. As an alternative, interpretable-by-design models such as
the Sparse BagNet [10] have recently been proposed for classification tasks on
retinal images. In this model architecture, evidence is gathered locally, repre-
sented in explicit evidence maps, which is then aggregated to predict the overall
class. The class evidence maps can be visualised as heatmaps overlaid on the
input image, highlighting the contribution of small local regions to the final pre-
diction. However, the Sparse BagNet was developed for a classification setting
and had not yet been adapted for disease progression modeling.

In this paper, we developed an interpretable-by-design model for AMD dis-
ease progression using a survival analysis framework. Survival analysis models

are a popular choice for time-to-event modeling which predict for each individual the probability to "survive" – i.e., to *not* observe – an event of interest within a time frame, such as progressing from early to late AMD. Classical survival analysis uses simple linear models such as the Cox proportional hazards (CoxPH) model [8], but deep survival models have been proposed that parameterise CoxPH models with neural networks [11,16,23]. In this work, we first integrated the current state-of-the-art model architectures for AMD progression modeling [4,25] into a deep survival model. We then replaced the backbone neural network with a Sparse BagNet, which provides an inherently interpretable evidence map for the predicted survival curve of an image, i.e. the risk of disease progression. To the best of our knowledge, this is the first image-based interpretable-by-design deep survival analysis model. We provide an overview over our method in Fig. 1.

## 2    Methods

### 2.1    Dataset and Preprocessing

We worked with data from the Age-Related Eye Disease Study (AREDS), a longitudinal study sponsored by the National Eye Institute, USA [1] which is available upon request[1] for research purposes. A total of 4,757 participants between 55 and 80 years of age were screened over a course of 12 years to study the natural progression of age-related eye diseases. The study was approved by the institutional review boards at all participating clinical sites and participants gave written consent [2]. We filtered out fundus images that were missing information such as the AMD severity score (1-9: increasing AMD severity; 10-12: neovascular AMD or central geographic atrophy, i.e., late AMD) and used the remaining 133,293 macula-centered fundus images for this project. Images came in pairs of photographs of the same fundus from slightly different angles to allow retinal specialists for depth impressions. We randomly selected one of the two views from each pair and excluded images after conversion to late AMD. Images were first resized to a height of 350 pixels and then cropped to a width of 350 pixels from the center. We further applied random resized cropping, flipping, color jitter, and rotation with the settings from [15], as these proved useful to enhance retinal disease detection. We split data into 60% training, 20% validation, and 20% test set, keeping a participant's records in the same split.

We defined the label targets *event* and *time* as follows: the event was 1 if the eye's AMD severity score reached 10 (late AMD) any time during the study and 0 otherwise (no late AMD). For this project, we defined time as the relative duration from an eye scan to the screening after which the eye was first diagnosed with late AMD or, if there was no such diagnosis, the time until the patient's last screening session. For time-dependent evaluation, we defined a "case" as an eye that converted to late AMD before or at that time.

---

[1] https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000001.v3.p1

## 2.2  Interpretable-by-design Deep Survival Model

**Deep Survival Analysis.** In survival analysis, CoxPH models provide a scalar output that can be used in a second step to estimate a set of risk predictions over time. We denote

$$\lambda(t|x) = \lambda_0(t) \cdot \exp h(x) \tag{1}$$

as the hazard function, where $\lambda_0(t)$ is the population wide baseline hazard function and $h(x)$ the individual log partial hazard based on a set of features $x$. A hazard at time $t$ is the instantaneous risk to observe an event of interest at (an infinitesimally small interval around) $t$. Classically, $h(x)$ is a linear function. For deep survival analysis, it can be parameterised by a neural network with any choice of architecture (ResNet-50, Inception-v3 or others) with one output node: we can apply a sigmoid function to the model's output logit and interpret the result as an estimate of the log partial hazard $\hat{h}(x)$. Here, this refers to the subject's relative risk to develop late AMD. This framework allows to directly work on fundus images as inputs $x$ and to train end-to-end. We estimated the baseline hazard function $\lambda_0(t)$ using Breslow's maximum likelihood approach [7] given the training data. The survival curve is given by

$$\hat{S}(t|x) = \exp - \int_0^t \hat{\lambda}(u|x)du \tag{2}$$

which represents the probabilities not to convert to late AMD up to time $t$.

**Adding Inherent Interpretability.** Instead of using a non-interpretable standard architecture to parameterise $h(x)$, we used a modified inherently interpretable model, the bag-of-local-features model, referred to as BagNet [6]. The BagNet is based on the ResNet-50 architecture, but changes in strides and the kernel sizes restrict the model to work on local image patches of size $q \times q$. After slight modification, the BagNet architecture can yield an explicit and local map of activations for each class: one pixel of each final feature map represents the local class evidence from an image patch in the input image [10]. Note that standard ResNets learn potentially global features and their interactions, while the BagNet learns the local evidence in an image patch. This eliminates the need for post-hoc saliency maps or the post-hoc analysis of convolutional filters. For the BagNet, we chose a receptive field size of $q = 33$ pixels and, following [10], we applied a sparsity constraint to the loss to avoid cluttered evidence maps (see Suppl. Fig. 1 for the selection of the sparsity coefficient). The model decision is then obtained by spatial average pooling, resulting in a final risk prediction logit. Correspondingly, the Sparse BagNet for survival analysis produces one class evidence map for the risk of disease progression and, as a result, allows a direct and intuitive interpretation of the survival predictions. In contrast, the baseline models are classification models that are trained once for each queried prediction time, resulting in a set of models, each of which has a saliency map to consider.

**Implementation Details.** We initialised the model with pre-trained weights from ImageNet and trained it for up to 50 epochs on the training set with a batch size of eight. Weights were updated based on the CoxPH loss [11,16], that is, the negative log likelihood loss of the log partial hazard $\hat{h}$, adapted from the auton-survival package [18]. We applied Breslow's method as implemented in scikit-survival [20]. We used the Adam optimiser at a learning rate of $1.6e$–5, as determined by a hyperparameter search. Training was conducted using a NVIDIA GeForce RTX 2080 Ti GPU and the PyTorch framework.

### 2.3   Baselines

Recently proposed end-to-end trained AMD progression models consist of one classification model for each queried time point [4,25,26]. As our baselines, we implemented classification models from Yan et al. [25] and Babenko et al. [4] that, similar to our model, are based only on fundus images. Both studies utilise Inception-v3 architectures to predict whether an eye converts to late AMD. However, Yan et al. use one fundus image per eye and screening as input, while Babenko et al. use both images of a stereo pair as inputs to separate Inception-v3 modules with shared weights and average their predictions after the activation function. Further, Yan et al. use the Adam optimiser and Babenko et al. use stochastic gradient descent. We re-implemented both models in our framework as close as possible to the originally reported settings. Training was conducted at a learning rate of $1e$–4, with the number of epochs, early stopping, batch size, image size and data augmentation set according to our proposed model. We used a binary cross-entropy loss and set the event label to 1 if the eye was first diagnosed with late AMD at or before the inquired year. If an eye did not convert to late AMD and the subject's last screening was before the inquired year, we could not extract classification labels and therefore had to exclude these records.

### 2.4   Evaluation Strategy

We evaluated the disease progression models using the area under the receiver operating characteristic curve (AUROC), Brier Score, and area under the precision-recall curve (AUPRC), adjusted for time-dependent predictions [17]. We chose the cumulative sensitivity/dynamic specificity approach to calculate AUROC using the scikit-survival package, which assesses the discrimination of eyes at higher risk from eyes at lower risk. Here, *cumulative cases* are subjects who have experienced the event up to a given time, while all eyes that have not (yet) experienced the event are referred to as *dynamic controls*. We computed Brier Scores as a measure of model calibration as implemented in scikit-survival and AUPRC as a performance measure focusing on cases using the R package time-ROC [5]. We additionally provide the metrics variants that are not adjusted for time-dependence in Suppl. Tab. 1. During training, we evaluated each epoch based on the Integrated Brier Score (IBS) of the survival probability predictions

**Table 1.** Performance of our interpretable-by-design model compared to black-box SOTA models on the unseen test data (see Methods for details).

| Model | Loss | AUROC ↑ | | Brier Score ↓ | | AUPRC ↑ | |
|---|---|---|---|---|---|---|---|
| | | Year 2 | Year 5 | Year 2 | Year 5 | Year 2 | Year 5 |
| Sparse BagNet (ours) | CoxPH | 0.941 | 0.938 | 0.036 | 0.054 | 0.642 | 0.676 |
| Babenko et al. [4] | Class. | 0.942 | 0.948 | 0.034 | 0.050 | 0.660 | 0.633 |
| Yan et al. [25] | Class. | 0.939 | 0.945 | 0.029 | 0.051 | 0.622 | 0.622 |

on the validation set. Training was stopped if the IBS did not improve within ten consecutive epochs. We selected the final model weights based on the epoch with the lowest IBS.
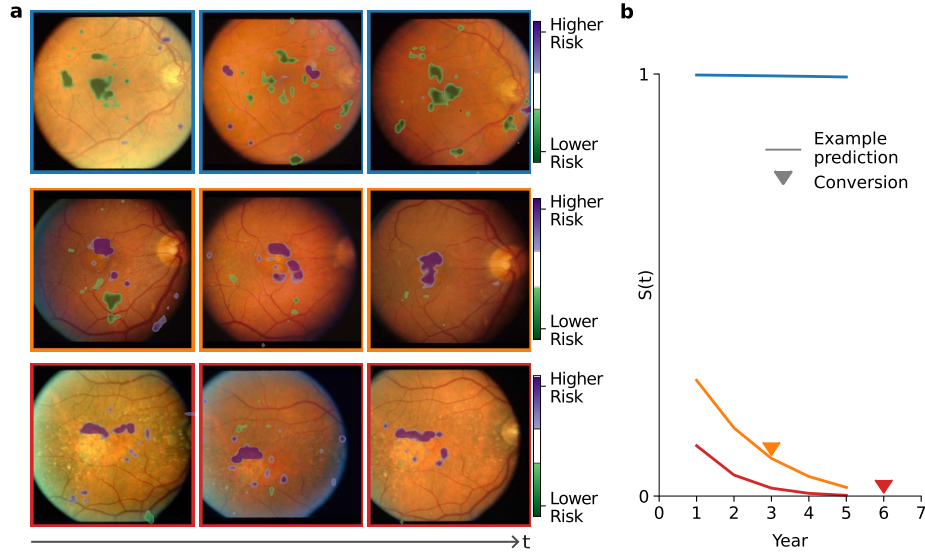
## 3   Results

We developed an architecture for disease progression modeling that achieves in-built interpretability through an explicit evidence map layer that contains the local evidence for disease progression. We applied the model to the prognosis of AMD progression and predicted the probability of not converting to late AMD within one to five years and compared it to state-of-the-art models.

### 3.1   Interpretable-by-design Model Achieves SOTA Performance

Our model performed comparably to state-of-the-art AMD progression models (Tab. 1) and only slightly worse in terms of AUROC and Brier Scores. We studied which of the components of our model were responsible for the performance difference compared to the baseline models. We found that training the Sparse BagNet as classification models reduced the performance, while replacing the Sparse BagNet with an unmodified ResNet-50 improved it (see Suppl. Fig. 2). This indicates that the CoxPH model training was not responsible for the slightly decreased performance, but rather the BagNet architecture.

### 3.2   Heatmaps Provide Faithful and Intuitive Model Interpretation

We extracted the evidence map from our model showing the local risk of conversion to late AMD – this map may show positive entries indicating higher risk of conversion in some regions or negative entries indicating lower risk. Crucially, the final risk prediction is simply formed by spatially averaging this evidence map so that it provides a faithful visualisation of the model's decision process when overlaid on an image (Fig. 2a). For example, the evidence map overlaid on fundus images of a non-converting participant at three subsequent screenings highlighted mostly regions indicating low risk of conversion (1<sup>st</sup> row of Fig. 2a) and the corresponding predicted survival curve stayed close to 1 (Fig. 2b). In contrast, evidence maps overlaid on fundus images of participants who converted to late AMD showed evidence for increased predicted risk in sizable portions
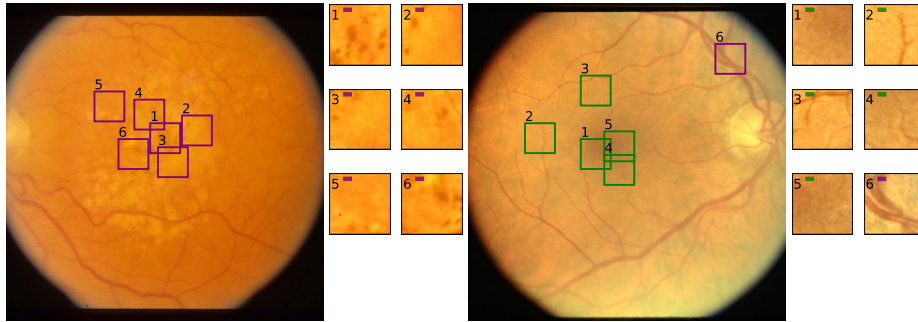
**Fig. 2. a.** Examples of predicted local evidence for the risk of AMD conversion for images from multiple consecutive screenings from a healthy eye (blue) and eyes that later converted to late AMD (orange, red). **b.** Example survival curves show predicted probabilities of surviving the event of interest – conversion to late AMD. An arrow indicates the time of late AMD onset.

of the fundus images (2nd and 3rd rows of Fig. 2a). Accordingly, the survival curve showed lower survival probability with conversion events indeed happening within a relatively short time frame of a few years (Fig. 2b). Remarkably, even though the image perspective shifted slightly between imaging time points the highlighted high risk regions were consistently identified (Fig. 2a).

Based on our evidence maps, the most important image regions can be bounded by boxes sized according to the model's receptive field and could be shown to clinicians to help them understand the prediction and potentially refine their own assessment (Fig. 3). In contrast, saliency maps for state-of-the-art models obtained using post-hoc gradient based techniques were often much less spatially confined (see Suppl. Fig. 3).

### 3.3    Heatmaps Capture Regions Known to Indicate AMD Progression

We next analysed to what extent the patches extracted from our evidence maps corresponded to known signs of conversion to late AMD. To this end, DM, a senior resident in Ophthalmology with experience in AMD research, annotated image patches. We included a random selection of the six patches with the highest predicted risk from 20 images from a pool of confident test set predictions for converters (for example Fig. 3). We found that our model focused mostly

**Fig. 3.** Two examples of interpretable model outputs. Patches show regions with high evidence for a risk to convert to late AMD (purple) and regions that provide evidence for a low AMD conversion risk (green). As the patch label number increases, the importance of the patch decreases.

on regions known to be associated with AMD, with 88.3% of the patches displaying known indications. Most prominently, 38/120 patches showed pigment mottling and 34/120 patches showed soft drusen, both indicative of AMD progression, and 32/120 patches already showed signs of atrophy. The patches shown in Fig. 3 (*left*) for instance contain either soft drusen, pigment mottling or both. In contrast, for non-converters, we could observe that the model focused on unremarkable retinal tissue in the macula or small blood vessels as indicative of low risk (Fig. 3, *right*).

## 4    Discussion

In this work, we introduced the first image-based interpretable-by-design deep survival model for modeling the risk of disease progression and applied it to the risk prediction of conversion to late AMD from fundus images. To this end, we combined a CoxPH survival model with a Sparse BagNet, which yielded highly localised evidence maps faithfully incorporated into the risk prediction process, which is desirable also for ethical reasons [14]. The high-risk areas identified by our model mostly corresponded to established signs of imminent conversion to late AMD. Our model uncovered the AMD risk areas without any prior knowledge from the AREDS dataset alone, indicating the model's potential for image-level biomarker discovery. In contrast, post-hoc saliency maps computed for state-of-the-art models were much less localised and do not yield faithful reflections on the model's decision making process [21]. Alternatives to the BagNet backbone include prototype models, which learn prototypical image parts and provide them for interpretability, and deep learning-based additive models such as the EPU-CNN [9], which provides contribution scores for colour and texture concepts along with their spatial relevance. These methods would be worth exploring in the context of disease prognosis. To date, however, prototype models

still suffer from imprecise explanations [24], and the EPU-CNN's interpretations require explicitly computed feature maps, typically based on luminance, color or frequency content. The Sparse BagNet, on the other hand, provides intuitive pixel importance for interpretation, making it well suited for clinical applications and finding novel image-based biomarkers. As the inductive bias of our model fits well to task structures involving small disease features, we expect it to generalise to other medical imaging tasks, including progression modeling for diabetic retinopathy. In summary, our interpretable-by-design deep survival model based on Sparse BagNets opens up new possibilities for trustworthy risk-modeling from medical images beyond ophthalmology and may help to identify new early indications of disease progression which could easily be overlooked by humans.

**Disclosure of Interests.** The authors declare no relevant competing interests.

# References

1. Age-Related Eye Disease Study Research Group: The Age-Related Eye Disease Study (AREDS): Design Implications AREDS Report No. 1. Controlled Clinical Trials **20**(6), 573–600 (Dec 1999). `https://doi.org/10.1016/S0197-2456(99)00031-8`

2. Age-Related Eye Disease Study Research Group: A Randomized, Placebo-Controlled, Clinical Trial of High-Dose Supplementation With Vitamins C and E, Beta Carotene, and Zinc for Age-Related Macular Degeneration and Vision Loss: AREDS Report No. 8. Archives of Ophthalmology **119**(10), 1417–1436 (Oct 2001). `https://doi.org/10.1001/archopht.119.10.1417`

3. Arun, N., Gaw, N., Singh, P., Chang, K., Aggarwal, M., Chen, B., Hoebel, K., Gupta, S., Patel, J., Gidwani, M., Adebayo, J., Li, M.D., Kalpathy-Cramer, J.: Assessing the Trustworthiness of Saliency Maps for Localizing Abnormalities in Medical Imaging. Radiology. Artificial Intelligence **3**(6), e200267 (Nov 2021). `https://doi.org/10.1148/ryai.2021200267`

4. Babenko, B., Balasubramanian, S., Blumer, K.E., Corrado, G.S., Peng, L., Webster, D.R., Hammel, N., Varadarajan, A.V.: Predicting Progression of Age-related Macular Degeneration from Fundus Images using Deep Learning (Apr 2019), arXiv:1904.05478 [cs.CV]

5. Blanche, P., Dartigues, J.F., Jacqmin-Gadda, H.: Estimating and comparing time-dependent areas under receiver operating characteristic curves for censored event times with competing risks. Statistics in Medicine **32**(30), 5381–5397 (2013). `https://doi.org/10.1002/sim.5958`

6. Brendel, W., Bethge, M.: Approximating CNNs with Bag-of-local-Features models works surprisingly well on ImageNet. In: International Conference on Learning Representations (Mar 2019)

7. Breslow, N.: Discussion of professor cox's paper. Journal of the Royal Statistical Society, Series B **34**, 216–217 (1972)

8. Cox, D.R.: Regression Models and Life-Tables. Journal of the Royal Statistical Society: Series B (Methodological) **34**(2), 187–202 (1972). `https://doi.org/10.1111/j.2517-6161.1972.tb00899.x`

9. Dimas, G., Cholopoulou, E., Iakovidis, D.K.: E pluribus unum interpretable convolutional neural networks. Scientific Reports **13**(1), 11421 (Jul 2023). `https://doi.org/10.1038/s41598-023-38459-1`, publisher: Nature Publishing Group

10. Djoumessi, K.R., Ilanchezian, I., Kühlewein, L., Faber, H., Baumgartner, C.F., Bah, B., Berens, P., Koch, L.M.: Sparse activations for interpretable disease grading. In: Medical Imaging with Deep Learning (2023)

11. Faraggi, D., Simon, R.: A neural network model for survival data. Statistics in Medicine **14**(1), 73–82 (1995). `https://doi.org/10.1002/sim.4780140108`

12. Fleckenstein, M., Keenan, T.D.L., Guymer, R.H., Chakravarthy, U., Schmitz-Valckenberg, S., Klaver, C.C., Wong, W.T., Chew, E.Y.: Age-related macular degeneration. Nature Reviews. Disease Primers **7**(1), 31 (May 2021). `https://doi.org/10.1038/s41572-021-00265-2`

13. Grassmann, F., Mengelkamp, J., Brandl, C., Harsch, S., Zimmermann, M.E., Linkohr, B., Peters, A., Heid, I.M., Palm, C., Weber, B.H.F.: A Deep Learning Algorithm for Prediction of Age-Related Eye Disease Study Severity Scale for Age-Related Macular Degeneration from Color Fundus Photography. Ophthalmology **125**(9), 1410–1420 (Sep 2018). `https://doi.org/10.1016/j.ophtha.2018.02.037`

14. Grote, T.: The Allure of Simplicity: On Interpretable Machine Learning Models in Healthcare. Philosophy of Medicine **4**(1) (Sep 2023). `https://doi.org/10.5195/pom.2023.139`, number: 1

15. Huang, Y., Lin, L., Cheng, P., Lyu, J., Tam, R., Tang, X.: Identifying the Key Components in ResNet-50 for Diabetic Retinopathy Grading from Fundus Images: A Systematic Investigation. Diagnostics **13**(10), 1664 (Jan 2023). `https://doi.org/10.3390/diagnostics13101664`, number: 10 Publisher: Multidisciplinary Digital Publishing Institute

16. Katzman, J.L., Shaham, U., Cloninger, A., Bates, J., Jiang, T., Kluger, Y.: DeepSurv: personalized treatment recommender system using a Cox proportional hazards deep neural network. BMC Medical Research Methodology **18**(1), 24 (Feb 2018). `https://doi.org/10.1186/s12874-018-0482-1`

17. Lambert, J., Chevret, S.: Summary measure of discrimination in survival models based on cumulative/dynamic time-dependent roc curves. Statistical methods in medical research **25**(5), 2088–2102 (2016)

18. Nagpal, C., Potosnak, W., Dubrawski, A.: auton-survival: an Open-Source Package for Regression, Counterfactual Estimation, Evaluation and Phenotyping with Censored Time-to-Event Data (Aug 2022). `https://doi.org/10.48550/arXiv.2204.07276`, arXiv:2204.07276 [cs, stat]

19. Peng, Y., Dharssi, S., Chen, Q., Keenan, T.D., Agrón, E., Wong, W.T., Chew, E.Y., Lu, Z.: DeepSeeNet: A Deep Learning Model for Automated Classification of Patient-based Age-related Macular Degeneration Severity from Color Fundus Photographs. Ophthalmology **126**(4), 565–575 (Apr 2019). `https://doi.org/10.1016/j.ophtha.2018.11.015`

20. Pölsterl, S.: scikit-survival: A library for time-to-event analysis built on top of scikit-learn. Journal of Machine Learning Research **21**(212), 1–6 (2020)
21. Rudin, C.: Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. Nature Machine Intelligence **1**(5), 206–215 (May 2019). `https://doi.org/10.1038/s42256-019-0048-x`, publisher: Nature Publishing Group
22. Saporta, A., Gui, X., Agrawal, A., Pareek, A., Truong, S.Q.H., Nguyen, C.D.T., Ngo, V.D., Seekins, J., Blankenberg, F.G., Ng, A.Y., Lungren, M.P., Rajpurkar, P.: Benchmarking saliency methods for chest X-ray interpretation. Nature Machine Intelligence **4**(10), 867–878 (Oct 2022). `https://doi.org/10.1038/s42256-022-00536-x`, number: 10 Publisher: Nature Publishing Group
23. Wiegrebe, S., Kopper, P., Sonabend, R., Bischl, B., Bender, A.: Deep Learning for Survival Analysis: A Review. Artificial Intelligence Review **57**(3), 65 (Feb 2024). `https://doi.org/10.1007/s10462-023-10681-3`
24. Xu-Darme, R., Quénot, G., Chihani, Z., Rousset, M.C.: Sanity checks and improvements for patch visualisation in prototype-based image classification (May 2023). `https://doi.org/10.48550/arXiv.2302.08508`, arXiv:2302.08508 [cs]
25. Yan, Q., Weeks, D.E., Xin, H., Swaroop, A., Chew, E.Y., Huang, H., Ding, Y., Chen, W.: Deep-learning-based Prediction of Late Age-Related Macular Degeneration Progression. Nature Machine Intelligence **2**(2), 141–150 (Feb 2020). `https://doi.org/10.1038/s42256-020-0154-9`
26. Yin, C., Moroi, S.E., Zhang, P.: Predicting Age-Related Macular Degeneration Progression with Contrastive Attention and Time-Aware LSTM. KDD: proceedings. International Conference on Knowledge Discovery & Data Mining **2022**, 4402–4412 (Aug 2022). `https://doi.org/10.1145/3534678.3539163`