

# Journal Pre-proof

Personalized model to predict keratoconus progression from demographic, topographic and genetic data

Howard P Maile , Ji-Peng Olivia Li , Mary D Fortune ,  
Patrick Royston , Marcello T Leucci , Ismail Moghul , Anita Szabo ,  
Konstantinos Balaskas , Bruce D Allan , Alison J Hardcastle ,  
Pirro Hysi , Nikolas Pontikos , Stephen J Tuft , Daniel M Gore

PII: S0002-9394(22)00150-7  
DOI: <https://doi.org/10.1016/j.ajo.2022.04.004>  
Reference: AJOPHT 12205

To appear in: *American Journal of Ophthalmology*

Received date: January 21, 2022  
Revised date: April 2, 2022  
Accepted date: April 13, 2022

Please cite this article as: Howard P Maile , Ji-Peng Olivia Li , Mary D Fortune , Patrick Royston , Marcello T Leucci , Ismail Moghul , Anita Szabo , Konstantinos Balaskas , Bruce D Allan , Alison J Hardcastle , Pirro Hysi , Nikolas Pontikos , Stephen J Tuft , Daniel M Gore , Personalized model to predict keratoconus progression from demographic, topographic and genetic data, *American Journal of Ophthalmology* (2022), doi: <https://doi.org/10.1016/j.ajo.2022.04.004>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2022 Published by Elsevier Inc.



**Personalized model to predict keratoconus progression from  
demographic, topographic and genetic data**

**Short title: Keratoconus progression**

Howard P Maile<sup>1\*</sup>; Ji-Peng Olivia Li<sup>2\*</sup>; Mary D Fortune<sup>3</sup>; Patrick Royston<sup>4</sup>;  
Marcello T Leucci<sup>2</sup>; Ismail Moghul<sup>1</sup>; Anita Szabo<sup>1</sup>; Konstantinos Balaskas<sup>2</sup>;  
Bruce D Allan<sup>2</sup>; Alison J Hardcastle<sup>1</sup>; Pirro Hysi<sup>5,6</sup>; Nikolas Pontikos<sup>1\*</sup>; Stephen J  
Tuft<sup>2\*</sup>; Daniel M Gore<sup>2\*§</sup>

<sup>1</sup> UCL Institute of Ophthalmology, 11-43 Bath Street, London EC1V 9EL

<sup>2</sup> Moorfields Eye Hospital NHS Foundation Trust, 162 City Road, London EC1V  
2PD

<sup>3</sup> MRC Biostatistics Unit, Cambridge Institute of Public Health, University of  
Cambridge, UK

<sup>4</sup> MRC Clinical Trials Unit at UCL, Aviation House, 125 Kingsway, London,  
WC2B 6NH

<sup>5</sup> Section of Ophthalmology, School of Life Course Sciences, King's College  
London

<sup>6</sup> Department of Twin Research and Genetic Epidemiology, King's College  
London

\* Authors contributed equally

§ Corresponding author: Mr Daniel Gore, Moorfields Eye Hospital, 162 City  
Road, London, EC1V 2PD, UK, [daniel.gore1@nhs.net](mailto:daniel.gore1@nhs.net), 0207 253 3411.

Keywords: Keratoconus, corneal cross-linking, keratoconus genetics, keratoconus prediction

Page Break

## Abstract

**Purpose:** To generate a prognostic model to predict keratoconus progression to corneal cross-linking (CXL).

**Design:** Retrospective cohort study.

**Methods:** We recruited 5025 patients (9341 eyes) with early keratoconus between January 2011 and November 2020. Genetic data from 926 patients was available. We investigated both keratometry or CXL as end-points for progression and used the Royston-Parmar method on the proportional hazards scale to generate a prognostic model. We calculated hazard ratios (HR) for each significant covariate, with explained variation and discrimination, and performed internal-external cross validation by geographic regions.

**Results:** After exclusions, model-fitting comprised 8701 eyes, of which 3232 underwent CXL. For early keratoconus, CXL provided a more robust prognostic model than keratometric progression. The final model explained 33% of the variation in time-to-event: age HR [95% confidence limits] 0.9 [0.90-0.91], maximum anterior keratometry (Kmax) 1.08 [1.07-1.09], and minimum corneal thickness 0.95 [0.93-0.96] as significant covariates. Single nucleotide polymorphisms (SNPs) associated with keratoconus (n=28) did not significantly contribute to the model. The predicted time-to-event curves closely followed the observed curves during internal-external validation. Differences in discrimination between geographic regions was low, suggesting the model maintained its predictive ability.

**Conclusions:** A prognostic model to predict keratoconus progression could aid patient empowerment, triage and service provision. Age at presentation is the most significant

predictor of progression risk. Candidate SNPs associated with keratoconus do not contribute to progression risk.

#### Precis

Corneal crosslinking is successful in halting keratoconus progression but providing patients with a personalized visual representation of risk to progression is desirable. This research presents a model to generate projected likelihood of having crosslinking from data collected at presentation. It was trained and validated from a large dataset of 8701 eyes from keratoconus patients. Univariable and multivariable analysis was performed to identify risk factors including single nucleotide polymorphisms associated with keratoconus.

#### Introduction

Keratoconus is a common corneal ectasia that causes irregular astigmatism, scarring and loss of vision. Thinning and steepening can progress through childhood and early adulthood, but the shape of most eyes stabilizes by the third or fourth decade. Without intervention, keratoconus can lead to severe visual loss, with approximately 10% of eyes eventually requiring corneal transplantation.<sup>1</sup> Corneal cross-linking (CXL) by topical application of riboflavin, followed by irradiation with UV-A light, can arrest progression of keratoconus in up to 88% to 100% of eyes even when there is relatively advanced disease.<sup>2-6</sup> The potential benefit of CXL is to prevent visual deterioration with a relatively low risk procedure that is cost effective for healthcare providers.<sup>7-9</sup> However, CXL is usually not offered to all patients at presentation because the disease may have already stabilized. In the recent KERALINK study 43% of children <17 years of age at presentation had not progressed after 18 months.<sup>10</sup> The definition of progression also varies with the severity of keratoconus, but for early disease a common threshold is either an increase in the maximum keratometry (Kmax) of >1 dioptre, a change in the manifest refractive spherical equivalent of >0.50 dioptre, or an increase in manifest refractive cylinder of >1 dioptre.<sup>2,11</sup> Depending on the rate of progression this threshold may be passed in a few months, years, or not at all. At the first assessment it can be a challenge to distinguish eyes that are at risk of rapid progression from those where it is safe to monitor. Unnecessary review visits are a burden to the patient and the care system.

We considered the date of numeric progression,<sup>6</sup> as well as the date when CXL was performed, as alternative end-points to define keratoconus progression. Although the use of keratometry as an end-point may appear the more objective method, there is variability on the definition of progression reported in the literature and conclusions may vary with the definition that is adopted.<sup>11-14</sup> Repeatability thresholds are not usually tailored to individual eyes (i.e. an increase in Kmax by 1 D is not significant in all eyes) although there is growing evidence on the variability of measurements in more advanced disease and the need for tailoring numerical progression definitions to the disease state, and distinguishing real progression from inherent variability of measurement modalities.<sup>15-17</sup> Finally, patients who receive CXL prior to progression must be censored from the dataset even though these eyes are likely to have been at risk of progression. This type of informative censoring creates a bias.<sup>18</sup> In contrast, the time to CXL depends on several variables that include numeric disease progression, but also incorporates patient-specific risk factors for future progression. Its strength is that it is an easily comprehensible and meaningful end-point for patients. It encompasses individual risk factors that are not considered when imaging is used in isolation and it has been used by others as defining the event of interest.<sup>19</sup>

For these reasons we have used demographic and serial tomography data from a large cohort of patients to generate a time-to-event model to predict the probability of an individual progressing to CXL. Because the Cox proportional hazards method does not generate smooth time-to-event curves, we used the Royston-Parmar model to achieve direction estimates of the hazard function.<sup>20</sup> We also performed a further analysis of a subset of patients who had genetic data in the form of single-nucleotide polymorphisms (SNP) generated as part of a study to determine keratoconus risk.<sup>21</sup>

## **Methods**

### **Cohort**

The study protocol was reviewed and approved by the Clinical Audit Assessment Committee of Moorfields Eye Hospital NHS Foundation Trust (reference CA17/CED/03). Institutional Review Board (IRB) approval was

obtained and individual patient consent was not required. The study conformed to the tenets of the Declaration of Helsinki. We identified from the Moorfields Eye Hospital electronic health record database (OpenEyes) patients aged 13 years and above diagnosed with clinical or suspected keratoconus who attended our Early Keratoconus Clinic (EKC) between January 2011 and November 2020. Clinical data included keratometry (Kmax, Front K1, Front K2, Back K1, Back K2), and pachymetry (minimum corneal thickness) captured by Scheimpflug tomography (Pentacam HR, Oculus GmbH, Wetzlar). We only included scans with a quality score of 'good' or 'ok', and where multiple scans were taken on the same day we used the mean value. The date of all CXL procedures was recorded. The protocol for offering CXL throughout the study period was, i) a documented history prior to referral to the EKC of our hospital of significant recent disease progression,<sup>6</sup> ii) a change in contemporary measurements of 95% above the repeatability limits of the baseline measurements as shown in Supplemental Table 1 (available at <http://www.ajo.com>),<sup>6</sup> or iii) a patient considered by a clinician to be at high risk of progression despite their not fulfilling the above two criteria. Exclusion criteria included pregnancy or breastfeeding, uncontrolled ocular surface disease or a minimum corneal thickness less than 375  $\mu\text{m}$ .

All the data used for model fitting started from the first appointment in the EKC. Patient demographics included age, gender, smoking status (current or ex/non-smoker) ethnicity and postcode. Ethnicity was coded as 1 for 'Black' or 'South Asian or South Asian British' and 0 for any other category (excluding missing values). Before model fitting, the pachymetry in microns was divided by 10 to generate a meaningful scale. For the primary analysis, eyes with any missing data were excluded. We also explored multiple imputation, which avoids data exclusion by generating multiple versions of the dataset with missing values replaced with values sampled from an appropriate distribution. To see whether genetic data can help predict keratoconus progression, we used 28 candidate SNPs from a recent keratoconus genome-wide association study that contained 926 patients from Moorfields Eye Hospital.<sup>21</sup> The SNP data was encoded as either 0 (homozygous reference genotype), 1 (heterozygous genotype), or 2 (homozygous variant genotype). We chose to use an additive encoding, thus

the risk of disease increases additively with the degree of genetic variation.<sup>22</sup> Anonymized data were then exported to Excel software for analysis (version 15.24 2016, Microsoft Corp.).

### **Model Fitting and Covariate Selection**

A Royston-Parmar flexible parametric survival model was fitted to the data to predict the probability of an eye progressing to CXL.<sup>23</sup> Initial analysis of the covariates was performed by univariate analysis using the same model characteristics as the multivariable model. When selecting covariates for the final multivariable model, we used backwards stepwise selection with a significance level of 0.05. We used linear covariates for ease of interpretation of our final model. To create a more parsimonious model we examined the effect on explained variation and discrimination of removing single variables from the model.

### **Keratometric Progression Sensitivity Analysis**

We included a sensitivity analysis in which we investigated keratometric progression as an alternative end-point. Keratometric progression was defined using thresholds from Gore et al 2021.<sup>6</sup> When using numerical thresholds to define progression, the appointments for eyes beyond the date of CXL cannot be used. However, censoring these eyes at the date of CXL represents informative censoring. Based on the recommendations of Clarke et al<sup>18</sup> for investigating the impact of informative censoring, we generated a 'best case' dataset where eyes were censored at the CXL date and a 'worst case' dataset where patients were assumed to progress at the CXL date. The corresponding Kaplan Meier curves were plotted to provide a visual comparison of the two datasets. A Royston-Parmar model was then fitted on both datasets. We used the same techniques (backwards stepwise selection, significance level of 0.05) as described in the previous section to fit the model and compare the explained variation and hazard ratios.

### **Multivariable Model Validation**

We validated the model using internal-external cross validation in which we split the dataset by geographical region.<sup>24,25</sup> For the  $k$ th region, the model is fitted on

the full dataset excluding region  $k$  and then Kaplan-Meier curves and predicted survival curves were generated for region  $k$ . Seven geographical regions were created based on the patient's postcode as shown in Supplemental Figure 1 (available at <http://www.ajo.com>). To quantitatively assess the validation, Royston and Sauerbrei's  $D$  statistic was calculated for both the model fitted from data excluding region  $k$  ( $D_{(-k)}$ ) and also the model applied to region  $k$  ( $D_k$ ).<sup>26</sup> The difference between these two discrimination metrics ( $D_k - D_{(-k)}$ ) was calculated with its corresponding standard error to assess the predictive ability of the model. To demonstrate how the model could be used in practice, we include three hypothetical patients' eyes with different progression risk profiles (high, medium, low risk) and plot the predicted time-to-event curve for each shown in Figure 2.

### Statistical Analysis

The event of interest was defined as the date that the eye underwent CXL. We calculated the time-to-event as the difference between the first appointment in our service and the date of CXL (or the last patient appointment in the case of censoring). Since we had paired observations (eyes), we used variance-corrected models to account for correlation between eyes and to ensure that robust standard errors were produced. The choice of scale and selection of degrees of freedom for the Royston-Parmar model was informed by inspecting the Akaike information criterion (AIC) and Bayes information criterion (BIC)<sup>20</sup> and the results of this were balanced with ease of interpretation. See Supplemental Table 2 and Supplemental Text 1 (available at <http://www.ajo.com>) for further explanation. Royston and Sauerbrei's  $D$  statistic was used as a measure of discrimination and  $R^2_D$  as a measure of explained variation (both calculated on the natural scale of the model). Although all of the primary results were generated from a complete case analysis, we performed an additional analysis using multiple chained imputation (predictive mean matching approach with 5 nearest neighbors). Model fitting was performed in Stata 13 (StataCorp LP, Texas, USA) and the Royston-Parmar model was fitted using the `stpm2` package from Stata 13.

### Results



## **Cohort**

From a potential of 9,341 eyes (4316 pairs of eyes and 709 individual eyes), the final model used 8,701 eyes of 4,823 patients, with 3,232 eyes that had CXL. The mean age was 28.3 years with standard deviation of 7.1 years. We excluded 640 eyes with missing data. Table 1 summarizes the available covariates along with missing data percentages. See Supplemental Text 2 and Supplemental Table 3 (available at <http://www.ajo.com>) for a description of the multiple imputation results.

## **Model Fitting and Covariate Selection (Genetic Data)**

We analyzed patients with genetic data separately because this data was only available for ~14% of patients. Of 926 patients (1852 eyes) with genetic data, 531 eyes were excluded with incomplete keratometry or CXL data, which left 1321 eyes, of which 665 had CXL. With univariate analysis of the 28 SNPs only rs72631889 was found to be significant ( $P=0.01$ ) (Supplemental Table 4 (available at <http://www.ajo.com>)). We then produced a multivariable model via backwards selection on this subset of eyes using corneal data, patient data and rs72631889 as an additional covariate as shown in Supplemental Table 5 (available at <http://www.ajo.com>). However rs72631889, although significant ( $P=0.005$ ), had a negligible contribution (0.3%) to the explained variation in the final model.

## **Model Fitting and Covariate Selection (Excluding Genetic Data)**

The results of the univariate time-to-event analysis on the hazards scale using a Royston-Parmar flexible parametric model is shown in Table 2. Genetic data was excluded from this analysis. All variables except smoking status were significant. The explained variation ( $R^2_b$ ) and discrimination (D) were highest for age (17%) and Kmax (15%) with Front K1, Front K2, Back K1, Back K2 and pachymetry each explaining 6-10% of the variation. Notably, gender and ethnicity, although significant in the univariate analysis, did not contribute to explained variation. The hazard ratios for significant covariates indicate that increasing age at presentation, greater pachymetry and flatter (less negative)

posterior keratometry values decrease risk of having CXL, whilst steeper anterior keratometry values and male gender increase the risk of having CXL.

When we fitted a multivariable model the significant covariates were age, Kmax, Front K1, Front K2 and pachymetry (Table 2). When we removed single variables from the model the effect this had on explained variation and discrimination is shown in Supplemental Table 6 (available at <http://www.ajo.com>). Age was the most important covariate (16.7%), with Kmax contributing ~5% of explained variation. K1, K2 and pachymetry had a small effect (<1%) when removed individually. We chose a model without K2 on the basis of parsimony, which was supported by the fact that K1 and K2 were highly correlated ( $R^2=0.91$ ) as shown in Supplemental Figure 2 (available at <http://www.ajo.com>). The final fitted model hazard ratios can be seen on the multivariable column of Table 2. It is notable that an increase in K1 now has a protective effect in the final model. The explained variation and discrimination for the final model were 32.7% and 1.43 respectively.<sup>27</sup> The opposing effect of Kmax and Front K1 can be explained by examining their regression coefficients before converting to hazard ratios; Kmax has a positive coefficient (0.0795) and Front K1 has a negative coefficient (-0.0749). This is logically similar to including the combined covariate (Kmax - Front K1) in the model which can be viewed clinically as a proxy for irregular astigmatism. We also investigated combining K1 and K2 into a single covariate as K2-K1 (standard definition of astigmatism), but the corresponding p value was not significant.

Figure 1 visually depicts the result of applying the final model to the original dataset. As expected, the predicted mean survival curves closely follow the Kaplan-Meier curves. To demonstrate the use of the model in clinical practice, survival curves for three hypothetical patients followed for five years are shown in Figure 2. We have also produced a web application from the model which can be accessed at <http://beta.moorfieldscxl.com>.

### **Keratometric Progression Sensitivity Analysis**

The results of the keratometric progression sensitivity analysis can be found in the Supplementary Material. By examining the Kaplan Meier curves in

Supplemental Figure 3, we can see that the best case time-to-event curve indicates a 40% survival probability at 5 years whilst the worst case curve indicates a 27% survival probability at 5 years. This 13% difference in survival probability at 5 years represents the upper bound of the discrepancy in survival probability within the data. After fitting the Royston-Parmar model, amongst the hazard ratios which overlap (age, Kmax, k2), there was reasonable similarity (Supplemental Tables 8 and 9). Most importantly, the model fitted to the best case had an explained variation of 11% compared to 23% for the worst case indicating a significant difference in model performance depending on the assumptions used for handling eyes which received CXL.

### **Multivariable Model Validation**

When performing validation using internal-external cross validation, Figure 3 shows the ability of our final model to predict keratoconus progression across different geographic regions. We did not identify any significant differences in prognostic factors across regions. The model prediction curves generally follow the Kaplan Maier curves. Notably, region 5 (South West Greater London) and region 7 (other regions) have a worse predictive performance than the other regions, indicating that these regions have different characteristics compared with the remainder of the dataset used for model fitting. This could be due to differing patient characteristics, such as complex cases that required referral to our tertiary referral centre rather than being managed locally. Overall, the prediction becomes less accurate over time, which is expected due to low numbers with follow-up beyond three years. Supplemental Table 7 displays quantitative validation results of the model using internal external validation. The difference column  $D_k - D_{(k)}$  is a measure of predictive ability. Region 7 (other regions outside of Greater London) has the greatest discrepancy in discrimination (-0.26) which indicates that the model fitted when excluding region 7 had greater discriminative ability than when applied to region 7 alone.

### **Discussion**

In this study we have incorporated demographic, keratometric, and genetic data to generate a prognostic model of keratoconus progression to CXL. We have

shown that parameters recorded at the first examination (age, Kmax, Front K1, minimum pachymetry) can produce a time-to-event curve to calculate a personalized risk for keratoconus progression. Although we chose time to CXL rather than keratometric progression as the end point for the time-to-event analysis, we performed a sensitivity analysis using keratometric progression, and found that a CXL model accounts for a much higher proportion of the explained variation (33%) compared to the keratometric model (11% or 23% for best and worst case respectively). The opposing effects of Kmax and Front K1 were unexpected, but similar to including the combined covariate (Kmax - Front K1) in the model; a possible explanation is that the opposing effect is the result of an increase in irregular astigmatism. Of the significant covariates in our model, younger age made the greatest contribution to our model. Thus, one should have a lower threshold for treatment in younger patients.

When applying internal-external cross validation, the survival curves closely followed the Kaplan Meier survival curves for each of the geographic regions, which indicates generalisability, and model discrimination between training and cross validation groups was similar, indicating that the predictive ability is well maintained. Finally, our SNP genetic data had limited additional predictive utility for keratoconus progression. However, the genetic dataset was relatively small (926 patients), and recruitment was based on the presence of keratoconus, as opposed to the severity of keratoconus, or any other index of risk of rapid progression.

The Royston-Parmar model has previously been used to predict the likelihood of the worst eye of patients with keratoconus progressing to corneal transplantation.<sup>28</sup> In their final model, Quartilho et al chose 3 significant covariates: Kmax, age and ethnicity. The reported covariate hazard ratios that overlap with our study (Kmax and age) were different in magnitude but in the same direction. When performing internal validation their model exhibited good predictive ability. They produced time-dependent receiver operating curves using the validation set and found one-year sensitivity and specificity to be 92.8% and 94.6% respectively. Using logistic regression, Kato et al. found that the two strongest factors associated with the requirement for CXL were age and Kmax, which is consistent with our findings.<sup>19</sup> Moreover the team went on to find that

age combined with corneal tomography maps was able to predict progression and need for crosslinking using deep learning.<sup>29</sup>

An ability to generate personalized time-to-event curves that predict progression to CXL (Figure 2) could directly inform clinical decisions that benefit patient care. Firstly, patients may better understand their own risk for progression and feel more confident in choosing their management options. Secondly, for both clinicians and patients, the prediction of progression may contribute to scheduling treatments, including prioritizing patients at high risk of early progression. For example, patients at high risk with a 98% probability of progressing to CXL at 5 years could be offered CXL at the point of first diagnosis without waiting to demonstrate keratometric progression. Medium risk patients may benefit from a period of clinician-led topographic monitoring. For the lowest risk patients, optometry-led monitoring in the community may be sufficient. This risk stratification could be tailored to regions and reflect local needs and resources such as provision of monitoring services in regions with lower risk and greater capacity for CXL in areas with more high risk patients. Finally, when a decision is made to postpone CXL for further monitoring, the time-to-event curve can contribute to decisions on the scheduling of future follow up reviews, with perhaps shorter time periods where the curve is steepest. Recommendations based on this model on clinical practice is yet to be evaluated.

Our study is subject to several limitations inherent to our dataset. First, if patients had CXL at another hospital, this may not be reliably recorded in the source database. This could lead to a very small number of patients being included in the analysis who have already had CXL. Second, ethnicity is a well established risk factor for keratoconus and keratoconus progression,<sup>27,30,31</sup> but ethnicity is now an optional field at patient registration at our institution and this information was unavailable for approximately 50% of our dataset. However, even when we restricted the dataset to those with ethnicity records, it was not found to be a significant covariate. Third, though the cohort used for univariable and multivariable analysis were identical the number of eyes where all covariates were available was lower than for univariable analysis due to

missing data. Finally, when we used multiple imputation to generate a multivariable model, ethnicity was still not found to be significant. In the model fitting process we chose to use a simple backwards selection as opposed to the multivariate fractional polynomial (MFP) method.<sup>32</sup> In our initial investigations, the results of MFP yielded nonlinear functional forms of the covariates and, whilst this method may have slightly increased the predictive power of the prognostic model, the resulting hazard ratios would be very hard to interpret. In addition, we did not examine time dependent effects for the covariates, which may provide a more accurate model fit, and future studies should examine this option. Finally, although no external validation dataset was available, internal external cross validation allowed us to confirm that our model is generalizable across geographical regions.

In conclusion, we have fitted a prognostic model for progression of keratoconus to CXL which generates a time-to-event curve using age, Kmax, Front K1, minimum pachymetry from time of presentation. Incorporation of a relatively small genetic dataset does not improve the explained variation of our model. Personalized modeling of risk may improve patients' understanding of their condition and the need for CXL. Such a model may help better improve patients and aid clinician decision making to CXL to achieve better outcomes and judicious use of healthcare resources.

### **Disclosures**

The authors have no financial disclosures.

Funding/Support: HM is funded by a Moorfields Eye Charity PhD Studentship (GR001147). NP is funded by a Moorfields Eye Charity Career Development Award (R190031A). Moorfields Eye Charity is supported in part by the National Institute for Health Research (NIHR) Biomedical Research Centre based at Moorfields Eye Hospital NHS Foundation Trust and UCL Institute of Ophthalmology. ST, BA and DG acknowledge that a proportion of their financial support is from the Department of Health through the award made by the National Institute for Health Research to Moorfields Eye Hospital NHS Foundation Trust and University College London Institute of Ophthalmology for

a Specialist Biomedical Research Centre for Ophthalmology. The sponsor or funding organization had no role in the design or conduct of this research.

Other acknowledgements: none.

## References

1. Gordon MO, Steger-May K, Szczotka-Flynn L, et al. Baseline factors predictive of incident penetrating keratoplasty in keratoconus. *Am J Ophthalmol.* 2006;142(6):923-930.
2. Wittig-Silva C, Chan E, Islam FMA, Wu T, Whiting M, Snibson GR. A randomized, controlled trial of corneal collagen cross-linking in progressive keratoconus: three-year results. *Ophthalmology.* 2014;121(4):812-821.
3. Caporossi A, Mazzotta C, Baiocchi S, Caporossi T. Long-term results of riboflavin ultraviolet A corneal collagen cross-linking for keratoconus in Italy: the Siena eye cross study. *Am J Ophthalmol.* 2010;149(4):585-593.
4. O'Brart DPS, Chan E, Samaras K, Patel P, Shah SP. A randomised, prospective study to investigate the efficacy of riboflavin/ultraviolet A (370 nm) corneal collagen cross-linkage to halt the progression of keratoconus. *Br J Ophthalmol.* 2011;95(11):1519-1524.
5. Koller T, Mrochen M, Seiler T. Complication and failure rates after corneal crosslinking. *J Cataract Refract Surg.* 2009;35(8):1358-1362.
6. Gore DM, Leucci MT, Koay SY, et al. Accelerated Pulsed High-Fluence Corneal Cross-Linking for Progressive Keratoconus. *Am J Ophthalmol.* 2021;221:9-16.
7. Salmon HA, Chalk D, Stein K, Frost NA. Cost effectiveness of collagen crosslinking for progressive keratoconus in the UK NHS. *Eye.* 2015;29(11):1504-1511.

8. Lindstrom RL, Berdahl JP, Donnenfeld ED, et al. Corneal cross-linking versus conventional management for keratoconus: a lifetime economic model. *J Med Econ*. Published online November 19, 2020:1.
9. Godefrooij DA, Mangen MJJ, Chan E, et al. Cost-Effectiveness Analysis of Corneal Collagen Crosslinking for Progressive Keratoconus. *Ophthalmology*. 2017;124(10):1485-1495.
10. Larkin DFP, Chowdhury K, Burr JM, et al. Effect of Corneal Cross-linking versus Standard Care on Keratoconus Progression in Young Patients: The KERALINK Randomized Controlled Trial. *Ophthalmology*. Published online April 20, 2021. doi:10.1016/j.ophtha.2021.04.019
11. Vinciguerra R, Belin MW, Borgia A, et al. Evaluating keratoconus progression prior to crosslinking: maximum keratometry vs the ABCD grading system. *J Cataract Refract Surg*. 2021;47(1):33-39.
12. Shajari M, Steinwender G, Herrmann K, et al. Evaluation of keratoconus progression. *Br J Ophthalmol*. 2019;103(4):551-557.
13. Ozalp O, Atalay E. Belin ABCD Progression Display Identifies Keratoconus Progression Earlier Than Conventional Metrics. *Am J Ophthalmol*. 2022;236:45-52.
14. Hashemi H, Panahi P, Asgari S, Emamian MH, Mehravaran S, Fotouhi A. Best Indicators for Detecting Keratoconus Progression in Children: A Report From the Shahroud Schoolchildren Eye Cohort Study. *Cornea*. 2022;41(4):450-455.
15. Flynn TH, Sharma DP, Bunce C, Wilkins MR. Differential precision of corneal Pentacam HR measurements in early and advanced keratoconus. *Br J Ophthalmol*. 2016;100(9):1183-1187.
16. Flockerzi E, Häfner L, Xanthopoulou K, et al. Reliability analysis of successive Corneal Visualization Scheimpflug Technology measurements in different keratoconus stages. *Acta Ophthalmol*. 2022;100(1):e83-e90.



17. Kreps EO, Jimenez-Garcia M, Issarti I, Claerhout I, Koppen C, Rozema JJ. Repeatability of the Pentacam HR in Various Grades of Keratoconus. *Am J Ophthalmol.* 2020;219:154-162.
18. Clark TG, Bradburn MJ, Love SB, Altman DG. Survival analysis part IV: further concepts and methods in survival analysis. *Br J Cancer.* 2003;89(5):781-786.
19. Kato N, Negishi K, Sakai C, Tsubota K. Baseline factors predicting the need for corneal crosslinking in patients with keratoconus. *PLoS One.* 2020;15(4):e0231439.
20. Patrick Royston PL. *Flexible Parametric Survival Analysis Using Stata: Beyond the Cox Model.* Stata Press; 2011.
21. Hardcastle AJ, Liskova P, Bykhovskaya Y, et al. A multi-ethnic genome-wide association study implicates collagen matrix integrity and cell differentiation pathways in keratoconus. *Commun Biol.* 2021;4(1):266.
22. Ding X, Guo X. A Survey of SNP Data Analysis. *Big Data Mining and Analytics.* 2018;1(3):173-190.
23. Royston P, Parmar MKB. Flexible parametric proportional-hazards and proportional-odds models for censored survival data, with application to prognostic modelling and estimation of treatment effects. *Stat Med.* 2002;21(15):2175-2197.
24. Baade PD, Royston P, Youl PH, Weinstock MA, Geller A, Aitken JF. Prognostic survival model for people diagnosed with invasive cutaneous melanoma. *BMC Cancer.* 2015;15:27.
25. Royston P, Parmar MKB, Sylvester R. Construction and validation of a prognostic model across several studies, with an application in superficial bladder cancer: CONSTRUCTION AND VALIDATION OF PROGNOSTIC MODEL. *Stat Med.* 2004;23(6):907-926.

26. Royston P, Sauerbrei W. A new measure of prognostic separation in survival data. *Stat Med*. 2004;23(5):723-748.
27. Tuft SJ, Moodaley LC, Gregory WM, Davison CR, Buckley RJ. Prognostic factors for the progression of keratoconus. *Ophthalmology*. 1994;101(3):439-447.
28. Quartilho A, Gore DM, Bunce C, Tuft SJ. Royston–Parmar flexible parametric survival model to predict the probability of keratoconus progression to corneal transplantation. *Eye* . 2020;34(4):657-662.
29. Kato N, Masumoto H, Tanabe M, et al. Predicting Keratoconus Progression and Need for Corneal Crosslinking Using Deep Learning. *J Clin Med Res*. 2021;10(4). doi:10.3390/jcm10040844
30. Pearson AR, Soneji B, Sarvananthan N. Does ethnic origin. *Eye* . 2000;14:625-628.
31. Georgiou T, Funnell CL, Cassels-Brown A, O’Conor R. Influence of ethnic origin on the incidence of keratoconus and associated atopic disease in Asians and white patients. *Eye* . 2004;18(4):379-383.
32. Royston P. *Multivariable Model-Building: A Pragmatic Approach to Regression Analysis Based on Fractional Polynomials for Continuous Variables*. John Wiley; 2008.

## Legend

**Figure 1:** Chart showing how the Royston-Parmar model fits the entire dataset. We split the eyes into 4 risk groups by their prognostic index: <25th centile (low risk), 25-50th centile (medium-low risk), 50-75th centile (medium-high risk), >75th centile (high risk). The number of eyes at risk corresponds to the Kaplan-Meier curves.

**Figure 2:** Time-to-event curves that predict the risk of progression to CXL for three hypothetical patient profiles. The blue line represents a high risk patient who has a 95% probability of progressing to CXL at 5 years. The red line is a medium risk patient who has a 48% probability of progressing to CXL at 5 years. The green line is a low risk patient who has a 14% probability of progressing to CXL at 5 years. The equation used to generate the curves is:  $S(t) = e^{-H(t)}$ , where  $H(t)$  is the cumulative hazard function and is commonly expressed as  $\ln(H(t)) = s(\ln(t)) + x\beta$ , where  $s(\ln(t))$  is a restricted cubic spline function of log time,  $\beta$  is the vector of coefficients and  $x$  is the vector of covariates. For further details of the derivation, we refer the reader to <sup>20</sup>.

Abbreviations: pachy, pachymetry

**Figure 3:** Predicted and observed survival curves for seven postal code regions of Greater London as shown in Supplemental Figure 1 (available at <http://www.ajo.com>) using IECV. We split the eyes into 4 risk groups by their prognostic index: <25th centile (low risk), 25-50th centile (medium-low risk), 50-75th centile (medium-high risk), >75th centile (high risk).

Abbreviations:

AIC: Akaike Information Criterion

Back K1: Flat posterior keratometry in the central 3 mm zone

Back K2: Steep posterior keratometry in the central 3 mm zone

BIC: Bayes Information Criterion

CXL: Corneal Cross-Linking

EKC: Early Keratoconus Clinic

EPR: Electronic Patient Record

Front K1: Flat anterior keratometry in the central 3 mm zone

Front K2: Steep anterior keratometry in the central 3 mm zone

HR: Hazard Ratio

IECV: Internal-external Cross Validation

Kmax: Maximum anterior keratometry

MFP: Multivariate Fractional Polynomial

SNP: Single-nucleotide Polymorphism

Journal Pre-proof

**Table 1:** Summary statistics for the available covariates at the first examination for 9341 eyes recorded at first visit.

Covariate	Type	Mean	SD	N	Missing No. (%)
Front K1 (D)	Numeric	45.31	3.86	8,813	528 (5.7)
Front K2 (D)	Numeric	48.39	4.85	8,839	502 (5.4)
Back K1 (D)	Numeric	-6.53	0.75	7,949	1392 (14.9)
Back K2 (D)	Numeric	-7.23	0.93	8,702	639 (6.8)
Kmax (D)	Numeric	54.14	8.01	8,834	507 (5.4)
Pachymetry (um)	Numeric	462.92	46.15	8,946	395 (4.2)
Age (years)	Numeric	28.28	7.10	9341	0 (0)
Genetic data <sup>a</sup>	Ordinal	N/A	N/A	1141	8020 (85.9)
Self-reported black or asian ethnicity <sup>b</sup>	Categorical (59.9% black or asian)	N/A	N/A	4889	4452 (47.7)
Male gender	Categorical (67% male)	N/A	N/A	9341	0 (0)
Smoker <sup>c</sup>	Categorical (4.5% smoker)	N/A	N/A	9341	0 (0)

Abbreviations: Front K1, flattest anterior keratometry; Front K2, steepest anterior keratometry; Back K1, flattest posterior keratometry; Back K2, steepest posterior keratometry; Kmax: maximum Keratometry; pachymetry: minimum corneal thickness; SD, Standard deviation; N, number of eyes; N/A, not applicable.

<sup>a</sup>Genetic data comprised of 28 SNPs and was encoded in an additive fashion (0,1,2).

<sup>b</sup>1=black or asian, 0=otherwise. <sup>c</sup>0=non-smoker/ex-smoker, 1=current smoker.

**Table 2:** Univariable and final multivariable model for all considered covariables excluding genetic data in the training dataset fitted on the hazards scale with 5 degrees of freedom.

Covariate	Univariable (N=9341)				Multivariable (N=8701)	
	Hazard Ratio [95% CI]	P Value	R <sup>2</sup> <sub>D</sub>	D	Hazard Ratio [95% CI]	P Value
Ethnicity	1.14 [1.02; 1.27]	0.02	0.4%	0.13	N/A	N/A
Smoker <sup>a</sup>	1.07 [0.9; 1.28]	0.46	0.1%	0.05	N/A	N/A
Male Gender	1.11 [1.01; 1.21]	0.02	0.2%	0.10	N/A	N/A
Age at presentation	0.91 [0.9; 0.92]	<0.001	16.7%	0.92	0.9 [0.90; 0.91]	<0.001
Kmax	1.06 [1.05; 1.06]	<0.001	14.9%	0.86	1.08 [1.07; 1.09]	<0.001
Front K1	1.09 [1.08; 1.1]	<0.001	7.0%	0.56	0.93 [0.91; 0.94]	<0.001
Front K2	1.08 [1.07; 1.08]	<0.001	9.8%	0.67	N/A	N/A
Back K1 <sup>c</sup>	0.67 [0.64; 0.71]	<0.001	5.9%	0.51	N/A	N/A
Back K2 <sup>c</sup>	0.7 [0.67; 0.72]	<0.001	8.4%	0.62	N/A	N/A
Pachymetry 10 <sup>b</sup>	0.93 [0.92; 0.94]	<0.001	7.5%	0.58	0.95 [0.93; 0.96]	<0.001

Abbreviations: N, number of eyes; R<sup>2</sup><sub>D</sub>, explained variation; D, Royston and Sauerbrei's D statistic (used as a measure of discrimination); CI, confidence interval; Kmax, maximum keratometry; Front K1, flattest anterior keratometry; Front K2, steepest anterior keratometry; Back K1, flattest posterior keratometry; Back K2, steepest posterior keratometry; pachymetry, minimum corneal thickness; N/A, not applicable due to this variable not being included in the final model.

<sup>a</sup>0=non-smoker/ex-smoker, 1=current smoker.

<sup>b</sup>Minimum pachymetry in steps of 10µm.

<sup>c</sup>Back K1 and Back K2 are negative values such that patients with advanced keratoconus are typically associated with large negative values. A hazard ratio below 1 indicates that as measurements become more positive, the risk of progression decreases.

Figure 1

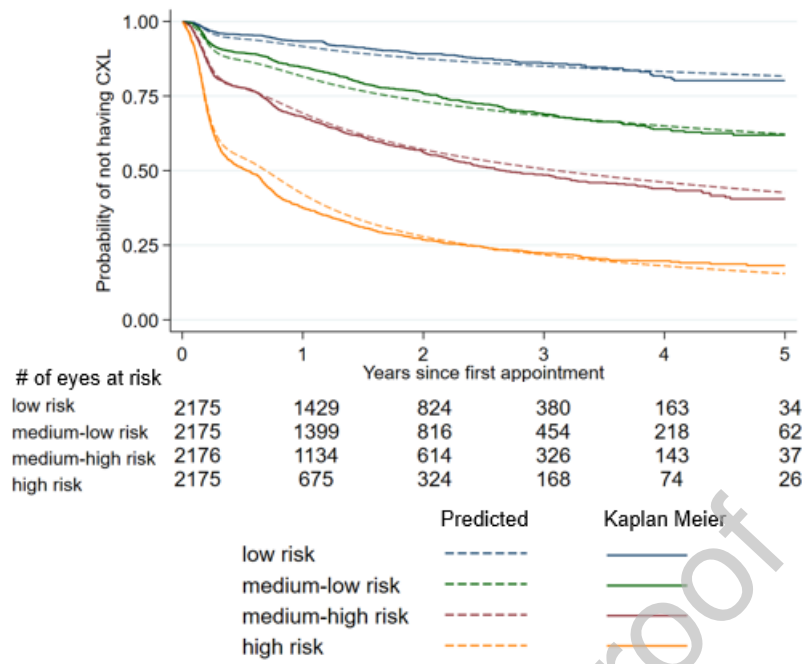


Figure 2

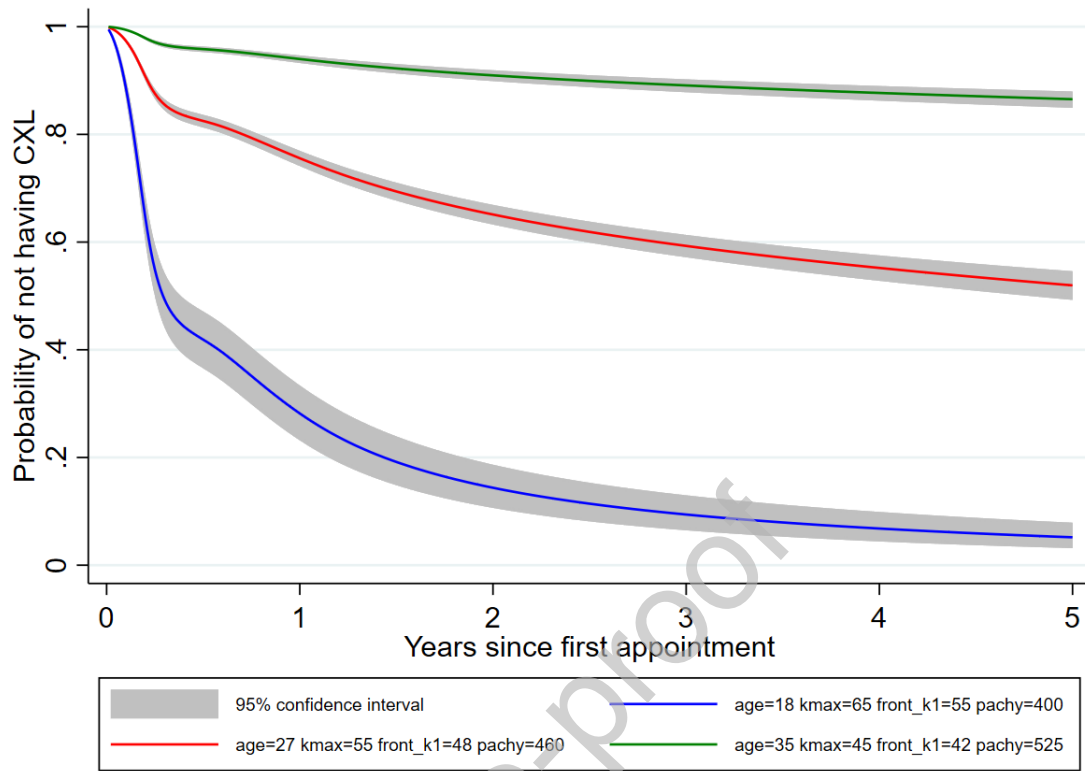




Figure 3

