

# BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email [info.bmjopen@bmj.com](mailto:info.bmjopen@bmj.com)

# BMJ Open

## Cohort profile: AlzEye: longitudinal record-level linkage of ophthalmic imaging and hospital admissions of 353,157 patients in London, United Kingdom

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-058552
Article Type:	Cohort profile
Date Submitted by the Author:	20-Oct-2021
Complete List of Authors:	Wagner, Siegfried; University College London Institute of Ophthalmology, Hughes, Fintan; Duke University Hospital Pontikos, Nikolas; University College London; Moorfields Eye Hospital NHS Foundation Trust Struyven, Robbert; University College London; Moorfields Eye Hospital NHS Foundation Trust Alexander, Daniel; University College London, Department of Computer Science Hindley, Jack; University College London, Department of Information Governance Keane, Pearse; University College London; Moorfields Eye Hospital NHS Foundation Trust Topol, Eric; Scripps Research Institute Cortina Borja, Mario; UCL, Institute of Child Health Liu, Xiaoxuan; University Hospitals Birmingham NHS Foundation Trust; University of Birmingham Montgomery, Hugh; University College London, Centre for Human Health and Performance Petersen, Steffen; Queen Mary University of London, William Harvey Research Institute; Barts Health NHS Trust, Balaskas, Konstantinos; Moorfields Eye Hospital NHS Foundation Trust; Moorfields Eye Hospital NHS Foundation Trust Petzold, Axel; UCL; Moorfields Eye Hospital NHS Foundation Trust Rahi, Jugnoo; UCL, Great Ormond Street Institute of Child Health; Great Ormond Street Hospital for Children NHS Foundation Trust Denniston, Alastair; Queen Elizabeth Hospital Birmingham, UK; University of Birmingham
Keywords:	OPHTHALMOLOGY, Medical ophthalmology < OPHTHALMOLOGY, Medical retina < OPHTHALMOLOGY, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™  
Manuscripts

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

1  
2  
3  
4  
5 Cohort profile: AlzEye: longitudinal record-level linkage of  
6  
7  
8 ophthalmic imaging and hospital admissions of 353,157 patients  
9  
10  
11 in London, United Kingdom  
12  
13  
14  
15

16 SK Wagner<sup>1, 2</sup>, F Hughes<sup>3</sup>, M Cortina-Borja<sup>4</sup>, N Pontikos<sup>1, 2</sup>, RR Struyven<sup>1, 2</sup>, X Liu<sup>5, 6</sup>, H  
17  
18 Montgomery<sup>7</sup>, DC Alexander<sup>8</sup>, E Topol<sup>9</sup>, S Petersen<sup>10, 11</sup>, K Balaskas<sup>1, 2</sup>, J Hindley<sup>12</sup>, A Petzold<sup>1, 2,</sup>  
19  
20  
21 13, 14, J Rahi<sup>1, 4, 15, 16</sup>, AK Denniston<sup>5, 6</sup>, PA Keane<sup>1, 2, \*</sup>  
22  
23  
24  
25

26 \*Corresponding author: Pearse A. Keane, MD FRCOphth, NIHR Biomedical Research Centre for Ophthalmology,  
27  
28 Moorfields Eye Hospital NHS Foundation Trust and UCL Institute of Ophthalmology, 162 City Road, London, United  
29  
30 Kingdom. Tel: +44 207 253 3411 Email: [pearse.keane1@nhs.net](mailto:pearse.keane1@nhs.net)  
31  
32

33 Word count (excluding title page, abstract, references, figures and tables): 4915 words  
34  
35

36 <sup>1</sup> Institute of Ophthalmology, University College London, London, UK,

37 <sup>2</sup> Moorfields Eye Hospital NHS Foundation Trust, London, UK,

38 <sup>3</sup> Duke University Hospital, Durham, North Carolina, USA,

39 <sup>4</sup> Great Ormond Street Institute of Child Health, University College London, London, UK,

40 <sup>5</sup> University of Birmingham, Birmingham, UK,

41 <sup>6</sup> University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK,

42 <sup>7</sup> Centre for Human Health and Performance, University College London, London, UK

43 <sup>8</sup> Centre for Medical Image Computing, Department of Computer Science, University College London, UK

44 <sup>9</sup> Department of Molecular Medicine, Scripps Research, La Jolla, CA, USA.

45 <sup>10</sup> William Harvey Research Institute, Queen Mary University of London, London, UK.

46 <sup>11</sup> St Bartholomew's Hospital, Barts Health National Health Service (NHS) Trust, London, UK..

47 <sup>12</sup> Department of Information Governance, School of Life and Medical Sciences, University College London, London, UK

48 <sup>13</sup> The National Hospital for Neurology and Neurosurgery (UCLH), Queen Square Institute of Neurology, University College London, London, UK.

49 <sup>14</sup> Neuro-ophthalmology Expert Centre, Amsterdam UMC, Departments of Neurology and Ophthalmology, Netherlands.

50 <sup>15</sup> Great Ormond Street Hospital NHS Foundation Trust, London, United Kingdom.

51 <sup>16</sup> Uiverscroft Vision Research Group, University College London, London, UK  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Abstract

Purpose: Retinal signatures of systemic disease ('oculomics') are increasingly being revealed through a combination of high-resolution ophthalmic imaging and sophisticated modelling strategies. Progress is currently limited not mainly by technical issues, but by the lack of large labelled datasets, a sine qua non for deep learning. Such data are derived from prospective epidemiological studies, in which retinal imaging is typically unimodal, cross-sectional, of modest number and relate to cohorts, which are not enriched with subpopulations of interest, such as those with systemic disease. We thus linked longitudinal multimodal retinal imaging from routinely-collected National Health Service data with systemic disease data from hospital admissions using a privacy-by-design third-party linkage approach.

Participants: Between January 1<sup>st</sup> 2008 and April 1<sup>st</sup> 2018, 353,157 participants aged 40 years or older, who attended Moorfields Eye Hospital NHS Foundation Trust, a tertiary ophthalmic institution incorporating a principal central site, four district hubs and five satellite clinics in and around London, United Kingdom serving a catchment population of approximately six million people.

1  
2  
3 Findings to date: Among the 353,157 individuals, 186,651 had a total of 1,337,711 Hospital  
4 Episode Statistics admitted patient care episodes. Systemic diagnoses recorded at these  
5 episodes include 12,022 patients with myocardial infarction, 11,735 with all-cause stroke and  
6 13,363 with all-cause dementia. A total of 6,261,931 retinal images of seven different modalities  
7 and across three manufacturers were acquired from 154,830 patients. The majority of retinal  
8 images were retinal photographs (n=1,874,175) followed by optical coherence tomography  
9 (n=1,567,358).  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

20 Future plans: AlzEye combines the world's largest single institution retinal imaging database  
21 with nationally collected systemic data to create an exceptional large-scale, enriched cohort that  
22 reflects the diversity of the population served. First analyses will address cardiovascular  
23 diseases and dementia, with a view to identifying hidden retinal signatures that may lead to  
24 earlier detection and risk management of these life-threatening conditions.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## STRENGTHS AND LIMITATIONS OF THIS STUDY

- Retinal signatures of systemic diseases, particularly cardiovascular and neurodegenerative diseases, have highlighted the eye as a potential window to risk-stratification.
- AlzEye is a large curated dataset linking more than six million routinely collected retinal images with National Health Service hospital admissions data over a ten-year period in 353,157 patients.
- The cohort is ethnically and socioeconomically diverse and consists of an enriched subpopulation of 12,022 individuals with myocardial infarction and 13,363 individuals with all-cause dementia.
- AlzEye provides a powerful foundation for the development and validation of both traditional and deep learning-based static and dynamic risk prediction models of systemic disease from retinal morphology.



## INTRODUCTION

Scientific discovery has increasingly been driven by the availability of large, diverse, high-dimensional datasets providing deeply phenotyping variables in health and disease[1–3]. Advances in healthcare informatics, hardware and statistical techniques have uncovered relationships previously unachievable through traditional methods of study design. One exemplar is genomics in which the rich and voluminous data afforded by modern sequencing technology has highlighted the contribution of the genome to the trajectory of specific diseases and therapeutic response[4,5]. Similarly, the dynamic field of radiomics seeks to explore meaningful relationships between quantitative data extracted from medical imaging and disease at substantial scale[6].

A crucial substrate for health data research has been the establishment of and accessibility to large prospective epidemiological studies, such as the United Kingdom Biobank (UKBB), the Rotterdam study, and the European Prospective Investigation into Cancer and Nutrition (EPIC) study[7–9]. While such studies undoubtedly represent an exceptionally powerful enabler for discovery science, as evidenced by the high number of publications drawing on them[10], they are potentially limited for investigations of specific subpopulations of interest (e.g. those with rare disease or specific sociodemographic groups) and, where they draw on volunteer participants, also prone to selection bias (e.g. over-representation of more healthy subjects). For example, participants in UKBB are less likely to be obese, smoke, or drink alcohol and accordingly, mortality rates for participants aged 70-74 in UKBB are 46.2% and 55.5% lower for men and women respectively compared to general UK population[11].

Healthcare in England is such that when a patient has a medical event requiring admission, in almost all cases they are admitted under the provisions of the National Health Service (NHS).

1  
2  
3 Following discharge, routinely collected healthcare administrative data during a patient's  
4 admission are subsequently translated into corresponding International Classification of  
5 Diseases (ICD) codes by clinical coders within the respective institution and submitted to the  
6  
7 Secondary Uses Service. This process is overseen by NHS Digital, who provide a unified  
8  
9 record-level national repository of Hospital Episode Statistics (HES) data relating admitted  
10  
11 patient care (APC) for almost all of the English population. While the original purpose of HES  
12  
13 was the monitoring of service activity and negotiation of financial reimbursement, it is  
14  
15 increasingly used for epidemiological research[12]. HES data are amenable to research as a  
16  
17 sole resource. However, using deterministic linkage where individual identifiers are matched in a  
18  
19 rules-based approach in contrast to probabilistic linkage[13], HES can enrich other datasets as  
20  
21 in the case of UKBB[14], the Avon Longitudinal Study of Parents and Children[15] or the  
22  
23 European Prospective Investigation into Cancer in Norfolk[16].  
24  
25  
26  
27  
28  
29

30 While these aforementioned studies demonstrate the value of enriching structured datasets  
31  
32 through linkage with HES, this has not yet been done at scale for routinely collected data of  
33  
34 high-dimensionality, such as imaging. Thus, we established AlzEye, a large linked dataset  
35  
36 which combines routinely collected retinal images and relevant ophthalmic data from an  
37  
38 unselected population attending Moorfields Eye Hospital NHS Foundation Trust (MEH), and  
39  
40 nationally collected systemic healthcare outcome data provided through the HES APC data  
41  
42 record of NHS Digital. MEH is a tertiary ophthalmic institution incorporating a principal central  
43  
44 site, four district hubs and five satellite clinics in London, United Kingdom (UK), providing care to  
45  
46 an ethnically and socioeconomically diverse population of six million people (9% of the UK  
47  
48 population)[17]. The primary aim of AlzEye is to characterise the association between imaging-  
49  
50 derived retinal biomarkers and chronic complex disorders of ageing, particularly dementia and  
51  
52 cardiovascular diseases. In addition to describing the characteristics of the AlzEye cohort and  
53  
54 dataset, in this paper, we describe the key governance, technical and ethical factors that need  
55  
56  
57  
58  
59  
60

1  
2  
3 to be addressed to enable large institution-led individual-level linkage of routinely collected  
4 multi-dimensional data. We also share an approach that has enabled us to create an  
5 exceptional trans-disciplinary resource that can support sophisticated modelling approaches,  
6 such as deep learning, so as to explore the retinal signatures of systemic disease.  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For peer review only

## COHORT DESCRIPTION

### Study population

The AlzEye project is a retrospective cohort study of patients aged 40 years and over who have attended MEH between January 1st 2008 and April 1st 2018. Patients were included if they had attended the glaucoma, retina, neuro-ophthalmology or emergency ophthalmic services and had valid NHS numbers. Those with invalid NHS numbers, dates of birth or who had previously opted out of their health data being used for purposes of research (described in the NHS as a 'Type 2 opt-out') were excluded. Ethnicity group was self-reported by the patient as i) Asian or Asian British, ii) Black or Black British, iii) Mixed, iv) Other Ethnic Group, v) White, or vi) Unknown. Socio-economic status was categorised using the Index of multiple deprivation (IMD) decile, which was estimated by permuting the IMD 2015 rank from the patient's postcode through Lower Super Output Areas followed by aggregation into deciles[18]. Mortality data were derived from the MEH database, which is updated on a two-weekly basis using reports extracted from the NHS National Spine, and is completed on an individual basis by the MEH Data Quality team to ensure accuracy. Data are completed on any patients who have ever attended MEH. Mortality data up to the end of the study period, 1st April 2018, were included.

## Approvals and process

The following key steps in the governance processes were required to establish AlzEye, and provide the necessary ongoing assurance within the research ethics framework of the NHS and the legal framework of the UK. In order to support other researchers who wish to establish similar linked cohorts, we provide an explanation of each stage which outlines the principle which that stage addresses, the UK framework that meets that principle, and finally any study-specific considerations that we undertook to not only meet but exceed those requirements (Figure 1).

- Funding

It was necessary to secure funding to deliver the study, and to provide assurance to the sponsor and others that the study would be completed and that the integrity of the study would not be compromised by inadequate resources. In AlzEye, the study was funded through a small grant awarded by Fight for Sight and Alzheimer's Research UK in October 2017 covering the costs of data storage and linkage fees. The funders had no role in the conception, design or analysis of the study.

- Sponsorship

It was necessary to secure a sponsor for the study that would take on 'overall responsibility for proportionate, effective arrangements being in place to set up, run and report on a research project'. In AlzEye, we sought sponsorship from the relevant NHS Trust (MEH) at which all patients had been seen. Sponsorship confirmation was only sought following internal consultation between data protection, information governance, information security and information technology (IT) teams at both MEH and UCL. MEH acted as the Sponsor, with UCL acting as the trusted third party linking retinal images and HES data and providing

1  
2  
3 computational facilities for data analysis. The study and data governance were approved on  
4  
5 24/05/2018 (internal reference: KEAP1004).  
6  
7

8  
9

- Health Research Authority (HRA) Approval

10  
11 For most research studies in England and Wales, including those limited to working with data for  
12  
13 specific projects, Health Research Authority (HRA) approval is required. HRA approvals involve  
14  
15 the assessment of governance and legal compliance of a research study with an independent  
16  
17 ethics review and opinion from the Research Ethics Committee (REC). Depending on other  
18  
19 study characteristics (e.g. the use of ionising radiation, gene therapy), additional applications  
20  
21 may be required to inform the HRA approval. In England, limited access to confidential patient  
22  
23 information without consent may be granted under the provisions of Section 251 of the National  
24  
25 Health Service Act 2006. This permits temporary lifting of the common law duty of confidentiality  
26  
27 around confidential patient information ‘in the public interest’ or ‘in the interests of improving  
28  
29 patient care’[19]. Obtaining section 251 support requires application to the Confidential Advisory  
30  
31 Group (CAG), an independent body providing expert advice to the HRA for research  
32  
33 applications and NHS Digital for data dissemination. Applications were accordingly made to the  
34  
35 Research Ethics Committee (18/LO/1163, approved 01/08/2018) and the CAG for Section 251  
36  
37 support (18/CAG/0111, approved 13/09/2018). The National Health Service Health Research  
38  
39 Authority gave final approval on 13/09/2018. Approvals thus far granted the legal basis for  
40  
41 submitting an application to the Data Access Request Service (DARS) of NHS Digital[20].  
42  
43  
44  
45

46  
47

- NHS Digital and the Data Access Request Service (DARS)

48  
49 Among its services as the national information and technology partner for the NHS, NHS Digital  
50  
51 oversees the Data Access Request Service (DARS), which administers and provides, upon  
52  
53 application, multiple England-wide datasets from disease-specific audits (e.g. National Diabetes  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Audit Core) to general admissions in secondary care (e.g. Hospital Episode Statistics).

4  
5 Applications to DARS require that the organisation have, at a minimum, the following:

- 6  
7 i) Data Sharing Framework Contract for Data Controllers  
8  
9 ii) Compliance with minimum-security standards for Data Processors and Data  
10  
11 Storage locations  
12  
13 iii) Adequate information security certification (e.g. ISO27001)  
14  
15 iv) A legal basis for data access (e.g. Section 251)  
16  
17  
18  
19

20 Applications are then reviewed with an assigned case officer, who will liaise with the applicant  
21 on project-specific items. For AlzEye, dialogue between the applicant and NHS Digital data  
22 production team revolved around confirmation of data fields and datasets (HES) and the  
23 pseudonymisation embedded within the linkage strategy.  
24  
25  
26  
27  
28  
29

- 30  
31 ● Independent Group Advising on the Release of Data (IGARD)

32 Following internal NHS Digital review and prior to data release, DARS applications are  
33 scrutinised by the Independent Group Advising on the Release of Data (IGARD) in line with  
34 Section 263(2) of the Health and Social Care Act 2012, the Code of Practice on confidential  
35 information. IGARD is an independent panel with a broad range of expertise, from legal to  
36 information governance to epidemiology. Support for AlzEye was given by IGARD in January  
37 and August 2019 citing that 'aspects of the application could be used as an exemplar by NHS  
38 Digital to help other researchers with their applications to the Data Access Request Service  
39 (DARS)'[21].  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49

- 50  
51 ● Data sharing agreements (DSA)

52 Prior to data receipt, data sharing agreements (DSA) must be signed between NHS Digital and  
53 the Data Controller and are overseen by their respective legal departments. AlzEye required an  
54  
55  
56  
57  
58  
59

1  
2  
3 additional DSA between MEH and UCL for the transfer of ophthalmic imaging and clinical data  
4  
5 between institutions outlining the purpose and legal basis for sharing.  
6  
7

- 8  
9 ● Data processing and transfer

10  
11 The dataset was finalised upon completion of engineering work parsing manufacturer-specific  
12  
13 file formats to non-proprietary data structures amenable to image analysis with appropriate  
14  
15 deidentification. A secure cloud-based informatics pipeline was used for transfer of images to  
16  
17 UCL from MEH, the establishment of which was delayed by the COVID-19 pandemic. Imaging  
18  
19 data was stored (with backup) across four dedicated discs on a network-attached storage  
20  
21 device within the UCL Institute of Ophthalmology, School of Life and Medical Sciences (SLMS)  
22  
23 and only accessible to members of the AlzEye research team. All data entities were listed within  
24  
25 the UCL SLMS Information Asset Register.  
26  
27  
28  
29  
30

## 31 **Patient and public involvement and engagement**

32  
33 Patient and public involvement and engagement (PPIE) support was provided by the National  
34  
35 Institute for Health Research (NIHR) Biomedical Research Centre (BRC) at Moorfields Eye  
36  
37 Hospital and UCL Institute of Ophthalmology and has been embedded throughout this project  
38  
39 from priority setting to plans for dissemination. Feedback has been sought through public  
40  
41 engagement events, survey of eye service users and reports within the media. Patients and the  
42  
43 public actively contributed to identifying the priority setting of dementia and the acceptability of  
44  
45 using routinely acquired eye scans for research purposes and without consent. In addition, two  
46  
47 members of the public will sit on the AlzEye working group to contribute to results interpretation  
48  
49 and co-authoring and dissemination of research outputs. The members will be supported in  
50  
51 selecting the results they find relevant and presenting them to wider patient communities.  
52  
53  
54  
55  
56  
57  
58  
59  
60



## Ophthalmic health variables

Patient-level ophthalmic variables were extracted from the MEH data warehouse, which aggregates information from the patient administration system (PAS), electronic health record and imaging database, all linked through a unique MEH hospital identification number. Sociodemographic data, including date of birth, sex, ethnicity and postcode as well as patient's clinic appointments and operation dates are housed within PAS. Surgical procedures were recorded in the electronic health record (EHR) at MEH from September 4th, 2012. Operation details, including procedure name, laterality and indication for surgery are contained within the MEH EHR and uploaded to the MEH data warehouse. A patient undergoing the most common operation in the United Kingdom, cataract extraction, would therefore have an entry for the typical procedure, (*phacoemulsification and intraocular lens implant*) operated eye (*right or left*) and indication (*cataract*).

Colour retinal photography (Figure 2A) and optical coherence tomography (OCT, Figure 2B) images, which represent the majority of retinal images within the database, have been processed through segmentation and feature extraction software. The Vascular Assessment and Measurement Platform for Images of the Retina (VAMPIRE) system provides fully automated segmentation and extraction of retinal fractal dimensions, a measure of the geometric complexity of retinal vasculature, from colour retinal photographs with future plans to derive retinal vascular calibre[22,23]. OCT scans are segmented and retinal sublayer thicknesses computed using the Topcon Advanced Biomedical Imaging Laboratory (TABIL) software[24]. Retinal sublayer analyses will conform to the consensus nomenclature outlined by the International Nomenclature for Optical Coherence Tomography Panel (Figure 2E)[25].

1  
2  
3 For the purposes of this report, four common ophthalmic diseases were described - cataract,  
4 glaucoma, neovascular age-related macular degeneration (AMD) and proliferative diabetic  
5 retinopathy (PDR).  
6  
7

8  
9 *Cataract* was defined as any operation code denoting phacoemulsification surgery and the  
10 indication of cataract. For the purposes of this report, only first eye cataract surgery was  
11 included.  
12  
13  
14

15  
16 *Glaucoma* was defined as any patient attending the glaucoma clinic three or more times with  
17 ongoing follow up from January 1st 2010. The first two years of the study period were excluded  
18 as this may have incorporated patients with previous diagnoses of glaucoma where the  
19 maximum follow up interval can approach two years; in contrast any patient being seen after 2  
20 years since study inception with no previous visit within that 2 year period can be assumed to  
21 have/carry a new diagnosis of glaucoma.  
22  
23  
24  
25  
26  
27

28  
29 *Diabetic eye disease* represents a special case due to audit procedures mandated by the NHS  
30 Diabetic Eye Screening Programme. Coding of eye disease secondary to diabetes mellitus is  
31 rigorously validated by a dedicated team within MEH with a dedicated database consisting of all  
32 patients with diabetic eye disease and their retinopathy and maculopathy grades according to  
33 the NHS Diabetic Eye Screening Programme criteria[26], at hospital appointment from  
34 September 12, 2013 onwards. Dates for onset of proliferative diabetic retinopathy dates were  
35 recorded as the first appointment for each patient where this diagnosis was first made.  
36  
37  
38  
39  
40  
41  
42

43  
44 *AMD* can be categorised into two major types - dry and neovascular ("wet"). Dry AMD is slowly  
45 progressive and with no active hospital intervention currently available; it is thus MEH standard  
46 practice for patients with features typical of dry AMD to be discharged with lifestyle and  
47 monitoring advice (self-monitoring and standard optometric review). In contrast, neovascular  
48 AMD requires treatment through intravitreal anti-VEGF injections, and therefore remains under  
49 active follow-up. The diagnostic codes for neovascular AMD were based on extensive previous  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 work in which all patients with neovascular AMD at MEH were manually validated up to  
4  
5 2018[27,28].  
6  
7  
8  
9

## 10 **Systemic health variables**

13 Systemic health data were derived from HES APC data, with a focus on cardiovascular disease  
14 and all-cause dementia. Diagnostic codes in HES APC are reported in line with the 10th revision  
15 of ICD[29]. In line with previous reports, myocardial infarction was defined as code I21 or  
16 I22[30–32]. Similarly, stroke was defined using stroke definitions from UK Biobank[33].  
17  
18 Dementia was defined as ICD codes E512 (Wernicke’s encephalopathy), F00 (Dementia in  
19 Alzheimer disease), F01 (Vascular dementia), F02 (Dementia in other diseases classified  
20 elsewhere), F03 (Unspecified dementia), F10.6 (Mental and behavioural disorders due to  
21 psychoactive substance use, Amnesic syndrome) F10.7 (Mental and behavioural disorders due  
22 to psychoactive substance use, Residual and late-onset psychotic disorder) , G30 (Alzheimer  
23 disease), or G31.0 (Other degenerative diseases of nervous system, not elsewhere classified),  
24 derived from previous work evaluating the agreement between HES admitted patient care data  
25 and primary care data, through general practitioner surveys and the Clinical Practice Research  
26 Datalink (CPRD)[34]. The unique number of patients specifically with Alzheimer’s disease (F00,  
27 G30) and vascular dementia (F01) are also described.  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46

## 47 **Data Linkage and Transfer**

48  
49 The linkage strategy was designed through collaboration between experts in information  
50 governance, information technology, computer scientists and clinicians based at MEH,  
51 University College London (UCL) and NHS Digital (Figure 3). A third-party linkage approach  
52 was used for two main reasons. Firstly, it enhanced privacy-preservation as the data originator,  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 MEH, never received HES admissions data and the third party, UCL, did not receive personally  
4 identifiable information. Secondly, it enabled the linked dataset to be accessible within a site  
5 with sufficient high-performance computing capability to undertake the proposed analyses, a  
6 function significantly beyond almost all NHS facilities. Patient link identifiers consisting of a  
7 unique NHS identification number, sex and date of birth originating from MEH were transferred  
8 to NHS Digital in conjunction with a unique study ID generated using a cryptographic hash  
9 function (random pseudonimisation). Ophthalmic covariates, mortality data, and patient  
10 sociodemographics with study ID were transferred to UCL. Ophthalmic imaging data pertaining  
11 to the patients within the study were extracted and de-identified during conversion from their  
12 proprietary format to Digital Imaging and Communications in Medicine (DICOM) format before  
13 transfer to UCL. Following linkage with HES, NHS Digital transferred HES data to the UCL Data  
14 Safe Haven, a “walled garden” trusted research environment certified and externally audited to  
15 ISO27001 information security standards[35,36].  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

### 33 **Statistical analysis**

34  
35  
36 Imaging-based studies within the AlzEye study are generally planned to take the form of nested  
37 case-control studies. To improve efficiency, controls may be matched with cases by age and  
38 sex, using conditional logistic regression for statistical modelling of binary outcomes and  
39 survival analysis for time-to-event data (e.g. Cox Proportional hazard modelling)[37]. In cases  
40 where the competing risk of death is prominent, subdistribution hazard ratios with 95%  
41 confidence intervals will be estimated as a sensitivity analysis using the multivariable  
42 semiparametric competing risks proportional hazards model approach described by Fine and  
43 Gray[38]. Alternative high-dimensional modelling approaches, such as convolutional neural  
44 networks and vision transformers, will also be explored. Prior to receipt of HES data from NHS  
45 Digital, sample size calculations were undertaken. Specifically, we evaluated the association  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 between OCT-derived peripapillary retinal nerve fibre layer and macular ganglion cell-inner  
4 plexiform layer thicknesses and dementia. Given an odds ratio of 1.4 with an alpha of 5% and a  
5 power of 90% on a 1:1 matched study design, a total sample size of 2106 is required[39].  
6  
7  
8  
9

10  
11 Figures for this report were designed in R version 4.1.0 (R Core Team, 2021. R Foundation for  
12 Statistical Computing, Vienna, Austria).  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## FINDINGS TO DATE

Extraction of unique patients attending MEH outpatient clinics between January 1st 2008 and April 1st 2018 generated a cohort of 353,191 unique patients. Thirty-four patients with invalid or missing linking variables were excluded leaving a total of 353,157. A breakdown of sociodemographic details by category of the cohort are provided in Table 1. Of the cohort, 190,494 were female (53.9%) and the mean age was 68.4 (+/- 13.9) years. Among those who self-reported ethnicity (n=271,293, 76.8%), 48,119 were South Asian, 31,614 Black, 2966 Mixed, 52,851 Other ethnic group and 135,743 White. Of the 353,157 patients, 186,651 had a total of 1,337,711 HES episodes in the study period. NHS Digital performs a hierarchical stepwise linkage approach providing a 'Match Rank' for each HES episode[40]. Among the 1,337,711 HES episodes matched, Match Rank was two for 1,337,482 episodes (exact NHS number, exact date of birth and exact sex linked), four for 46 episodes (exact NHS number, exact sex and partial date of birth), and eight for 183 episodes (exact NHS number).

An illustration of the major common ophthalmic diseases within the cohort is shown in a CONSORT-style diagram in Figure 4. Following the case definition, and exclusion of invalid dates, a total of 59,102 patients had first eye cataract surgery, 31,060 glaucoma, 7214 neovascular AMD and 2494 PDR.

	Characteristic	N (%)
	All	353,157
Sex	Female	190,494 (53.9)
	Male	162,663 (46.1)
Age group*	40-49 years	35,262 (10.0)
	50-59 years	66,101 (18.7)
	60-69 years	79,018 (22.4)
	70-79 years	84,942 (24.1)
	80+ years	87,834 (24.9)
Ethnicity	Black	31,614 (9.0)
	White	135,743 (38.4)
	South Asian	48,119 (13.6)
	Other/Unknown	137,681 (39.0)
Index of multiple deprivation decile	1 (most deprived)	18,194 (5.2)
	2	50,443 (14.3)
	3	50,869 (14.4)
	4	42,603 (12.1)
	5	38,964 (11.0)
	6	36,906 (10.5)
	7	31,317 (8.9)
	8	28,180 (8.0)
	9	29,906 (8.5)
	10 (least deprived)	24,610 (7.0)
	Unknown	1165 (0.3)

Table 1: Baseline sociodemographic characteristics of the AlzEye cohort. Data are shown as n(%).

\*Age is taken as that of April 1st, 2018.

Among the 187,811 patients with recorded HES episodes, 12,022 patients had episodes with coded myocardial infarction, 11,735 patients with all-cause stroke and 13,363 with dementia. Within the dementia group, 4487 patients had codes that were specific for Alzheimer's dementia and 3381 for vascular dementia (Table 2).

Group	Disease	ICD code(s)	Number of patients
Cardiovascular	Acute coronary syndrome	I21, I22	12,022
	Heart failure	I50	24,034
	Atrial Fibrillation	I48	32,848
	Hypertension	I10, I15	151,937
	Subarachnoid haemorrhage	I60	642
	Intracerebral hemorrhage	I61	1865
	Ischaemic stroke	I63-I64	9996
	All stroke	I60, I61, I63, I64	11,735
Neurodegenerative	Alzheimer's disease	F00, G30	4487
	Vascular dementia	F01	3381
	Parkinson's disease	G20	3211
	All-cause dementia	E12, F00, F01, F02, F03, F106, F107, G30, G310	13,363
Other	Diabetes mellitus (Type 1 and 2)	E10, E11	71,570

Table 2: Number of patients by selected examples of specified 10th revision of International Classification of Diseases (ICD) codes relating to diabetes mellitus, cardiovascular and neurodegenerative diseases.

## Imaging

During the study period, a total of 6,261,931 images were acquired from 154,830 patients. The three leading image modalities were colour retinal photographs (n=1,874,175), OCT (n=1,567,358) and red-free photographs (n=1,147,487). The distribution of imaging modalities across the three vendors used for retinal imaging at MEH - Topcon (Topcon corp, Tokyo, Japan), Heidelberg (Heidelberg Engineering, Heidelberg, Germany), and Optos (Optos,



Dunfermline, UK) - are shown in Figure 5 and Table 2. Most images were acquired on the Topcon system (n=5,553,826, 88.7%). Number of images by year is shown in Figure 6. During the study period, annual imaging acquisition increased from 229,868 scans in 2008 to 1,021,904 in 2017. For 2018, collection stopped on April 1st precluding a complete annual figure.

Vendor	Modality	Number of images	Number of patients
<b>Topcon</b>	Angiography	1,128,723	21,225
	Autofluorescence	11,761	2078
	Colour photography	1,874,175	139,307
	Red-free	1,146,854	122,453
	OCT	1,391,826	138,911
	Other	487	48
<b>Heidelberg</b>	Angiography	89,264	4061
	Autofluorescence	94,533	16,863
	Infrared	192,634	21,676
	OCT	175,532	21,191
	Other	19,781	2439
<b>Optos</b>	Angiography	77,813	2215
	Autofluorescence	18,590	5666
	Pseudocolour photography	39,958	6887

Table 3: Retinal imaging within the AlzEye dataset by vendor and imaging modality. OCT: Optical coherence tomography. Angiography refers to dye-based techniques (fluorescein and indocyanine green).

## STRENGTHS AND LIMITATIONS

To our knowledge, we have created the world's largest retinal imaging research dataset available presently, linking secondary healthcare ophthalmic data from 353,157 patients seen over a ten-year period with information on general health and key systemic diseases, as captured through admissions to any hospital within the NHS of England. This comprises 6,261,931 images, obtained using seven different modalities from three different manufacturers, in 154,830 patients. The current large-scale UK health database, UKBB, provides useful context for AlzEye. Cross-sectional data are available in UKBB with two retinal imaging modalities (colour retinal photography and OCT) obtained using technology from one manufacturer (Topcon), and at a single time point in 67,321 people. Notwithstanding the recognised limitations (see subsection 'Limitations of the AlzEye cohort') of real-world datasets and the coding within the HES database, AlzEye provides a number of distinct advantages beyond purely scale. Retinal imaging data are longitudinal, highly multimodal, and pertain to an ethnically and socioeconomically diverse cohort representative of the adult population with eye disease. Moreover, AlzEye has demonstrated relatively low cost. The study is funded through a charity small grant award and NIHR Biomedical Research Centre support amounting to £20,000.

### **Comparison with other resources**

UKBB is the major comparator for AlzEye, being the largest of the prospective epidemiological cohort datasets which provide cross-sectional retinal imaging in association with systemic disease variables[41]. One of the limitations of UKBB is that, unlike AlzEye, it does not provide longitudinal retinal images. Another prospective epidemiological cohort study, the Rotterdam study, does collect longitudinal retinal imaging data from approximately 15,000 participants, of which 5065 participants were eligible for OCT scanning as of September 2017[42]. The

1  
2  
3 Rotterdam Study has uncovered several landmark findings, particularly in regards to causal  
4 determinants, but its cohort remains relatively small in comparison to UKBB and AlzEye with the  
5 majority of participants recruited from one district within Rotterdam, Netherlands[7,10]. The  
6  
7 Singapore Epidemiology of Eye Disease (SEED), is one longitudinal multimodal retinal imaging  
8 initiative which is underway, in which 10,033 participants of Chinese, Indian and Malay ethnicity  
9 have been recruited to undergo six-yearly retinal imaging[43]. Whilst the importance of such  
10 initiatives is often highlighted, this does not appear to have yet translated into actual research-  
11 ready resources. A recent review of ophthalmic imaging datasets did not reveal any additional  
12 relevant publicly available datasets that included linked systemic health data[44]. Additionally,  
13 our own review of the literature has not identified any examples of large-scale linked real-world  
14 datasets (i.e. including those with restricted access) which include linked systemic health data. The  
15 scarcity of such resources suggests that the construction of such datasets is challenging to  
16 undertake, presumably due to factors such as cost, required duration and delayed output, retention  
17 of participants, and concerns over technological redundancy. The AlzEye approach is an important  
18 alternative model in this context.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35

### 36 **Potential research impact from the novel AlzEye cohort**

37  
38 Several epidemiological opportunities arise with AlzEye. Firstly, it provides a real-world  
39 snapshot of ophthalmic secondary care use, representing approximately 1.2% of the UK  
40 population aged 40 and above (27,858,459 in 2011). This is a powerful tool for informing public  
41 health and policymaking in eye services[45], and is exceptional in characterising the potential  
42 impact that may arise from the intersect between disabling diseases such as stroke and  
43 proliferative diabetic retinopathy.  
44  
45  
46  
47  
48  
49  
50  
51  
52

53 Secondly, it allows the identification and exploration of relationships between newly diagnosed  
54 ophthalmic disease (or newly referred to hospital eye services) and emerging systemic events  
55  
56  
57  
58  
59  
60

1  
2  
3 and accruing multimorbidity. Patients tend to respond early to issues with their sight and an  
4 understanding of how an ophthalmic presentation is linked to an increased likelihood of serious  
5 systemic disease may provide an opportunity for earlier intervention in those diseases[46].  
6  
7  
8  
9

10  
11 Thirdly, nested case-control studies evaluating retinal-based oculomic biomarkers in those with  
12 systemic diseases (e.g. dementia) can provide insight into their value in either static or dynamic  
13 risk prediction. Newer modelling approaches have highlighted the potential utility of the retina in  
14 screening for and risk stratification of cardiovascular, neurodegenerative, renal, hepatic and  
15 haematological diseases[47–52];[53].  
16  
17  
18  
19  
20  
21  
22  
23

24 Finally, by its magnitude and wealth of high-quality labels, both ophthalmic and systemic,  
25 AlzEye provides a powerful catalyst for high-dimensional model development, echoing that of  
26 Imagenet, a database currently exceeding 14 million images, which propelled deep learning and  
27 computer vision research forward a decade ago[54].  
28  
29  
30  
31  
32  
33  
34

### 35 **Lessons learned from the AlzEye approach**

36 AlzEye highlights an opportunity for maximising the value of routinely-collected data to support  
37 research for patient benefit. However, there are a number of governance and technical  
38 challenges when undertaking large scale investigator-led data linkage[55]. In AlzEye, early  
39 cross-disciplinary dialogue between experts in information governance, information technology  
40 and data protection at both institutional parties (MEH and UCL) as well as the data production  
41 team at NHS Digital established a privacy-by-design linkage approach, which enhanced privacy  
42 preservation while maintaining computational feasibility[56][57]. At its worst, an intrusion of the  
43 identifiable data during the development of AlzEye would have informed the violator that a given  
44 individual had visited MEH at some point between 2008 and 2018. Due to the novel approach of  
45 AlzEye within our centre, the greatest governance hurdle was securing study sponsorship, a  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 process which took nearly eight months. Once approved, permissions from the bodies of the  
4 Health Research Authority, and overall approval were given within eight weeks. Linking with  
5 high-dimensional imaging data also posed a number of technical obstacles. As highlighted  
6 recently by the American Academy of Ophthalmology (AAO), ophthalmic imaging technologies  
7 suffer from limited interoperability and low compliance to standardised formats, such as DICOM,  
8 the AAO also commenting that 'even so-called DICOM-compliant devices fail to meet DICOM  
9 standards with significant limitations, such as the embedding of patient identifiers on the  
10 image'[58]. A key undertaking within AlzEye was thus the secure and robust but efficient fully-  
11 automated processing of raw ophthalmic imaging data from its proprietary file format with  
12 associated personal metadata to standard DICOM form with the identifiers stripped. Fortunately,  
13 while this operation requires significant technical and engineering input, most medical imaging  
14 modalities already benefit from standardisation among vendors obviating this step for other  
15 researchers seeking to emulate our approach. Finally, a key objective of AlzEye is the  
16 development of clinical prediction models using deep learning approaches, which require  
17 significant computing capacity. Provisions for four graphics processing units (GPU) housed  
18 within UCL enable this step however other centres may additionally consider recent guidance  
19 on the safeguards required for locating health data within cloud environments and the  
20 implications this brings for accessing virtual GPUs[59].  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

### 43 **Limitations of the AlzEye cohort**

44  
45 Despite the opportunities afforded by AlzEye, there are a number of limitations to this kind of  
46 approach and potential sources of bias. Firstly, caution must be paid to the validity of HES  
47 diagnostic coding[60]. Although previous validation studies have concluded that discharge  
48 coding within HES is sufficiently robust for research purposes[34,61], sizable proportions of  
49 cases may be missed when using individual sources[62]. For example, recent work linking the  
50 electronic health records of 54.4 million people in England showed that HES captured 80.5%  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 and 65% of myocardial infarctions and stroke/transient ischaemic attacks respectively when  
4 compared to linkage additionally incorporating Death Registry and primary care records[63].  
5  
6 One mitigation strategy for this source of bias for real world data is therefore linking to multiple  
7 sources. In terms of selection bias, as a hospital-attending cohort, the individuals within the  
8 AlzEye cohort are likely to have greater medical comorbidity than the general population,  
9 limiting the external validity of any findings. In addition, by the very nature of the dataset,  
10 patients within the AlzEye cohort will have definite or suspected ophthalmic disease, particularly  
11 among those with repeated retinal imaging. The risk of under-recording of potentially important  
12 variables such as smoking may also lead to residual confounding.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24 The enrichment of multimodal health data acquired as part of a patient's routine clinical care  
25 with nationally-held databases provides a powerful foundation for discovery science and  
26 epidemiological research. We highlight key considerations and challenges for those seeking to  
27 link high-dimensional data sources, from high-resolution imaging to waveform data, with locally-  
28 held specialist data. Additionally, we provide the cohort profile for AlzEye, a powerful platform  
29 for oculomic discovery, specifically evaluating the association between retinal morphology, and  
30 both cardiovascular diseases and dementia. Beyond discovery, the ALzEye cohort is  
31 anticipated to become an important resource for the development and validation of deep  
32 learning-based clinical prediction models that may enable earlier intervention for patients at risk  
33 of these life-threatening conditions.  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## LEGENDS

Figure 1: Schematic of the key milestones, prerequisites and approvals with their corresponding achievement dates for the AlzEye dataset. REC: Research Ethics Committee, CAG: Confidential Advisory Group, HRA: Health Research Authority, NHS: National Health Service, IGARD: Independent Group Advising on the Release of Data, DSA: Data Sharing Agreement.

Figure 2: Composite figure showing the major retinal imaging modalities within AlzEye. A: Colour fundus photograph, B: Red-free photograph, C: Fundus autofluorescence (widefield), D: Pseudocolour photography (widefield) and E: Optical coherence tomography of the central macula illustrating segmentation of the individual sublayers. Consensus nomenclature for the retinal sublayers is indicated. NFL= nerve fibre layer, GCL = ganglion cell layer, IPL: inner plexiform layer, INL: inner nuclear layer, OPL: outer plexiform layer, ONL = outer nuclear layer, MZ = myoid zone, RPE-IZ = RPE-interdigitation zone, BM = Bruch's membrane.

Figure 3: Linkage approach of AlzEye. Moorfields Eye Hospital NHS Foundation Trust securely transfers a spreadsheet of identifiers with a study ID to NHS Digital and separately transfers the study ID with ophthalmic data, including diagnoses and retinal images, to University College London. NHS Digital links the identifiers with the Hospital Episode Statistics (HES) database and returns the admissions data with the study ID (and no identifiable data) to UCL. UCL links the ophthalmic data from Moorfields Eye Hospital with HES data from NHS Digital using the study ID.

Figure 4: CONSORT style flow chart illustrating the distribution of cataract, glaucoma, neovascular age-related macular degeneration (AMD) and proliferative diabetic retinopathy (PDR) within the AlzEye dataset.

Figure 5: Parallel sets diagram illustrating the imaging modality across vendors within AlzEye. The majority of images were acquired on the Topcon system and the most frequent modalities were colour photography and optical coherence tomography. Designed using the `networkD3` package.

Figure 6: Stacked bar chart of the annual number of images acquired during the study period for the three leading device vendors at Moorfields Eye Hospital. Data for 2018 represents 3 months only prior to the study end-date.

## ACKNOWLEDGEMENTS

We thank Karen Bonstein and Andi Skilton for support with patient and public involvement and engagement, Menachem Katz, Ben Ward, Ross Green, Maxim Daniline and Simon St John-Green for information technology support, Anthony Peacock for advice on use of Trusted Research Environments and Antonio de la Plaza Larrea and Richard Macmillan for legal guidance on data sharing agreements. We also are grateful to Anthony Khawaja and Cathie Sudlow for feedback on study design.

## CONTRIBUTORS

SW, AD and PK wrote the first draft of the manuscript, which was critically revised by FH, MCB, RS, NP, XL, HM, DA, ET, SP, KB, JH, AP, and JR. SW, FH, NP, HM, JH, AP, JR, AD and PK were involved in the original design of the study. Computer science expertise was through RS, NP and DA, information governance through JH. MCB, AP, JR and AD provided statistical and epidemiological guidance. All authors have approved the final version of this manuscript.

## FUNDING

The study was funded through a small grant awarded by Fight for Sight and Alzheimer's Research UK. Infrastructural support was through the NIHR Biomedical Research Centre at Moorfields Eye Hospital and the UCL Institute of Ophthalmology. The funders had no role in the conception, design or analysis of the study.



## COLLABORATION

National and international collaborations are welcomed though restrictions on access to the cohort mean that only the AlzEye researchers can directly analyse individual-level systemic health data. Interested researchers should contact [s.wagner@ucl.ac.uk](mailto:s.wagner@ucl.ac.uk).

## DATA AVAILABILITY STATEMENT

The data are subject to the contractual restrictions of the DSA between NHS Digital, Moorfields Eye Hospital and University College London and are therefore not available for access beyond the AlzEye research team.

## COMPETING INTERESTS STATEMENT

SW is funded through a Medical Research Council Clinical Research Training Fellowship (MR/TR000953/1).

FH: None.

MCB: None.

NP is funded by a Moorfields Eye Charity Career Development Award (R190031A).

RS: None.

XL: None.

HM is supported by the National Institute of Health Research's Comprehensive Biomedical Research Centre (BRC) at University College London Hospitals (UCLH).

DA: None

ET: None

1  
2  
3 SP receives support from the National Institute for Health Research Biomedical Research  
4  
5 Centre at Barts.

6  
7 KB has received speaker fees from Novartis, Bayer, Alimera, Allergan, Roche, and Heidelberg;  
8  
9 meeting or travel fees from Novartis and Bayer; compensation for being on an advisory board  
10  
11 from Novartis and Bayer; consulting fees from Novartis and Roche; and research support from  
12  
13 Apellis, Novartis, and Bayer.

14  
15 JH: None

16  
17 AP receives financial support from the National Institute for Health Research Biomedical  
18  
19 Research Centre based at Moorfields Eye Hospital NHS Foundation Trust and UCL Institute of  
20  
21 Ophthalmology. AP is part of the steering committee of the ANGI network which is sponsored by  
22  
23 ZEISS, steering committee of the OCTiMS study which is sponsored by Novartis, and reports  
24  
25 speaker fees from Heidelberg Engineering.

26  
27 JR receives support from the National Institute for Health Research as a Senior Investigator and  
28  
29 via the National Institute for Health Research Biomedical Research Centres at Moorfields Eye  
30  
31 Hospital and Great Ormond Street Hospital.

32  
33 AD is Director of INSIGHT, the HDRUK Health Data Research Hub for Eye Health.

34  
35 PK is supported by a Moorfields Eye Charity Career Development Award (R190028A) and a UK  
36  
37 Research & Innovation Future Leaders Fellowship (MR/T019050/1); receives research support  
38  
39 from Apellis; is a consultant for DeepMind, Roche, Novartis, Apellis, and BitFount; is an equity  
40  
41 owner in Big Picture Medical; and has received speaker fees from Heidelberg Engineering,  
42  
43 Topcon, Allergan, Roche, and Bayer; meeting or travel fees from Novartis and Bayer; and  
44  
45 compensation for being on an advisory board from Novartis and Bayer  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## REFERENCES

- 1 Munevar S. Unlocking Big Data for better health. *Nat Biotechnol* 2017;**35**:684–6.
- 2 Shilo S, Rossman H, Segal E. Axes of a revolution: challenges and promises of big data in healthcare. *Nat Med* 2020;**26**:29–38.
- 3 Obermeyer Z, Emanuel EJ. Predicting the Future - Big Data, Machine Learning, and Clinical Medicine. *N Engl J Med* 2016;**375**:1216–9.
- 4 Bhardwaj R, Sethi A, Nambiar R. Big data in genomics: An overview. In: *2014 IEEE International Conference on Big Data (Big Data)*. IEEE 2014. doi:10.1109/bigdata.2014.7004392
- 5 He KY, Ge D, He MM. Big Data Analytics for Genomic Medicine. *Int J Mol Sci* 2017;**18**. doi:10.3390/ijms18020412
- 6 Gillies RJ, Kinahan PE, Hricak H. Radiomics: Images Are More than Pictures, They Are Data. *Radiology* 2016;**278**:563–77.
- 7 Hofman A, Breteler MMB, van Duijn CM, *et al*. The Rotterdam Study: objectives and design update. *Eur J Epidemiol* 2007;**22**:819–29.
- 8 Riboli E, Hunt KJ, Slimani N, *et al*. European Prospective Investigation into Cancer and Nutrition (EPIC): study populations and data collection. *Public Health Nutr* 2002;**5**:1113–24.
- 9 Bycroft C, Freeman C, Petkova D, *et al*. The UK Biobank resource with deep phenotyping and genomic data. *Nature* 2018;**562**:203–9.
- 10 Ikram MA, Brusselle GGO, Murad SD, *et al*. The Rotterdam Study: 2018 update on objectives, design and main results. *Eur J Epidemiol* 2017;**32**:807–50.
- 11 Fry A, Littlejohns TJ, Sudlow C, *et al*. Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank Participants With Those of the General Population. *Am J Epidemiol* 2017;**186**:1026–34.
- 12 Chaudhry Z, Mannan F, Gibson-White A, *et al*. Research Outputs of England's Hospital Episode Statistics (HES) Database: Bibliometric Analysis. *Journal of Innovation in Health Informatics*. 2017;**24**:329. doi:10.14236/jhi.v24i4.949
- 13 Zhu Y, Matsuyama Y, Ohashi Y, *et al*. When to conduct probabilistic linkage vs. deterministic linkage? A simulation study. *J Biomed Inform* 2015;**56**:80–6.
- 14 Hospital information boosts UK Biobank resource. <https://www.ukbiobank.ac.uk/2013/09/20000-participants-return-for-a-repeat-assessment/> (accessed 13 Oct 2020).
- 15 Boyd A, Golding J, Macleod J, *et al*. Cohort Profile: The 'Children of the 90s'—the index offspring of the Avon Longitudinal Study of Parents and Children. *International Journal of Epidemiology*. 2013;**42**:111–27. doi:10.1093/ije/dys064

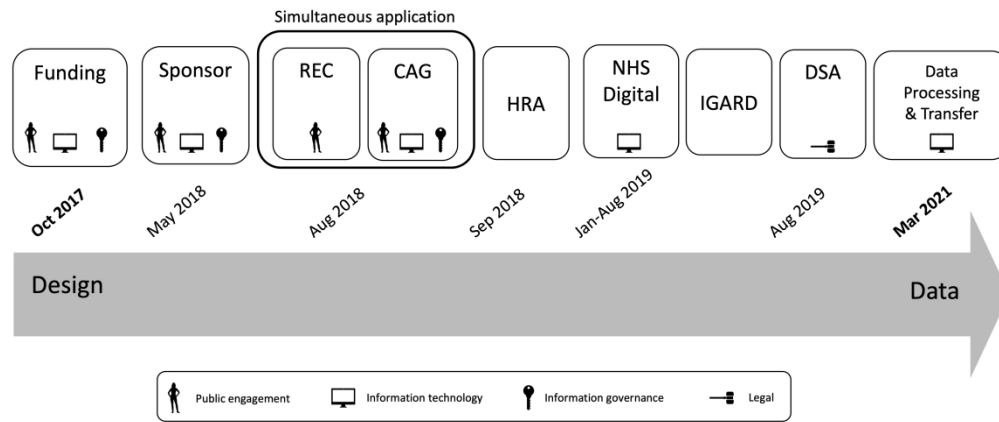
- 16 Luben R, Hayat S, Khawaja A, *et al*. Residential area deprivation and risk of subsequent hospital admission in a British population: the EPIC-Norfolk cohort. *BMJ Open*. 2019;**9**:e031251. doi:10.1136/bmjopen-2019-031251
- 17 Large P. Annual mid-year population estimates, UK - Office for National Statistics. 2014. <https://www.ons.gov.uk/peoplepopulationandcommunity/populationandmigration/populationestimates/bulletins/annualmidyearpopulationestimates/2014-06-26> (accessed 23 Jul 2021).
- 18 Ministry of Housing, Communities & Local Government. English indices of deprivation 2015. 2015. <https://www.gov.uk/government/statistics/english-indices-of-deprivation-2015> (accessed 18 Sep 2020).
- 19 Great Britain. *National Health Service Act 2006*. The Stationery Office 2006.
- 20 Data Access Request Service (DARS). <https://digital.nhs.uk/services/data-access-request-service-dars> (accessed 7 Jul 2021).
- 21 Independent Group Advising on the Release of Data. <https://webarchive.nationalarchives.gov.uk/20200706184649/https://digital.nhs.uk/about-nhs-digital/corporate-information-and-documents/independent-group-advising-on-the-release-of-data> (accessed 7 Jul 2021).
- 22 Perez-Rovira A, MacGillivray T, Trucco E, *et al*. VAMPIRE: Vessel assessment and measurement platform for images of the RETina. *Conf Proc IEEE Eng Med Biol Soc* 2011;**2011**:3391–4.
- 23 Liew G, Wang JJ, Cheung N, *et al*. The retinal vasculature as a fractal: methodology, reliability, and relationship to blood pressure. *Ophthalmology* 2008;**115**:1951–6.
- 24 Keane PA, Grossi CM, Foster PJ, *et al*. Optical Coherence Tomography in the UK Biobank Study - Rapid Automated Analysis of Retinal Thickness for Large Population-Based Studies. *PLoS One* 2016;**11**:e0164095.
- 25 Staurenghi G, Sadda S, Chakravarthy U, *et al*. Proposed lexicon for anatomic landmarks in normal posterior segment spectral-domain optical coherence tomography: the IN•OCT consensus. *Ophthalmology* 2014;**121**:1572–8.
- 26 Peate I. The NHS diabetic eye screening programme. *British Journal of Healthcare Assistants*. 2019;**13**:596–9. doi:10.12968/bjha.2019.13.12.596
- 27 Fasler K, Fu DJ, Moraes G, *et al*. Moorfields AMD database report 2: fellow eye involvement with neovascular age-related macular degeneration. *Br J Ophthalmol* 2020;**104**:684–90.
- 28 Fasler K, Moraes G, Wagner S, *et al*. One- and two-year visual outcomes from the Moorfields age-related macular degeneration database: a retrospective cohort study and an open science resource. *BMJ Open* 2019;**9**:e027441.
- 29 ICD-10 Version:2010. <https://icd.who.int/browse10/2010/en> (accessed 22 Jul 2021).
- 30 Asaria P, Elliott P, Douglass M, *et al*. Acute myocardial infarction hospital admissions and deaths in England: a national follow-back and follow-forward record-linkage study. The

- 1  
2  
3 Lancet Public Health. 2017;**2**:e191–201. doi:10.1016/s2468-2667(17)30032-4  
4  
5 31 Metcalfe A, Neudam A, Forde S, *et al.* Case Definitions for Acute Myocardial Infarction in  
6 Administrative Databases and Their Impact on In-Hospital Mortality Rates. *Health Services*  
7 *Research*. 2013;**48**:290–318. doi:10.1111/j.1475-6773.2012.01440.x  
8  
9 32 McCormick N, Lacaille D, Bhole V, *et al.* Validity of myocardial infarction diagnoses in  
10 administrative databases: a systematic review. *PLoS One* 2014;**9**:e92286.  
11  
12 33 [No title]. [https://biobank.ndph.ox.ac.uk/showcase/showcase/docs/alg\\_outcome\\_stroke.pdf](https://biobank.ndph.ox.ac.uk/showcase/showcase/docs/alg_outcome_stroke.pdf)  
13 (accessed 23 Aug 2021).  
14  
15 34 Brown A, Kirichek O, Balkwill A, *et al.* Comparison of dementia recorded in routinely  
16 collected hospital admission data in England with dementia recorded in primary care.  
17 *Emerg Themes Epidemiol* 2016;**13**:11.  
18  
19 35 Lea NC, Nicholls J, Dobbs C, *et al.* Data Safe Havens and Trust: Toward a Common  
20 Understanding of Trusted Research Platforms for Governing Secure and Ethical Health  
21 Research. *JMIR Medical Informatics*. 2016;**4**:e22. doi:10.2196/medinform.5571  
22  
23 36 Certificate Client Directory search results. [https://www.bsigroup.com/en-GB/our-](https://www.bsigroup.com/en-GB/our-services/certification/certificate-and-client-directory/search-results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bLondondondon&licencenumber=IS%20612909)  
24 [services/certification/certificate-and-client-directory/search-](https://www.bsigroup.com/en-GB/our-services/certification/certificate-and-client-directory/search-results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bLondondondon&licencenumber=IS%20612909)  
25 [results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bL](https://www.bsigroup.com/en-GB/our-services/certification/certificate-and-client-directory/search-results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bLondondondon&licencenumber=IS%20612909)  
26 [ondon&licencenumber=IS%20612909](https://www.bsigroup.com/en-GB/our-services/certification/certificate-and-client-directory/search-results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bLondondondon&licencenumber=IS%20612909) (accessed 26 Jul 2021).  
27  
28 37 Cox DR. Regression Models and Life-Tables. *Journal of the Royal Statistical Society:*  
29 *Series B (Methodological)*. 1972;**34**:187–202. doi:10.1111/j.2517-6161.1972.tb00899.x  
30  
31 38 Fine JP, Gray RJ. A Proportional Hazards Model for the Subdistribution of a Competing  
32 Risk. *Journal of the American Statistical Association*. 1999;**94**:496–509.  
33 doi:10.1080/01621459.1999.10474144  
34  
35 39 Mutlu U, Colijn JM, Ikram MA, *et al.* Association of Retinal Neurodegeneration on Optical  
36 Coherence Tomography With Dementia: A Population-Based Study. *JAMA Neurol*  
37 2018;**75**:1256–63.  
38  
39 40 Harper G. Linkage of Maternity Hospital Episode Statistics data to birth registration and  
40 notification records for births in England 2005–2014: Quality assurance of linkage of routine  
41 data for singleton and multiple births. *BMJ Open*. 2018;**8**:e017898. doi:10.1136/bmjopen-  
42 2017-017898  
43  
44 41 Chua SYL, Thomas D, Allen N, *et al.* Cohort profile: design and methods in the eye and  
45 vision consortium of UK Biobank. *BMJ Open* 2019;**9**:e025077.  
46  
47 42 Mutlu U, Bonnemaier PWM, Ikram MA, *et al.* Retinal neurodegeneration and brain MRI  
48 markers: the Rotterdam Study. *Neurobiol Aging* 2017;**60**:183–91.  
49  
50 43 Majithia S, Tham Y-C, Chee M-L, *et al.* Cohort Profile: The Singapore Epidemiology of Eye  
51 Diseases study (SEED). *Int J Epidemiol* 2021;**50**:41–52.  
52  
53 44 Khan SM, Liu X, Nath S, *et al.* A global review of publicly available datasets for  
54 ophthalmological imaging: barriers to access, usability, and generalisability. *Lancet Digit*  
55 *Health* 2021;**3**:e51–66.  
56  
57  
58  
59

- 1  
2  
3 45 Age groups. 2018. <https://www.ethnicity-facts-figures.service.gov.uk/uk-population-by-ethnicity/demographics/age-groups/latest#:~:text=by%20ethnicity%20Summary-,The%20data%20shows%20that%3A,aged%2060%20years%20and%20over> (accessed 15 Jul 2021).
- 4  
5  
6  
7  
8 46 Enoch J, McDonald L, Jones L, *et al.* Evaluating Whether Sight Is the Most Valued Sense. *JAMA Ophthalmol* 2019;**137**:1317–20.
- 9  
10  
11 47 Sabanayagam C, Xu D, Ting DSW, *et al.* A deep learning algorithm to detect chronic kidney disease from retinal photographs in community-based populations. *Lancet Digit Health* 2020;**2**:e295–302.
- 12  
13  
14  
15 48 Mitani A, Huang A, Venugopalan S, *et al.* Detection of anaemia from retinal fundus images via deep learning. *Nat Biomed Eng* 2020;**4**:18–27.
- 16  
17  
18 49 Poplin R, Varadarajan AV, Blumer K, *et al.* Prediction of cardiovascular risk factors from retinal fundus photographs via deep learning. *Nat Biomed Eng* 2018;**2**:158–64.
- 19  
20  
21 50 Wisely CE, Wang D, Henao R, *et al.* Convolutional neural network to identify symptomatic Alzheimer's disease using multimodal retinal imaging. *Br J Ophthalmol* Published Online First: 26 November 2020. doi:10.1136/bjophthalmol-2020-317659
- 22  
23  
24  
25 51 Cheung CY, Xu D, Cheng C-Y, *et al.* A deep-learning system for the assessment of cardiovascular disease risk via the measurement of retinal-vessel calibre. *Nat Biomed Eng* 2021;**5**:498–508.
- 26  
27  
28  
29 52 Xiao W, Huang X, Wang JH, *et al.* Screening and identifying hepatobiliary diseases through deep learning using ocular images: a prospective, multicentre study. *Lancet Digit Health* 2021;**3**:e88–97.
- 30  
31  
32  
33 53 Wagner SK, Fu DJ, Faes L, *et al.* Insights into Systemic Disease through Retinal Imaging-Based Oculomics. *Transl Vis Sci Technol* 2020;**9**:6.
- 34  
35  
36  
37 54 Deng J, Dong W, Socher R, *et al.* ImageNet: A large-scale hierarchical image database. 2009 IEEE Conference on Computer Vision and Pattern Recognition. 2009. doi:10.1109/cvprw.2009.5206848
- 38  
39  
40  
41 55 Harron K, Dibben C, Boyd J, *et al.* Challenges in administrative data linkage for research. *Big Data & Society*. 2017;**4**:205395171774567. doi:10.1177/2053951717745678
- 42  
43  
44 56 Data protection by design and default. Published Online First: 9 February 2021. <https://ico.org.uk/for-organisations/guide-to-data-protection/guide-to-the-general-data-protection-regulation-gdpr/accountability-and-governance/data-protection-by-design-and-default/> (accessed 13 Jul 2021).
- 45  
46  
47  
48  
49 57 [No title]. <https://www.ipc.on.ca/wp-content/uploads/Resources/7foundationalprinciples.pdf> (accessed 13 Jul 2021).
- 50  
51  
52 58 Lee AY, Campbell JP, Hwang TS, *et al.* Recommendations for Standardization of Images in Ophthalmology. *Ophthalmology* 2021;**128**:969–70.
- 53  
54  
55 59 NHS and social care data: off-shoring and the use of public cloud services. <https://digital.nhs.uk/data-and-information/looking-after-information/data-security-and->
- 56  
57  
58  
59

1  
2  
3 information-governance/nhs-and-social-care-data-off-shoring-and-the-use-of-public-cloud-  
4 services (accessed 13 Jul 2021).  
5

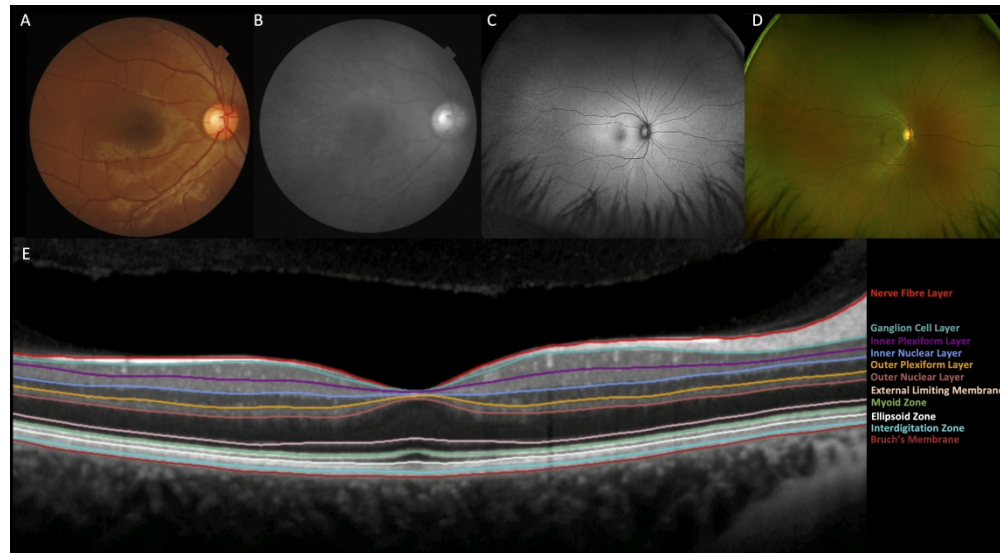
- 6 60 Sinha S, Peach G, Poloniecki JD, *et al*. Studies using English administrative data (Hospital  
7 Episode Statistics) to assess health-care outcomes—systematic review and  
8 recommendations for reporting. *European Journal of Public Health*. 2013;**23**:86–92.  
9 doi:10.1093/eurpub/cks046  
10
- 11 61 Burns EM, Rigby E, Mamidanna R, *et al*. Systematic review of discharge coding accuracy.  
12 *J Public Health* 2012;**34**:138–48.  
13
- 14 62 Herrett E, Shah AD, Boggon R, *et al*. Completeness and diagnostic validity of recording  
15 acute myocardial infarction events in primary care, hospital care, disease registry, and  
16 national mortality records: cohort study. *BMJ* 2013;**346**:f2350.  
17
- 18 63 Wood A, Denholm R, Hollings S, *et al*. Linked electronic health records for research on a  
19 nationwide cohort of more than 54 million people in England: data resource. *BMJ*  
20 2021;**373**:n826.  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



Schematic of the key milestones, prerequisites and approvals with their corresponding achievement dates for the AlzEye dataset. REC: Research Ethics Committee, CAG: Confidential Advisory Group, HRA: Health Research Authority, NHS: National Health Service, IGARD: Independent Group Advising on the Release of Data, DSA: Data Sharing Agreement.

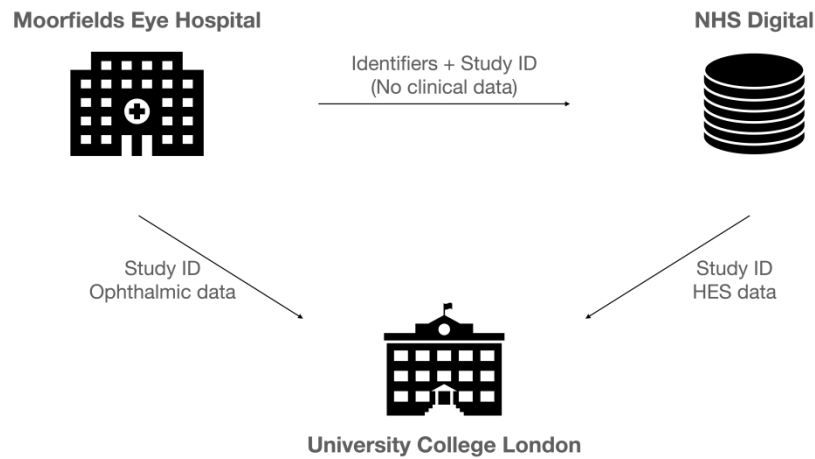
1327x627mm (72 x 72 DPI)





Composite figure showing the major retinal imaging modalities within AlzEye. A: Colour fundus photograph, B: Red-free photograph, C: Fundus autofluorescence (widefield), D: Pseudocolour photography (widefield) and E: Optical coherence tomography of the central macula illustrating segmentation of the individual sublayers. Consensus nomenclature for the retinal sublayers is indicated. NFL= nerve fibre layer, GCL = ganglion cell layer, IPL: inner plexiform layer, INL: inner nuclear layer, OPL: outer plexiform layer, ONL = outer nuclear layer, MZ = myoid zone, RPE-IZ = RPE-interdigitation zone, BM = Bruch's membrane.

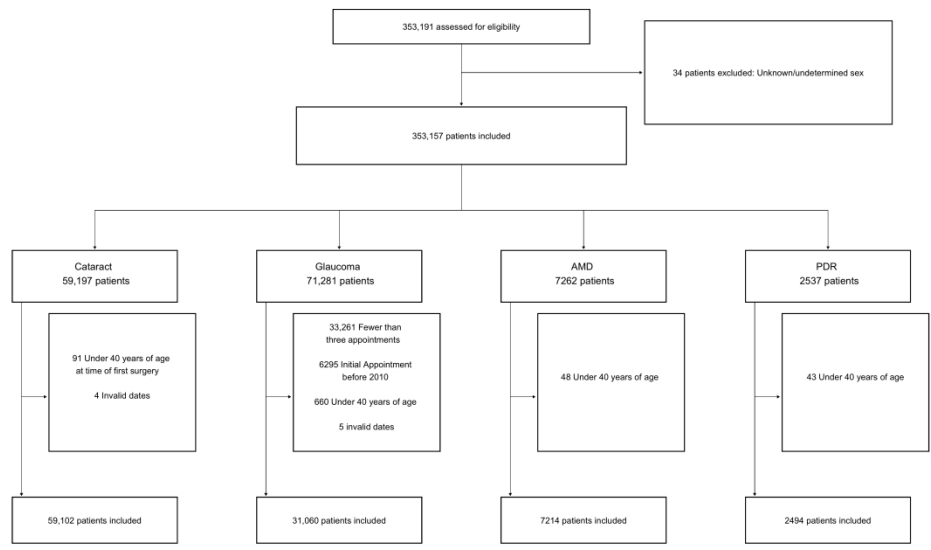
1456x803mm (72 x 72 DPI)



Linkage approach of AlzEye. Moorfields Eye Hospital NHS Foundation Trust securely transfers a spreadsheet of identifiers with a study ID to NHS Digital and separately transfers the study ID with ophthalmic data, including diagnoses and retinal images, to University College London. NHS Digital links the identifiers with the Hospital Episode Statistics (HES) database and returns the admissions data with the study ID (and no identifiable data) to UCL. UCL links the ophthalmic data from Moorfields Eye Hospital with HES data from NHS Digital using the study ID.

2822x1587mm (72 x 72 DPI)

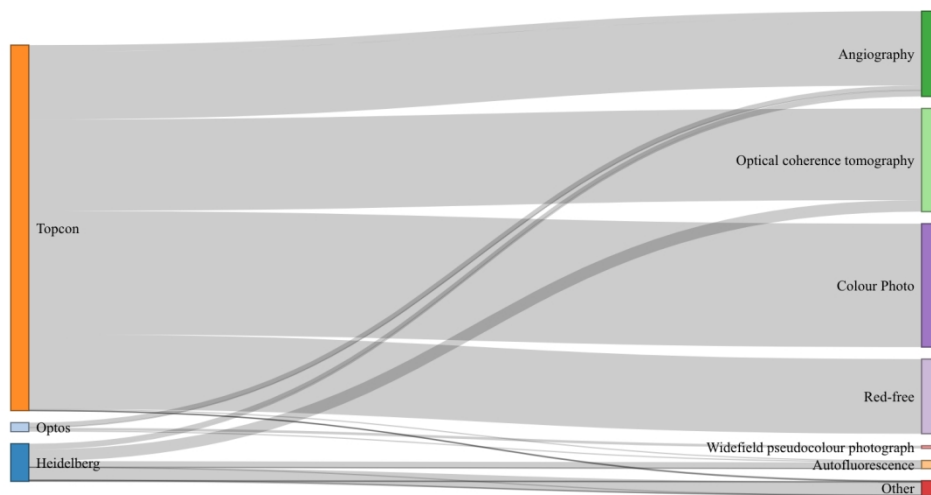
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



CONSORT style flow chart illustrating the distribution of cataract, glaucoma, neovascular age-related macular degeneration (AMD) and proliferative diabetic retinopathy (PDR) within the AlzEye dataset.

1693x1058mm (72 x 72 DPI)

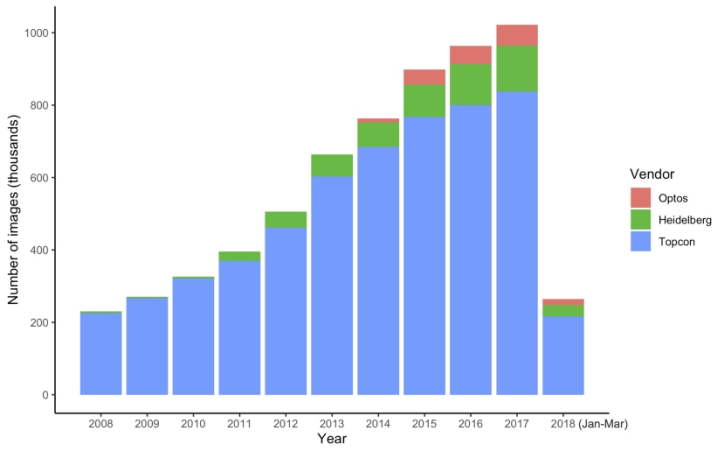
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



Parallel sets diagram illustrating the imaging modality across vendors within AlzEye. The majority of images were acquired on the Topcon system and the most frequent modalities were colour photography and optical coherence tomography. Designed using the networkD3 package.

584x324mm (72 x 72 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



Stacked bar chart of the annual number of images acquired during the study period for the three leading device vendors at Moorfields Eye Hospital. Data for 2018 represents 3 months only prior to the study end-date

1411x793mm (72 x 72 DPI)

# BMJ Open

## Cohort profile: AlzEye: longitudinal record-level linkage of ophthalmic imaging and hospital admissions of 353,157 patients in London, United Kingdom

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-058552.R1
Article Type:	Cohort profile
Date Submitted by the Author:	18-Feb-2022
Complete List of Authors:	Wagner, Siegfried; University College London Institute of Ophthalmology, Hughes, Fintan; Duke University Hospital Cortina Borja, Mario; UCL, Institute of Child Health Pontikos, Nikolas; University College London; Moorfields Eye Hospital NHS Foundation Trust Struyven, Robbert; University College London; Moorfields Eye Hospital NHS Foundation Trust Liu, Xiaoxuan; University Hospitals Birmingham NHS Foundation Trust; University of Birmingham Montgomery, Hugh; University College London, Centre for Human Health and Performance Alexander, Daniel; University College London, Department of Computer Science Topol, Eric; Scripps Research Institute Petersen, Steffen; Queen Mary University of London, William Harvey Research Institute; Barts Health NHS Trust, Balaskas, Konstantinos; Moorfields Eye Hospital NHS Foundation Trust; Moorfields Eye Hospital NHS Foundation Trust Hindley, Jack; University College London, Department of Information Governance Petzold, Axel; UCL; Moorfields Eye Hospital NHS Foundation Trust Rahi, Jugnoo; UCL, Great Ormond Street Institute of Child Health; Great Ormond Street Hospital for Children NHS Foundation Trust Denniston, Alastair; Queen Elizabeth Hospital Birmingham, UK; University of Birmingham Keane, Pearse; University College London; Moorfields Eye Hospital NHS Foundation Trust
<b>Primary Subject Heading</b>:	Ophthalmology
Secondary Subject Heading:	Cardiovascular medicine, Epidemiology, Neurology, Ophthalmology
Keywords:	OPHTHALMOLOGY, Medical ophthalmology < OPHTHALMOLOGY, Medical retina < OPHTHALMOLOGY, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60





I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.



1  
2  
3  
4  
5 Cohort profile: AlzEye: longitudinal record-level linkage of  
6  
7  
8 ophthalmic imaging and hospital admissions of 353,157 patients  
9  
10  
11 in London, United Kingdom  
12  
13  
14  
15

16 SK Wagner<sup>1, 2</sup>, F Hughes<sup>3</sup>, M Cortina-Borja<sup>4</sup>, N Pontikos<sup>1, 2</sup>, RR Struyven<sup>1, 2</sup>, X Liu<sup>5, 6</sup>, H  
17  
18 Montgomery<sup>7</sup>, DC Alexander<sup>8</sup>, E Topol<sup>9</sup>, S Petersen<sup>10, 11</sup>, K Balaskas<sup>1, 2</sup>, J Hindley<sup>12</sup>, A Petzold<sup>1, 2,</sup>  
19  
20  
21 13, 14, J Rahi<sup>1, 4, 15, 16</sup>, AK Denniston<sup>5, 6</sup>, PA Keane<sup>1, 2, \*</sup>  
22  
23  
24  
25

26 \*Corresponding author: Pearse A. Keane, MD FRCOphth, NIHR Biomedical Research Centre for Ophthalmology,  
27  
28 Moorfields Eye Hospital NHS Foundation Trust and UCL Institute of Ophthalmology, 162 City Road, London, United  
29  
30 Kingdom. Tel: +44 207 253 3411 Email: [pearse.keane1@nhs.net](mailto:pearse.keane1@nhs.net)  
31  
32

33 Word count (excluding title page, abstract, references, figures and tables): 4473 words  
34  
35

36 <sup>1</sup> Institute of Ophthalmology, University College London, London, UK,  
37

38 <sup>2</sup> Moorfields Eye Hospital NHS Foundation Trust, London, UK,  
39

40 <sup>3</sup> Duke University Hospital, Durham, North Carolina, USA,  
41

42 <sup>4</sup> Great Ormond Street Institute of Child Health, University College London, London, UK,  
43

44 <sup>5</sup> University of Birmingham, Birmingham, UK,  
45

46 <sup>6</sup> University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK,  
47

48 <sup>7</sup> Centre for Human Health and Performance, University College London, London, UK  
49

50 <sup>8</sup> Centre for Medical Image Computing, Department of Computer Science, University College London, UK  
51

52 <sup>9</sup> Department of Molecular Medicine, Scripps Research, La Jolla, CA, USA.  
53

54 <sup>10</sup> William Harvey Research Institute, Queen Mary University of London, London, UK.  
55

56 <sup>11</sup> St Bartholomew's Hospital, Barts Health National Health Service (NHS) Trust, London, UK..  
57

58 <sup>12</sup> Department of Information Governance, School of Life and Medical Sciences, University College London, London, UK  
59

60 <sup>13</sup> The National Hospital for Neurology and Neurosurgery (UCLH), Queen Square Institute of Neurology, University College London, London, UK.  
61

62 <sup>14</sup> Neuro-ophthalmology Expert Centre, Amsterdam UMC, Departments of Neurology and Ophthalmology, Netherlands.  
63

64 <sup>15</sup> Great Ormond Street Hospital NHS Foundation Trust, London, United Kingdom.  
65

66 <sup>16</sup> Uiverscroft Vision Research Group, University College London, London, UK  
67  
68  
69  
70

## Abstract

Purpose: Retinal signatures of systemic disease ('oculomics') are increasingly being revealed through a combination of high-resolution ophthalmic imaging and sophisticated modelling strategies. Progress is currently limited not mainly by technical issues, but by the lack of large labelled datasets, a sine qua non for deep learning. Such data are derived from prospective epidemiological studies, in which retinal imaging is typically unimodal, cross-sectional, of modest number and relate to cohorts, which are not enriched with subpopulations of interest, such as those with systemic disease. We thus linked longitudinal multimodal retinal imaging from routinely-collected National Health Service data with systemic disease data from hospital admissions using a privacy-by-design third-party linkage approach.

Participants: Between January 1<sup>st</sup> 2008 and April 1<sup>st</sup> 2018, 353,157 participants aged 40 years or older, who attended Moorfields Eye Hospital NHS Foundation Trust, a tertiary ophthalmic institution incorporating a principal central site, four district hubs and five satellite clinics in and around London, United Kingdom serving a catchment population of approximately six million people.

1  
2  
3 Findings to date: Among the 353,157 individuals, 186,651 had a total of 1,337,711 Hospital  
4 Episode Statistics admitted patient care episodes. Systemic diagnoses recorded at these  
5 episodes include 12,022 patients with myocardial infarction, 11,735 with all-cause stroke and  
6 13,363 with all-cause dementia. A total of 6,261,931 retinal images of seven different modalities  
7 and across three manufacturers were acquired from 154,830 patients. The majority of retinal  
8 images were retinal photographs (n=1,874,175) followed by optical coherence tomography  
9 (n=1,567,358).  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

20 Future plans: AlzEye combines the world's largest single institution retinal imaging database  
21 with nationally collected systemic data to create an exceptional large-scale, enriched cohort that  
22 reflects the diversity of the population served. First analyses will address cardiovascular  
23 diseases and dementia, with a view to identifying hidden retinal signatures that may lead to  
24 earlier detection and risk management of these life-threatening conditions.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## STRENGTHS AND LIMITATIONS OF THIS STUDY

- AlzEye is a large retrospective cohort dataset linking ophthalmic data from Moorfields Eye Hospital NHS Foundation Trust in London, United Kingdom with National Health Service hospital admissions data over a ten-year period in 353,157 patients.
- The dataset consists of more than six million routinely collected retinal images of seven different modalities across three vendors deterministically linked to prevalent and incident cardiovascular and neurodegenerative disease.
- Actively informed by ongoing patient and public engagement, the project leverages a privacy-by-design approach using third-party linkage to facilitate access to high-performance computing while mitigating risks to data privacy.

## INTRODUCTION

Scientific discovery has increasingly been driven by the availability of large, diverse, high-dimensional datasets providing deeply phenotyping variables in health and disease[1–3]. Advances in healthcare informatics, hardware and statistical techniques have uncovered relationships previously unachievable through traditional methods of study design. Thus, rich and voluminous genome sequencing data has provided insight into disease pathogenesis and therapeutic targets [4,5] while radiomic analysis has supported exploration of relationships between quantitative data extracted from medical imaging and disease [6].

Crucial to health data research has been the establishment of and accessibility to large prospective epidemiological studies, such as the United Kingdom Biobank (UKBB), the Rotterdam study, and the European Prospective Investigation into Cancer and Nutrition (EPIC) study[7–9]. While such studies represent exceptionally powerful enablers for discovery science, they are potentially limited for investigations of specific subpopulations of interest (e.g. those with rare disease or specific sociodemographic groups) and, where they draw on volunteer participants, also prone to selection bias (e.g. over-representation of more healthy subjects). Participants in UKBB are less likely to be obese, smoke, or drink alcohol and accordingly, mortality rates for participants aged 70-74 in UKBB are 46.2% and 55.5% lower for men and women respectively compared to general UK population[10].

Healthcare in England is such that when a patient has a medical event requiring admission, in almost all cases they are admitted under the provisions of the National Health Service (NHS). Routinely collected healthcare administrative data during a patient's admission are subsequently translated into corresponding International Classification of Diseases (ICD) codes by clinical coders, submitted to the Secondary Uses Service and aggregated by NHS Digital into

1  
2  
3 a unified record-level national repository of Hospital Episode Statistics (HES) data relating  
4 admitted patient care (APC). While the original purpose of HES was the monitoring of service  
5 activity and negotiation of financial reimbursement, it is increasingly used for epidemiological  
6 research[11]. HES data are amenable to research as a sole resource. However, using  
7 deterministic linkage where identifiers are matched in a rules-based approach in contrast to  
8 probabilistic linkage[12], HES can enrich other datasets as in the case of UKBB[13] or the  
9 European Prospective Investigation into Cancer in Norfolk[14].  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

20 While these aforementioned studies demonstrate the value of enriching structured datasets  
21 through HES linkage, this has not yet been done at scale for routinely collected data of high-  
22 dimensionality, such as imaging. Thus, we established AlzEye, a large dataset which links  
23 routinely collected retinal images and relevant ophthalmic data from an unselected population  
24 attending Moorfields Eye Hospital NHS Foundation Trust (MEH) with nationally collected  
25 systemic healthcare outcome data provided through the HES APC database. MEH is a tertiary  
26 ophthalmic institution incorporating a principal central site, four district hubs and five satellite  
27 clinics in London, United Kingdom (UK), providing care to a sociodemographically diverse  
28 population of six million people (9% of the UK population)[15]. The aim of AlzEye is to  
29 characterise the association between retinal biomarkers and chronic disorders of ageing,  
30 particularly dementia and cardiovascular diseases. In addition to describing the characteristics  
31 of the AlzEye cohort, we outline the key governance, technical and ethical factors that need to  
32 be addressed to support large institution-led individual-level linkage of routinely collected multi-  
33 dimensional data and have enabled us to create an exceptional trans-disciplinary resource to  
34 explore the retinal signatures of systemic disease.  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## COHORT DESCRIPTION

### Study population

The AlzEye project is a retrospective cohort study of patients aged 40 years and over who have attended MEH between January 1st 2008 and April 1st 2018. Patients were included if they had attended the glaucoma, retina, neuro-ophthalmology or emergency ophthalmic services and had valid NHS numbers. Those with invalid NHS numbers, dates of birth or who had previously opted out of their health data being used for purposes of research (described in the NHS as a 'Type 2 opt-out') were excluded. Ethnicity group was self-reported by the patient as i) Asian or Asian British, ii) Black or Black British, iii) Mixed, iv) Other Ethnic Group, v) White, or vi) Unknown. Socio-economic status was categorised using the Index of multiple deprivation (IMD) decile, which was estimated by permuting the IMD 2015 rank from the patient's postcode through Lower Super Output Areas followed by aggregation into deciles[16]. Mortality data were derived from the MEH database, which is updated on a two-weekly basis using reports extracted from the NHS National Spine, and is completed on an individual basis by the MEH Data Quality team to ensure accuracy. Data are completed on any patients who have ever attended MEH. Mortality data up to the end of the study period, 1st April 2018, were included.

## Approvals and process

The following key steps in the governance processes were required to provide the necessary ongoing assurance within the research ethics framework of the NHS and the legal framework of the UK. In order to support other researchers wishing to establish similar linked cohorts, we provide an explanation of each stage which outlines the principle which that stage addresses, the UK framework that meets that principle, and finally any study-specific considerations that we undertook to not only meet but exceed those requirements (Figure 1).

- Funding

It was necessary to secure funding to deliver the study, and to provide assurance to the sponsor and others that the study would be completed and that the integrity of the study would not be compromised by inadequate resources. In AlzEye, the study was funded through a small grant awarded by Fight for Sight and Alzheimer's Research UK in October 2017 covering the costs of data storage and linkage fees. The funders had no role in the conception, design or analysis of the study.

- Sponsorship

It was necessary to secure a sponsor for the study that would take on 'overall responsibility for proportionate, effective arrangements being in place to set up, run and report on a research project'. In AlzEye, we sought sponsorship from the relevant NHS Trust (MEH) at which all patients had been seen. Sponsorship confirmation was only sought following internal consultation between data protection, information governance, information security and information technology (IT) teams at both MEH and UCL. MEH acted as the Sponsor, with UCL acting as the trusted third party linking retinal images and HES data and providing



1  
2  
3 computational facilities for data analysis. The study and data governance were approved on  
4  
5 24/05/2018 (internal reference: KEAP1004).  
6  
7

8  
9

- Health Research Authority (HRA) Approval

10  
11 For most research studies in England and Wales, including those limited to working with data for  
12  
13 specific projects, Health Research Authority (HRA) approval is required. HRA approvals involve  
14  
15 the assessment of governance and legal compliance of a research study with an independent  
16  
17 ethics review and opinion from the Research Ethics Committee (REC). Depending on other  
18  
19 study characteristics (e.g.gene therapy), additional applications may be required to inform HRA  
20  
21 approval. In England, limited access to confidential patient information without consent may be  
22  
23 granted under the provisions of Section 251 of the National Health Service Act 2006, permitting  
24  
25 temporary lifting of the common law duty of confidentiality around confidential patient  
26  
27 information ‘in the public interest’ or ‘in the interests of improving patient care’[17]. Obtaining  
28  
29 section 251 support requires application to the Confidential Advisory Group (CAG), an  
30  
31 independent body providing expert advice to the HRA for research applications and NHS Digital  
32  
33 for data dissemination. Applications were accordingly made to the Research Ethics Committee  
34  
35 (18/LO/1163, approved 01/08/2018) and the CAG for Section 251 support (18/CAG/0111,  
36  
37 approved 13/09/2018). The National Health Service Health Research Authority gave final  
38  
39 approval on 13/09/2018. Approvals thus far granted the legal basis for submitting an application  
40  
41 to the Data Access Request Service (DARS) of NHS Digital[18].  
42  
43  
44  
45

46  
47

- NHS Digital and the Data Access Request Service (DARS)

48  
49 NHS Digital oversees the Data Access Request Service (DARS), which administers and  
50  
51 provides, upon application, multiple England-wide datasets from disease-specific audits (e.g.  
52  
53 National Diabetes Audit Core) to general admissions in secondary care (e.g. HES). Applications  
54  
55 to DARS require that the organisation have, at a minimum, the following:  
56  
57  
58  
59

- 1
- 2
- 3 i) Data Sharing Framework Contract for Data Controllers
- 4
- 5 ii) Compliance with minimum-security standards for Data Processors and Data
- 6
- 7 Storage locations
- 8
- 9 iii) Adequate information security certification (e.g. ISO27001)
- 10
- 11 iv) A legal basis for data access (e.g. Section 251)
- 12
- 13
- 14
- 15

16 Applications are then reviewed with an assigned case officer, who will liaise with the applicant  
17 on project-specific items. For AlzEye, dialogue between the applicant and NHS Digital data  
18 production team revolved around confirmation of data fields and datasets (HES) and the  
19 pseudonymisation embedded within the linkage strategy.  
20  
21  
22  
23

- 24
- 25
- 26 ● Independent Group Advising on the Release of Data (IGARD)
- 27

28 Following internal NHS Digital review and prior to data release, DARS applications are  
29 scrutinised by the Independent Group Advising on the Release of Data (IGARD) in line with  
30 Section 263(2) of the Health and Social Care Act 2012, the Code of Practice on confidential  
31 information. IGARD is an independent panel with a broad range of expertise, from legal to  
32 information governance to epidemiology. Support for AlzEye was given by IGARD in January  
33 and August 2019 citing that 'aspects of the application could be used as an exemplar by NHS  
34 Digital to help other researchers with their applications to the Data Access Request Service  
35 (DARS)'<sup>[19]</sup>.  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45

- 46
- 47 ● Data sharing agreements (DSA)
- 48

49 Prior to data receipt, data sharing agreements (DSA) must be signed between NHS Digital and  
50 the Data Controller and are overseen by their respective legal departments. AlzEye required an  
51 additional DSA between MEH and UCL for the transfer of ophthalmic imaging and clinical data  
52 between institutions outlining the purpose and legal basis for sharing.  
53  
54  
55  
56  
57  
58  
59  
60

- Data processing and transfer

The dataset was finalised upon completion of engineering work parsing manufacturer-specific file formats to non-proprietary data structures amenable to image analysis with appropriate deidentification. A secure cloud-based informatics pipeline was used for transfer of images to UCL from MEH, the establishment of which was delayed by the COVID-19 pandemic. Imaging data was stored (with backup) across dedicated network-attached storage device within the UCL School of Life and Medical Sciences (SLMS) and only accessible to members of the AlzEye research team. All data entities were listed within the UCL SLMS Information Asset Register.

## **Patient and public involvement and engagement**

Patient and public involvement and engagement (PPIE) support was provided by the National Institute for Health Research (NIHR) Biomedical Research Centre (BRC) at Moorfields Eye Hospital and UCL Institute of Ophthalmology and has been embedded throughout this project from priority setting to plans for dissemination. Feedback has been sought through public engagement events, survey of eye service users and reports within the media. Patients and the public actively contributed to identifying the priority setting of dementia and the acceptability of using routinely acquired eye scans for research purposes and without consent. In addition, two members of the public will sit on the AlzEye working group to contribute to results interpretation and co-authoring and dissemination of research outputs. The members will be supported in selecting the results they find relevant and presenting them to wider patient communities.

## Ophthalmic health variables

Patient-level ophthalmic variables were extracted from the MEH data warehouse, which aggregates information from the patient administration system (PAS), electronic health record and imaging database, all linked through a unique MEH hospital identification number.

Sociodemographic data, including date of birth, sex, ethnicity and postcode as well as patients' clinic appointments and operation dates are housed within PAS. Surgical procedures were recorded in the electronic health record (EHR) at MEH from September 4th, 2012. Operation details, including procedure name, laterality and indication for surgery are contained within the MEH EHR and uploaded to the MEH data warehouse. A patient undergoing the most common operation in the UK, cataract extraction, would therefore have an entry for the typical procedure, (*phacoemulsification and intraocular lens implant*) operated eye (*right or left*) and indication (*cataract*).

Colour retinal photography (Figure 2A) and optical coherence tomography (OCT, Figure 2B) images, which represent the majority of retinal images within the database, have been processed through segmentation and feature extraction software. The Vascular Assessment and Measurement Platform for Images of the Retina (VAMPIRE) system provides fully automated segmentation and extraction of retinal vascular indices[20,21]. OCT scans are segmented and retinal sublayer thicknesses computed using the Topcon Advanced Biomedical Imaging Laboratory (TABIL) software[22].

For the purposes of this report, four common ophthalmic diseases were described - cataract, glaucoma, neovascular age-related macular degeneration (AMD) and proliferative diabetic retinopathy (PDR).

1  
2  
3 *Cataract* was defined as any operation code denoting phacoemulsification surgery and the  
4 indication of cataract. For the purposes of this report, only first eye cataract surgery was  
5 included.  
6  
7  
8

9 *Glaucoma* was defined as any patient attending the glaucoma clinic three or more times with  
10 ongoing follow up from January 1st 2010. The first two years of the study period were excluded  
11 as this may have incorporated patients with previous diagnoses of glaucoma where the  
12 maximum follow up interval can approach two years; in contrast any patient being seen after 2  
13 years since study inception with no previous visit within that 2 year period can be assumed to  
14 have/carry a new diagnosis of glaucoma.  
15  
16  
17  
18  
19  
20  
21

22 *Diabetic eye disease* represents a special case due to audit procedures mandated by the NHS  
23 Diabetic Eye Screening Programme. Coding of eye disease secondary to diabetes mellitus is  
24 rigorously validated by a dedicated team within MEH according to the NHS Diabetic Eye  
25 Screening Programme criteria[23], at hospital appointment from September 12, 2013 onwards.  
26 Dates for onset of proliferative diabetic retinopathy dates were recorded as the first appointment  
27 for each patient where this diagnosis was first made.  
28  
29  
30  
31  
32  
33

34 *AMD* can be categorised into two major types - dry and neovascular ("wet"). Given dry AMD is  
35 slowly progressive and has no active hospital intervention currently available, it is MEH standard  
36 practice for patients to be discharged with lifestyle and monitoring advice (self-monitoring and  
37 standard optometric review). In contrast, neovascular AMD requires treatment through  
38 intravitreal anti-VEGF injections, and therefore remains under active follow-up. The diagnostic  
39 codes for neovascular AMD were based on extensive previous work in which all patients with  
40 neovascular AMD at MEH were manually validated up to 2018[24,25].  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Systemic health variables

Systemic health data were derived from HES APC data, with a focus on cardiovascular disease and all-cause dementia. Diagnostic codes in HES APC are reported in line with the 10th revision of ICD[26]. In line with previous reports, myocardial infarction was defined as code I21 or I22[27–29]. Similarly, stroke was defined using stroke definitions from UK Biobank[30].

Dementia was defined as ICD codes E512 (Wernicke’s encephalopathy), F00 (Dementia in Alzheimer disease), F01 (Vascular dementia), F02 (Dementia in other diseases classified elsewhere), F03 (Unspecified dementia), F10.6 (Mental and behavioural disorders due to psychoactive substance use, Amnesic syndrome) F10.7 (Mental and behavioural disorders due to psychoactive substance use, Residual and late-onset psychotic disorder), G30 (Alzheimer disease), or G31.0 (Other degenerative diseases of nervous system, not elsewhere classified), derived from previous work evaluating the agreement between HES admitted patient care data and primary care data, through general practitioner surveys and the Clinical Practice Research Datalink (CPRD)[31].

## Data Linkage and Transfer

The linkage strategy was designed through collaboration between experts in information governance, information technology, computer scientists and clinicians based at MEH, University College London (UCL) and NHS Digital (Figure 3). A third-party linkage approach was used for two main reasons. Firstly, it enhanced privacy-preservation as the data originator, MEH, never received HES admissions data and the third party, UCL, did not receive personally identifiable information. Secondly, it enabled the linked dataset to be accessible within a site with sufficient high-performance computing capability to undertake the proposed analyses, a function significantly beyond almost all NHS facilities. Patient link identifiers consisting of a

1  
2  
3 unique NHS identification number, sex and date of birth originating from MEH were transferred  
4  
5 to NHS Digital in conjunction with a unique study ID generated using a cryptographic hash  
6  
7 function (random pseudonymisation). Ophthalmic covariates, mortality data, and patient  
8  
9 sociodemographics with study ID were transferred to UCL. Ophthalmic imaging data pertaining  
10  
11 to the patients within the study were extracted and de-identified during conversion from their  
12  
13 proprietary format to Digital Imaging and Communications in Medicine (DICOM) format before  
14  
15 transfer to UCL. Following linkage with HES, NHS Digital transferred HES data to the UCL Data  
16  
17 Safe Haven, a “walled garden” trusted research environment certified and externally audited to  
18  
19 ISO27001 information security standards[32,33].  
20  
21  
22  
23  
24

## 25 **Statistical analysis**

26  
27  
28 Imaging-based studies within the AlzEye study are generally planned to take the form of nested  
29  
30 case-control studies. To improve efficiency, controls may be matched with cases, using  
31  
32 conditional logistic regression for statistical modelling of binary outcomes and survival analysis  
33  
34 for time-to-event data (e.g. Cox Proportional hazard modelling)[34]. In cases where the  
35  
36 competing risk of death is prominent, subdistribution hazard ratios with 95% confidence  
37  
38 intervals will be estimated as a sensitivity analysis [35]. Alternative high-dimensional modelling  
39  
40 approaches, such as vision transformers, will also be explored. Prior to receipt of HES data from  
41  
42 NHS Digital, sample size calculations were undertaken. Specifically, we evaluated the  
43  
44 association between OCT-derived peripapillary retinal nerve fibre layer and macular ganglion  
45  
46 cell-inner plexiform layer thicknesses and dementia. Given an odds ratio of 1.4 with an alpha of  
47  
48 5% and a power of 90% on a 1:1 matched study design, a total sample size of 2106 is  
49  
50 required[36].  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Figures for this report were designed in R version 4.1.0 (R Core Team, 2021. R Foundation for  
4  
5 Statistical Computing, Vienna, Austria).  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



## FINDINGS TO DATE

Extraction of unique patients attending MEH outpatient clinics between January 1st 2008 and April 1st 2018 generated a cohort of 353,157 unique patients. A breakdown of sociodemographic details by category of the cohort are provided in Table 1. Of the cohort, 190,494 were female (53.9%) and the mean age was 68.4 (+/- 13.9) years. Of the 353,157 patients, 186,651 had a total of 1,337,711 HES episodes in the study period. NHS Digital performs a hierarchical stepwise linkage approach providing a 'Match Rank' for each HES episode[37]. Among the 1,337,711 HES episodes matched, Match Rank was two for 1,337,482 episodes (exact NHS number, exact date of birth and exact sex linked), four for 46 episodes (exact NHS number, exact sex and partial date of birth), and eight for 183 episodes (exact NHS number).

An illustration of the major common ophthalmic diseases within the cohort is shown in a CONSORT-style diagram in Figure 4. Following the case definition, and exclusion of invalid dates, a total of 59,102 patients had first eye cataract surgery, 31,060 glaucoma, 7214 neovascular AMD and 2494 PDR.

	Characteristic	N (%)
--	----------------	-------

	<b>All</b>	353,157
<b>Sex</b>	<b>Female</b>	190,494 (53.9)
	<b>Male</b>	162,663 (46.1)
<b>Age group*</b>	<b>40-49 years</b>	35,262 (10.0)
	<b>50-59 years</b>	66,101 (18.7)
	<b>60-69 years</b>	79,018 (22.4)
	<b>70-79 years</b>	84,942 (24.1)
	<b>80+ years</b>	87,834 (24.9)
<b>Ethnicity</b>	<b>Black</b>	31,614 (9.0)
	<b>White</b>	135,743 (38.4)
	<b>South Asian</b>	48,119 (13.6)
	<b>Other/Unknown</b>	137,681 (39.0)
<b>Index of multiple deprivation decile</b>	<b>1 (most deprived)</b>	18,194 (5.2)
	<b>2</b>	50,443 (14.3)
	<b>3</b>	50,869 (14.4)
	<b>4</b>	42,603 (12.1)
	<b>5</b>	38,964 (11.0)
	<b>6</b>	36,906 (10.5)
	<b>7</b>	31,317 (8.9)
	<b>8</b>	28,180 (8.0)
	<b>9</b>	29,906 (8.5)
	<b>10 (least deprived)</b>	24,610 (7.0)
	<b>Unknown</b>	1165 (0.3)

Table 1: Baseline sociodemographic characteristics of the AlzEye cohort. Data are shown as n(%).

\*Age is taken as that of April 1st, 2018.

Among the 187,811 patients with recorded HES episodes, 12,022 patients had episodes with coded myocardial infarction, 11,735 patients with all-cause stroke and 13,363 with dementia. Within the dementia group, 4487 patients had codes that were specific for Alzheimer's dementia and 3381 for vascular dementia (Table 2).

Group	Disease	ICD code(s)	Number of patients
Cardiovascular	Acute coronary syndrome	I21, I22	12,022
	Heart failure	I50	24,034
	Atrial Fibrillation	I48	32,848
	Hypertension	I10, I15	151,937
	Subarachnoid haemorrhage	I60	642
	Intracerebral hemorrhage	I61	1865
	Ischaemic stroke	I63-I64	9996
	All stroke	I60, I61, I63, I64	11,735
Neurodegenerative	Alzheimer's disease	F00, G30	4487
	Vascular dementia	F01	3381
	Parkinson's disease	G20	3211
	All-cause dementia	E12, F00, F01, F02, F03, F106, F107, G30, G310	13,363
Other	Diabetes mellitus (Type 1 and 2)	E10, E11	71,570

Table 2: Number of patients by selected examples of specified 10th revision of International Classification of Diseases (ICD) codes relating to diabetes mellitus, cardiovascular and neurodegenerative diseases.

## Imaging

During the study period, a total of 6,261,931 images were acquired from 154,830 patients. The two leading image modalities were colour retinal photographs (n=1,874,175) and OCT (n=1,567,358). The distribution of imaging modalities across the three vendors used for retinal imaging at MEH - Topcon (Topcon corp, Tokyo, Japan), Heidelberg (Heidelberg Engineering, Heidelberg, Germany), and Optos (Optos, Dunfermline, UK) - are shown in Figure 5 and Table

2. Most images were acquired on the Topcon system (n=5,553,826, 88.7%). Number of images by year is shown in Figure 6. During the study period, annual imaging acquisition increased from 229,868 scans in 2008 to 1,021,904 in 2017. For 2018, collection stopped on April 1st precluding a complete annual figure. Example images of the major ophthalmic and systemic disease outcomes are shown in Figure 7.

Vendor	Modality	Number of images	Number of patients
<b>Topcon</b>	Angiography	1,128,723	21,225
	Autofluorescence	11,761	2078
	Colour photography	1,874,175	139,307
	Red-free	1,146,854	122,453
	OCT	1,391,826	138,911
	Other	487	48
<b>Heidelberg</b>	Angiography	89,264	4061
	Autofluorescence	94,533	16,863
	Infrared	192,634	21,676
	OCT	175,532	21,191
	Other	19,781	2439
<b>Optos</b>	Angiography	77,813	2215
	Autofluorescence	18,590	5666
	Pseudocolour photography	39,958	6887

Table 3: Retinal imaging within the AlzEye dataset by vendor and imaging modality. OCT: Optical coherence tomography.

Angiography refers to dye-based techniques (fluorescein and indocyanine green).

## STRENGTHS AND LIMITATIONS

To our knowledge, we have created the world's largest retinal imaging research dataset available presently, linking secondary healthcare ophthalmic data from 353,157 patients seen over a ten-year period with information on general health and key systemic diseases, as captured through admissions to any hospital within the NHS of England. This comprises 6,261,931 images, obtained using seven different modalities from three different manufacturers, in 154,830 patients. The current large-scale UK cohort, UKBB, provides useful context for AlzEye. Cross-sectional data are available in UKBB with two retinal imaging modalities (colour retinal photography and OCT) obtained using technology from one manufacturer (Topcon), and at a single time point in 67,321 people. Notwithstanding the recognised limitations (see 'Limitations' section) of real-world datasets and the coding within the HES database, AlzEye provides some distinct advantages beyond purely scale. Imaging data are longitudinal, highly multimodal, and pertain to an ethnically and socioeconomically diverse cohort representative of the adult population with eye disease. Moreover, AlzEye has demonstrated relatively low cost. The study is funded through a charity small grant award and NIHR Biomedical Research Centre support amounting to £20,000.

### **Comparison with other resources**

UKBB is the major comparator for AlzEye, being the largest of the prospective epidemiological cohort datasets which provide cross-sectional retinal imaging in association with systemic disease variables[38]. One of the limitations of UKBB is that, unlike AlzEye, it provides minimal longitudinal retinal images. Another prospective cohort study, the Rotterdam study, does collect longitudinal retinal imaging data from approximately 15,000 participants, of which 5065 participants were eligible for OCT scanning in 2017[39]. The Rotterdam Study has uncovered several landmark findings, particularly in regards to causal determinants, but its cohort remains

1  
2  
3 relatively small in comparison to UKBB and AlzEye with the majority of participants recruited  
4 from one district within Rotterdam, Netherlands[7,40]. The Singapore Epidemiology of Eye  
5 Disease (SEED), is one longitudinal multimodal retinal imaging initiative which is underway, in  
6 which 10,033 participants of Chinese, Indian and Malay ethnicity have been recruited to  
7 undergo six-yearly retinal imaging[41]. A recent review of ophthalmic imaging datasets did not  
8 reveal any additional relevant publicly available datasets that included linked systemic health  
9 data[42]. Additionally, our own review of the literature has not identified any examples of large-scale  
10 linked real-world datasets (i.e. including those with restricted access) which include linked systemic  
11 health data. The scarcity of such resources suggests that the construction of such datasets is  
12 challenging to undertake, presumably due to factors such as cost, required duration and delayed  
13 output, retention of participants, and concerns over technological redundancy. The AlzEye approach  
14 is an important alternative model in this context.  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

### 31 **Potential research impact from the novel AlzEye cohort**

32  
33 Several epidemiological opportunities arise with AlzEye. Firstly, it provides a real-world  
34 snapshot of ophthalmic secondary care use, representing approximately 1.2% of the UK  
35 population aged 40 and above (27,858,459 in 2011) [43]. This is a powerful tool for informing  
36 public health and policymaking in eye services, and is exceptional in characterising the potential  
37 impact that may arise from the intersect between disabling diseases such as stroke and PDR.  
38  
39  
40  
41  
42  
43  
44  
45

46 Secondly, it allows the identification and exploration of relationships between newly diagnosed  
47 ophthalmic disease (or newly referred to hospital eye services) and emerging systemic events  
48 and accruing multimorbidity. Patients tend to respond early to issues with their sight and an  
49 understanding of how an ophthalmic presentation is linked to an increased likelihood of serious  
50 systemic disease may provide an opportunity for earlier intervention in those diseases[44].  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5 Thirdly, nested case-control studies evaluating retinal-based oculomic biomarkers in those with  
6 systemic diseases (e.g. dementia) can provide insight into their value in either static or dynamic  
7 risk prediction. Newer modelling approaches have highlighted the potential utility of the retina in  
8 screening for and risk stratification of cardiovascular, neurodegenerative, renal, hepatic and  
9 haematological diseases[45–50];[51].  
10  
11  
12  
13  
14  
15

16  
17  
18 Finally, by its magnitude and wealth of high-quality labels, both ophthalmic and systemic,  
19 AlzEye provides a powerful catalyst for high-dimensional model development, echoing that of  
20 Imagenet, a database currently exceeding 14 million images, which propelled deep learning and  
21 computer vision research forward a decade ago[52].  
22  
23  
24  
25  
26  
27

### 28 **Lessons learned from the AlzEye approach**

29  
30 AlzEye highlights an opportunity for maximising the value of routinely-collected data to support  
31 research for patient benefit. However, there are a number of governance and technical  
32 challenges when undertaking large scale investigator-led data linkage[53]. In AlzEye, early  
33 dialogue between experts in information governance, information technology and data  
34 protection at both institutional parties (MEH and UCL) as well as the data production team at  
35 NHS Digital established a privacy-by-design linkage approach, which enhanced privacy  
36 preservation while maintaining computational feasibility[30,54]. At its worst, an intrusion of the  
37 identifiable data during the development of AlzEye would have informed the violator that a given  
38 individual had visited MEH at some point between 2008 and 2018. Due to the novel approach of  
39 AlzEye within our centre, the greatest governance hurdle was securing study sponsorship, a  
40 process which took nearly eight months. Once approved, permissions from the bodies of the  
41 Health Research Authority, and overall approval were given within eight weeks. Linking with  
42 high-dimensional imaging data also posed several technical obstacles. As highlighted recently  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 by the American Academy of Ophthalmology (AAO), ophthalmic imaging technologies suffer  
4 from limited interoperability and low compliance to standardised formats, such as DICOM [55]. A  
5 key undertaking within AlzEye was thus the secure and robust but efficient fully-automated  
6 processing of raw ophthalmic imaging data from its proprietary file format with associated  
7 metadata to standard DICOM form with the identifiers stripped. Fortunately, while this operation  
8 requires significant technical and engineering input, most medical imaging modalities already  
9 benefit from standardisation among vendors obviating this step for other researchers seeking to  
10 emulate our approach. Finally, a key objective of AlzEye is the development of clinical prediction  
11 models using deep learning approaches, which require significant computing capacity.  
12 Provisions for graphics processing units (GPU) housed within UCL enable this step however  
13 others may consider recent guidance on the safeguards required for locating health data within  
14 cloud environments and the implications this brings for accessing virtual GPUs[56].  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

### 31 **Limitations of the AlzEye cohort**

32 Despite the opportunities afforded by AlzEye, there are several limitations to this kind of  
33 approach and potential sources of bias. Firstly, caution must be paid to the validity of HES  
34 diagnostic coding[57]. Although previous validation studies have concluded that discharge  
35 coding within HES is sufficiently robust for research purposes[31,58], sizable proportions of  
36 cases may be missed when using individual sources[59]. For example, recent work linking the  
37 electronic health records of 54.4 million people in England showed that HES captured 80.5%  
38 and 65% of myocardial infarctions and stroke/transient ischaemic attacks respectively when  
39 compared to linkage additionally incorporating Death Registry and primary care records[60].  
40 One mitigation strategy for this source of bias for real world data is therefore linking to multiple  
41 sources. In terms of selection bias, as a hospital-attending cohort, the individuals within the  
42 AlzEye cohort are likely to have greater medical comorbidity than the general population,  
43 limiting the external validity of any findings. In addition, by the very nature of the dataset,  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 patients within the AlzEye cohort will have definite or suspected ophthalmic disease, particularly  
4 among those with repeated retinal imaging. The risk of under-recording of potentially important  
5 variables such as smoking may also lead to residual confounding.  
6  
7  
8  
9

10  
11 The enrichment of multimodal health data acquired as part of a patient's routine clinical care  
12 with nationally-held databases provides a powerful foundation for discovery science and  
13 epidemiological research. We highlight key considerations and challenges for those seeking to  
14 link high-dimensional data sources, from high-resolution imaging to waveform data, with locally-  
15 held specialist data. Additionally, we provide the cohort profile for AlzEye, a powerful platform  
16 for oculomic discovery, specifically evaluating the association between retinal morphology, and  
17 both cardiovascular diseases and dementia. Beyond discovery, the AlzEye cohort is anticipated  
18 to become an important resource for the development and validation of deep learning-based  
19 clinical prediction models that may enable earlier intervention for patients at risk of these life-  
20 threatening conditions.  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## LEGENDS

Figure 1: Schematic of the key milestones, prerequisites and approvals with their corresponding achievement dates for the AlzEye dataset. REC: Research Ethics Committee, CAG: Confidential Advisory Group, HRA: Health Research Authority, NHS: National Health Service, IGARD: Independent Group Advising on the Release of Data, DSA: Data Sharing Agreement.

Figure 2: Composite figure showing the major retinal imaging modalities within AlzEye. A: Colour fundus photograph, B: Red-free photograph, C: Fundus autofluorescence (widefield), D: Pseudocolour photography (widefield) and E: Optical coherence tomography of the central macula illustrating segmentation of the individual sublayers. Consensus nomenclature for the retinal sublayers is indicated. NFL= nerve fibre layer, GCL = ganglion cell layer, IPL: inner plexiform layer, INL: inner nuclear layer, OPL: outer plexiform layer, ONL = outer nuclear layer, MZ = myoid zone, RPE-IZ = RPE-interdigitation zone, BM = Bruch's membrane.

Figure 3: Linkage approach of AlzEye. Moorfields Eye Hospital NHS Foundation Trust securely transfers a spreadsheet of identifiers with a study ID to NHS Digital and separately transfers the study ID with ophthalmic data, including diagnoses and retinal images, to University College London. NHS Digital links the identifiers with the Hospital Episode Statistics (HES) database and returns the admissions data with the study ID (and no identifiable data) to UCL. UCL links the ophthalmic data from Moorfields Eye Hospital with HES data from NHS Digital using the study ID.

Figure 4: CONSORT style flow chart illustrating the distribution of cataract, glaucoma, neovascular age-related macular degeneration (AMD) and proliferative diabetic retinopathy (PDR) within the AlzEye dataset.

Figure 5: Parallel sets diagram illustrating the imaging modality across vendors within AlzEye. The majority of images were acquired on the Topcon system and the most frequent modalities were colour photography and optical coherence tomography. Designed using the `networkD3` package.

Figure 6: Stacked bar chart of the annual number of images acquired during the study period for the three leading device vendors at Moorfields Eye Hospital. Data for 2018 represents 3 months only prior to the study end-date.

1  
2  
3 Figure 7: Example colour retinal photographs of ophthalmic and systemic diseases within AlzEye. A: age-related  
4 macular degeneration, B: cataract, C: glaucoma, D: proliferative diabetic retinopathy E: prevalent Alzheimer's  
5 disease, F: incident ischaemic stroke, G: incident myocardial infarction, H: prevalent vascular dementia.  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## ACKNOWLEDGEMENTS

We thank Karen Bonstein and Andi Skilton for support with patient and public involvement and engagement, Menachem Katz, Ben Ward, Ross Green, Maxim Daniline and Simon St John-Green for information technology support, Anthony Peacock for advice on use of Trusted Research Environments and Antonio de la Plaza Larrea and Richard Macmillan for legal guidance on data sharing agreements. We also are grateful to Anthony Khawaja and Cathie Sudlow for feedback on study design.

## CONTRIBUTORS

SW, AD and PK wrote the first draft of the manuscript, which was critically revised by FH, MCB, RS, NP, XL, HM, DA, ET, SP, KB, JH, AP, and JR. SW, FH, NP, HM, JH, AP, JR, AD and PK were involved in the original design of the study. Computer science expertise was through RS, NP and DA, information governance through JH. MCB, AP, JR and AD provided statistical and epidemiological guidance. All authors have approved the final version of this manuscript.

## FUNDING STATEMENT

The study was funded through a small grant awarded by Fight for Sight (grant reference: 24AZ171). Infrastructural support was through the NIHR Biomedical Research Centre at Moorfields Eye Hospital and the UCL Institute of Ophthalmology. The funders had no role in the conception, design or analysis of the study.

## COLLABORATION

National and international collaborations are welcomed though restrictions on access to the cohort mean that only the AlzEye researchers can directly analyse individual-level systemic health data. Interested researchers should contact [s.wagner@ucl.ac.uk](mailto:s.wagner@ucl.ac.uk).

## DATA AVAILABILITY STATEMENT

No additional data available. The data are subject to the contractual restrictions of the DSA between NHS Digital, Moorfields Eye Hospital and University College London and are therefore not available for access beyond the AlzEye research team.

## COMPETING INTERESTS STATEMENT

SW is funded through a Medical Research Council Clinical Research Training Fellowship (MR/TR000953/1).

FH: None.

MCB: None.

NP is funded by a Moorfields Eye Charity Career Development Award (R190031A).

RS: None.

XL: None.

HM is supported by the National Institute of Health Research's Comprehensive Biomedical Research Centre (BRC) at University College London Hospitals (UCLH).

DA: None

ET: None

1  
2  
3 SP receives support from the National Institute for Health Research Biomedical Research  
4  
5 Centre at Barts.

6  
7 KB has received speaker fees from Novartis, Bayer, Alimera, Allergan, Roche, and Heidelberg;  
8  
9 meeting or travel fees from Novartis and Bayer; compensation for being on an advisory board  
10  
11 from Novartis and Bayer; consulting fees from Novartis and Roche; and research support from  
12  
13 Apellis, Novartis, and Bayer.

14  
15 JH: None

16  
17 AP receives financial support from the National Institute for Health Research Biomedical  
18  
19 Research Centre based at Moorfields Eye Hospital NHS Foundation Trust and UCL Institute of  
20  
21 Ophthalmology. AP is part of the steering committee of the ANGI network which is sponsored by  
22  
23 ZEISS, steering committee of the OCTiMS study which is sponsored by Novartis, and reports  
24  
25 speaker fees from Heidelberg Engineering.

26  
27 JR receives support from the National Institute for Health Research as a Senior Investigator and  
28  
29 via the National Institute for Health Research Biomedical Research Centres at Moorfields Eye  
30  
31 Hospital and Great Ormond Street Hospital.

32  
33 AD is Director of INSIGHT, the HDRUK Health Data Research Hub for Eye Health.

34  
35 PK is supported by a Moorfields Eye Charity Career Development Award (R190028A) and a UK  
36  
37 Research & Innovation Future Leaders Fellowship (MR/T019050/1); receives research support  
38  
39 from Apellis; is a consultant for DeepMind, Roche, Novartis, Apellis, and BitFount; is an equity  
40  
41 owner in Big Picture Medical; and has received speaker fees from Heidelberg Engineering,  
42  
43 Topcon, Allergan, Roche, and Bayer; meeting or travel fees from Novartis and Bayer; and  
44  
45 compensation for being on an advisory board from Novartis and Bayer.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## ETHICAL APPROVAL STATEMENT

This study was approved by the Research Ethics Committee (18/LO/1163, approved 01/08/2018) and the CAG for Section 251 support (18/CAG/0111, approved 13/09/2018). The National Health Service Health Research Authority gave final approval on 13/09/2018.

## REFERENCES

- 1 Munevar S. Unlocking Big Data for better health. *Nat Biotechnol* 2017;**35**:684–6.
- 2 Shilo S, Rossman H, Segal E. Axes of a revolution: challenges and promises of big data in healthcare. *Nat Med* 2020;**26**:29–38.
- 3 Obermeyer Z, Emanuel EJ. Predicting the Future - Big Data, Machine Learning, and Clinical Medicine. *N Engl J Med* 2016;**375**:1216–9.
- 4 Bhardwaj R, Sethi A, Nambiar R. Big data in genomics: An overview. In: *2014 IEEE International Conference on Big Data (Big Data)*. IEEE 2014. doi:10.1109/bigdata.2014.7004392
- 5 He KY, Ge D, He MM. Big Data Analytics for Genomic Medicine. *Int J Mol Sci* 2017;**18**. doi:10.3390/ijms18020412
- 6 Gillies RJ, Kinahan PE, Hricak H. Radiomics: Images Are More than Pictures, They Are Data. *Radiology* 2016;**278**:563–77.
- 7 Hofman A, Breteler MMB, van Duijn CM, *et al*. The Rotterdam Study: objectives and design update. *Eur J Epidemiol* 2007;**22**:819–29.
- 8 Riboli E, Hunt KJ, Slimani N, *et al*. European Prospective Investigation into Cancer and Nutrition (EPIC): study populations and data collection. *Public Health Nutr* 2002;**5**:1113–24.
- 9 Bycroft C, Freeman C, Petkova D, *et al*. The UK Biobank resource with deep phenotyping and genomic data. *Nature* 2018;**562**:203–9.
- 10 Fry A, Littlejohns TJ, Sudlow C, *et al*. Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank Participants With Those of the General Population. *Am J Epidemiol* 2017;**186**:1026–34.
- 11 Chaudhry Z, Mannan F, Gibson-White A, *et al*. Research Outputs of England's Hospital Episode Statistics (HES) Database: Bibliometric Analysis. *Journal of Innovation in Health Informatics*. 2017;**24**:329. doi:10.14236/jhi.v24i4.949
- 12 Zhu Y, Matsuyama Y, Ohashi Y, *et al*. When to conduct probabilistic linkage vs.

- deterministic linkage? A simulation study. *J Biomed Inform* 2015;**56**:80–6.
- 13 Hospital information boosts UK Biobank resource.  
<https://www.ukbiobank.ac.uk/2013/09/20000-participants-return-for-a-repeat-assessment/>  
(accessed 13 Oct 2020).
- 14 Luben R, Hayat S, Khawaja A, *et al*. Residential area deprivation and risk of subsequent hospital admission in a British population: the EPIC-Norfolk cohort. *BMJ Open*. 2019;**9**:e031251. doi:10.1136/bmjopen-2019-031251
- 15 Large P. Annual mid-year population estimates, UK - Office for National Statistics. 2014.<https://www.ons.gov.uk/peoplepopulationandcommunity/populationandmigration/populationestimates/bulletins/annualmidyearpopulationestimates/2014-06-26> (accessed 23 Jul 2021).
- 16 Ministry of Housing, Communities & Local Government. English indices of deprivation 2015. 2015.<https://www.gov.uk/government/statistics/english-indices-of-deprivation-2015> (accessed 18 Sep 2020).
- 17 Great Britain. *National Health Service Act 2006*. The Stationery Office 2006.
- 18 Data Access Request Service (DARS). <https://digital.nhs.uk/services/data-access-request-service-dars> (accessed 7 Jul 2021).
- 19 Independent Group Advising on the Release of Data.  
<https://webarchive.nationalarchives.gov.uk/20200706184649/https://digital.nhs.uk/about-nhs-digital/corporate-information-and-documents/independent-group-advising-on-the-release-of-data> (accessed 7 Jul 2021).
- 20 Perez-Rovira A, MacGillivray T, Trucco E, *et al*. VAMPIRE: Vessel assessment and measurement platform for images of the REtina. *Conf Proc IEEE Eng Med Biol Soc* 2011;**2011**:3391–4.
- 21 Liew G, Wang JJ, Cheung N, *et al*. The retinal vasculature as a fractal: methodology, reliability, and relationship to blood pressure. *Ophthalmology* 2008;**115**:1951–6.
- 22 Keane PA, Grossi CM, Foster PJ, *et al*. Optical Coherence Tomography in the UK Biobank Study - Rapid Automated Analysis of Retinal Thickness for Large Population-Based Studies. *PLoS One* 2016;**11**:e0164095.
- 23 Peate I. The NHS diabetic eye screening programme. *British Journal of Healthcare Assistants*. 2019;**13**:596–9. doi:10.12968/bjha.2019.13.12.596
- 24 Fasler K, Fu DJ, Moraes G, *et al*. Moorfields AMD database report 2: fellow eye involvement with neovascular age-related macular degeneration. *Br J Ophthalmol* 2020;**104**:684–90.
- 25 Fasler K, Moraes G, Wagner S, *et al*. One- and two-year visual outcomes from the Moorfields age-related macular degeneration database: a retrospective cohort study and an open science resource. *BMJ Open* 2019;**9**:e027441.
- 26 ICD-10 Version:2010. <https://icd.who.int/browse10/2010/en> (accessed 22 Jul 2021).



- 1  
2  
3 27 Asaria P, Elliott P, Douglass M, *et al.* Acute myocardial infarction hospital admissions and  
4 deaths in England: a national follow-back and follow-forward record-linkage study. *The*  
5 *Lancet Public Health.* 2017;**2**:e191–201. doi:10.1016/s2468-2667(17)30032-4  
6  
7 28 Metcalfe A, Neudam A, Forde S, *et al.* Case Definitions for Acute Myocardial Infarction in  
8 Administrative Databases and Their Impact on In-Hospital Mortality Rates. *Health Services*  
9 *Research.* 2013;**48**:290–318. doi:10.1111/j.1475-6773.2012.01440.x  
10  
11 29 McCormick N, Lacaille D, Bhole V, *et al.* Validity of myocardial infarction diagnoses in  
12 administrative databases: a systematic review. *PLoS One* 2014;**9**:e92286.  
13  
14 30 [No title]. [https://biobank.ndph.ox.ac.uk/showcase/showcase/docs/alg\\_outcome\\_stroke.pdf](https://biobank.ndph.ox.ac.uk/showcase/showcase/docs/alg_outcome_stroke.pdf)  
15 (accessed 23 Aug 2021).  
16  
17 31 Brown A, Kirichek O, Balkwill A, *et al.* Comparison of dementia recorded in routinely  
18 collected hospital admission data in England with dementia recorded in primary care.  
19 *Emerg Themes Epidemiol* 2016;**13**:11.  
20  
21 32 Lea NC, Nicholls J, Dobbs C, *et al.* Data Safe Havens and Trust: Toward a Common  
22 Understanding of Trusted Research Platforms for Governing Secure and Ethical Health  
23 Research. *JMIR Medical Informatics.* 2016;**4**:e22. doi:10.2196/medinform.5571  
24  
25 33 Certificate Client Directory search results. [https://www.bsigroup.com/en-GB/our-](https://www.bsigroup.com/en-GB/our-services/certification/certificate-and-client-directory/search-results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bLondondondon&licencenumber=IS%20612909)  
26 [services/certification/certificate-and-client-directory/search-](https://www.bsigroup.com/en-GB/our-services/certification/certificate-and-client-directory/search-results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bLondondondon&licencenumber=IS%20612909)  
27 [results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bL](https://www.bsigroup.com/en-GB/our-services/certification/certificate-and-client-directory/search-results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bLondondondon&licencenumber=IS%20612909)  
28 [ondon&licencenumber=IS%20612909](https://www.bsigroup.com/en-GB/our-services/certification/certificate-and-client-directory/search-results/?searchkey=licence%3dIS%2b612909%26company%3duniversity%2bcollege%2bLondondondon&licencenumber=IS%20612909) (accessed 26 Jul 2021).  
29  
30  
31 34 Cox DR. Regression Models and Life-Tables. *Journal of the Royal Statistical Society:*  
32 *Series B (Methodological).* 1972;**34**:187–202. doi:10.1111/j.2517-6161.1972.tb00899.x  
33  
34 35 Fine JP, Gray RJ. A Proportional Hazards Model for the Subdistribution of a Competing  
35 Risk. *Journal of the American Statistical Association.* 1999;**94**:496–509.  
36 doi:10.1080/01621459.1999.10474144  
37  
38 36 Mutlu U, Colijn JM, Ikram MA, *et al.* Association of Retinal Neurodegeneration on Optical  
39 Coherence Tomography With Dementia: A Population-Based Study. *JAMA Neurol*  
40 2018;**75**:1256–63.  
41  
42 37 Harper G. Linkage of Maternity Hospital Episode Statistics data to birth registration and  
43 notification records for births in England 2005–2014: Quality assurance of linkage of routine  
44 data for singleton and multiple births. *BMJ Open.* 2018;**8**:e017898. doi:10.1136/bmjopen-  
45 2017-017898  
46  
47 38 Chua SYL, Thomas D, Allen N, *et al.* Cohort profile: design and methods in the eye and  
48 vision consortium of UK Biobank. *BMJ Open* 2019;**9**:e025077.  
49  
50 39 Mutlu U, Bonnemaier PWM, Ikram MA, *et al.* Retinal neurodegeneration and brain MRI  
51 markers: the Rotterdam Study. *Neurobiol Aging* 2017;**60**:183–91.  
52  
53 40 Ikram MA, Brusselle GGO, Murad SD, *et al.* The Rotterdam Study: 2018 update on  
54 objectives, design and main results. *Eur J Epidemiol* 2017;**32**:807–50.  
55  
56 41 Majithia S, Tham Y-C, Chee M-L, *et al.* Cohort Profile: The Singapore Epidemiology of Eye  
57  
58  
59  
60

- 1  
2  
3 Diseases study (SEED). *Int J Epidemiol* 2021;**50**:41–52.  
4
- 5 42 Khan SM, Liu X, Nath S, *et al*. A global review of publicly available datasets for  
6 ophthalmological imaging: barriers to access, usability, and generalisability. *Lancet Digit*  
7 *Health* 2021;**3**:e51–66.  
8
- 9 43 Age groups. 2018. <https://www.ethnicity-facts-figures.service.gov.uk/uk-population-by-ethnicity/demographics/age-groups/latest#:~:text=by%20ethnicity%20Summary-,The%20data%20shows%20that%3A,aged%2060%20years%20and%20over> (accessed 15  
10 Jul 2021).  
11  
12  
13
- 14 44 Enoch J, McDonald L, Jones L, *et al*. Evaluating Whether Sight Is the Most Valued Sense.  
15 *JAMA Ophthalmol* 2019;**137**:1317–20.  
16
- 17 45 Sabanayagam C, Xu D, Ting DSW, *et al*. A deep learning algorithm to detect chronic kidney  
18 disease from retinal photographs in community-based populations. *Lancet Digit Health*  
19 2020;**2**:e295–302.  
20
- 21 46 Mitani A, Huang A, Venugopalan S, *et al*. Detection of anaemia from retinal fundus images  
22 via deep learning. *Nat Biomed Eng* 2020;**4**:18–27.  
23
- 24 47 Poplin R, Varadarajan AV, Blumer K, *et al*. Prediction of cardiovascular risk factors from  
25 retinal fundus photographs via deep learning. *Nat Biomed Eng* 2018;**2**:158–64.  
26
- 27 48 Wisely CE, Wang D, Henao R, *et al*. Convolutional neural network to identify symptomatic  
28 Alzheimer's disease using multimodal retinal imaging. *Br J Ophthalmol* Published Online  
29 First: 26 November 2020. doi:10.1136/bjophthalmol-2020-317659  
30
- 31 49 Cheung CY, Xu D, Cheng C-Y, *et al*. A deep-learning system for the assessment of  
32 cardiovascular disease risk via the measurement of retinal-vessel calibre. *Nat Biomed Eng*  
33 2021;**5**:498–508.  
34
- 35 50 Xiao W, Huang X, Wang JH, *et al*. Screening and identifying hepatobiliary diseases through  
36 deep learning using ocular images: a prospective, multicentre study. *Lancet Digit Health*  
37 2021;**3**:e88–97.  
38
- 39 51 Wagner SK, Fu DJ, Faes L, *et al*. Insights into Systemic Disease through Retinal Imaging-  
40 Based Oculomics. *Transl Vis Sci Technol* 2020;**9**:6.  
41
- 42 52 Deng J, Dong W, Socher R, *et al*. ImageNet: A large-scale hierarchical image database.  
43 2009 IEEE Conference on Computer Vision and Pattern Recognition. 2009.  
44 doi:10.1109/cvprw.2009.5206848  
45
- 46 53 Harron K, Dibben C, Boyd J, *et al*. Challenges in administrative data linkage for research.  
47 *Big Data & Society*. 2017;**4**:205395171774567. doi:10.1177/2053951717745678  
48  
49
- 50 54 Data protection by design and default. Published Online First: 9 February  
51 2021. [https://ico.org.uk/for-organisations/guide-to-data-protection/guide-to-the-general-data-  
52 protection-regulation-gdpr/accountability-and-governance/data-protection-by-design-and-  
53 default/](https://ico.org.uk/for-organisations/guide-to-data-protection/guide-to-the-general-data-protection-regulation-gdpr/accountability-and-governance/data-protection-by-design-and-default/) (accessed 13 Jul 2021).  
54
- 55 55 Lee AY, Campbell JP, Hwang TS, *et al*. Recommendations for Standardization of Images in  
56 Ophthalmology. *Ophthalmology* 2021;**128**:969–70.  
57  
58  
59

- 1  
2  
3 56 NHS and social care data: off-shoring and the use of public cloud services.  
4 [https://digital.nhs.uk/data-and-information/looking-after-information/data-security-and-](https://digital.nhs.uk/data-and-information/looking-after-information/data-security-and-information-governance/nhs-and-social-care-data-off-shoring-and-the-use-of-public-cloud-services)  
5 [information-governance/nhs-and-social-care-data-off-shoring-and-the-use-of-public-cloud-](https://digital.nhs.uk/data-and-information/looking-after-information/data-security-and-information-governance/nhs-and-social-care-data-off-shoring-and-the-use-of-public-cloud-services)  
6 [services](https://digital.nhs.uk/data-and-information/looking-after-information/data-security-and-information-governance/nhs-and-social-care-data-off-shoring-and-the-use-of-public-cloud-services) (accessed 13 Jul 2021).  
7
- 8 57 Sinha S, Peach G, Poloniecki JD, *et al*. Studies using English administrative data (Hospital  
9 Episode Statistics) to assess health-care outcomes—systematic review and  
10 recommendations for reporting. *European Journal of Public Health*. 2013;**23**:86–92.  
11 doi:10.1093/eurpub/cks046  
12
- 13 58 Burns EM, Rigby E, Mamidanna R, *et al*. Systematic review of discharge coding accuracy.  
14 *J Public Health* 2012;**34**:138–48.  
15
- 16 59 Herrett E, Shah AD, Boggon R, *et al*. Completeness and diagnostic validity of recording  
17 acute myocardial infarction events in primary care, hospital care, disease registry, and  
18 national mortality records: cohort study. *BMJ* 2013;**346**:f2350.  
19
- 20 60 Wood A, Denholm R, Hollings S, *et al*. Linked electronic health records for research on a  
21 nationwide cohort of more than 54 million people in England: data resource. *BMJ*  
22 2021;**373**:n826.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

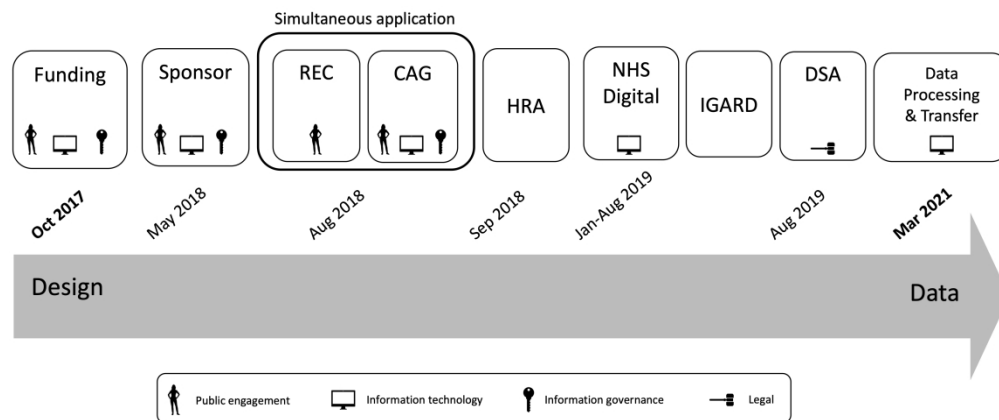


Figure 1: Schematic of the key milestones, prerequisites and approvals with their corresponding achievement dates for the AlzEye dataset. REC: Research Ethics Committee, CAG: Confidential Advisory Group, HRA: Health Research Authority, NHS: National Health Service, IGARD: Independent Group Advising on the Release of Data, DSA: Data Sharing Agreement.

1327x627mm (72 x 72 DPI)

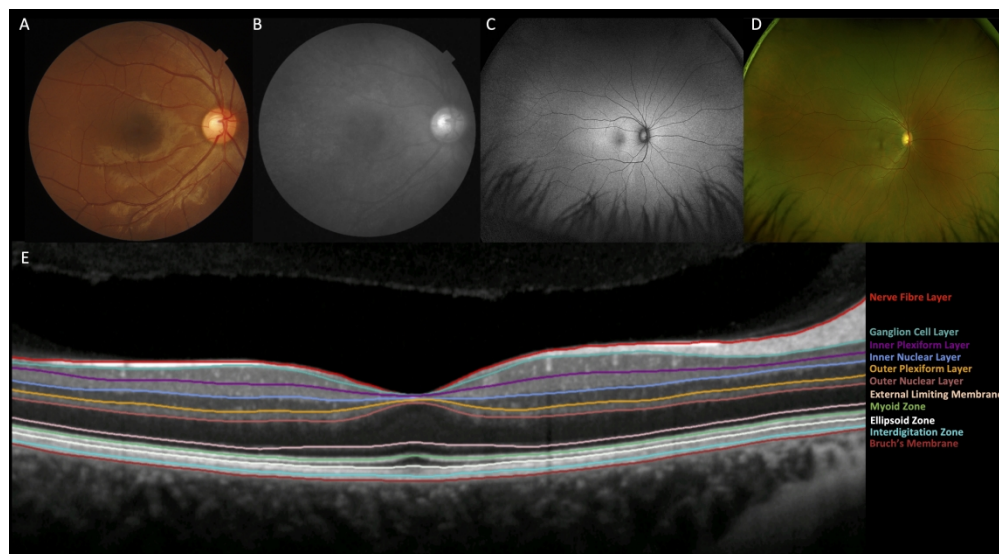


Figure 2: Composite figure showing the major retinal imaging modalities within AlzEye. A: Colour fundus photograph, B: Red-free photograph, C: Fundus autofluorescence (widefield), D: Pseudocolour photography (widefield) and E: Optical coherence tomography of the central macula illustrating segmentation of the individual sublayers. Consensus nomenclature for the retinal sublayers is indicated. NFL= nerve fibre layer, GCL = ganglion cell layer, IPL: inner plexiform layer, INL: inner nuclear layer, OPL: outer plexiform layer, ONL = outer nuclear layer, MZ = myoid zone, RPE-IZ = RPE-interdigitation zone, BM = Bruch's membrane.

1456x803mm (72 x 72 DPI)

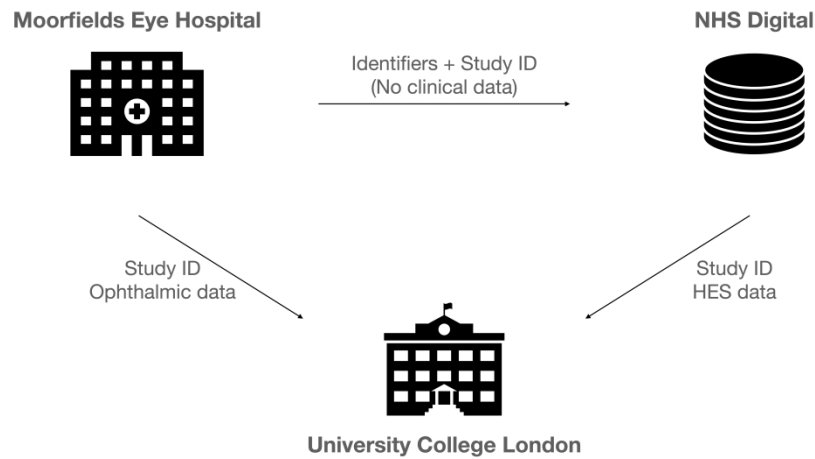


Figure 3: Linkage approach of AlzEye. Moorfields Eye Hospital NHS Foundation Trust securely transfers a spreadsheet of identifiers with a study ID to NHS Digital and separately transfers the study ID with ophthalmic data, including diagnoses and retinal images, to University College London. NHS Digital links the identifiers with the Hospital Episode Statistics (HES) database and returns the admissions data with the study ID (and no identifiable data) to UCL. UCL links the ophthalmic data from Moorfields Eye Hospital with HES data from NHS Digital using the study ID.

2822x1587mm (72 x 72 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

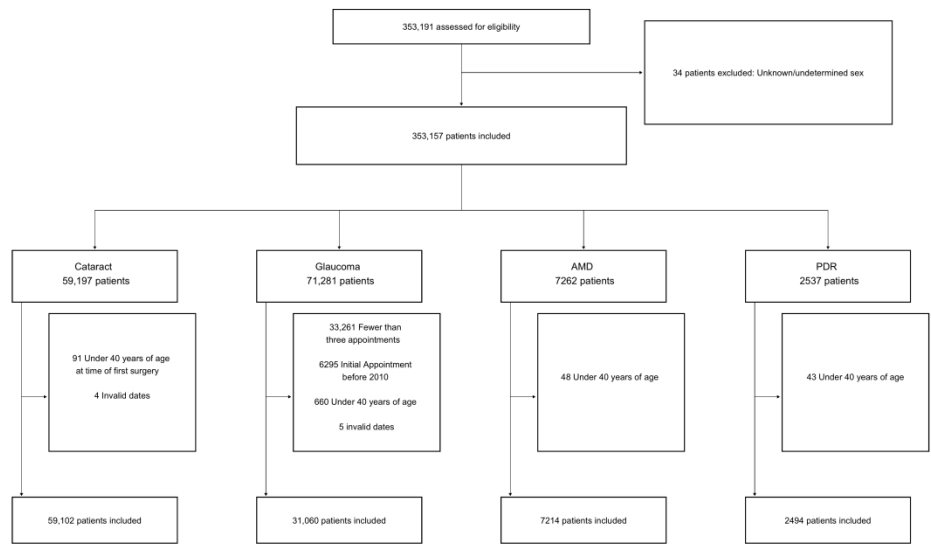


Figure 4: CONSORT style flow chart illustrating the distribution of cataract, glaucoma, neovascular age-related macular degeneration (AMD) and proliferative diabetic retinopathy (PDR) within the AlzEye dataset.

1693x1058mm (72 x 72 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

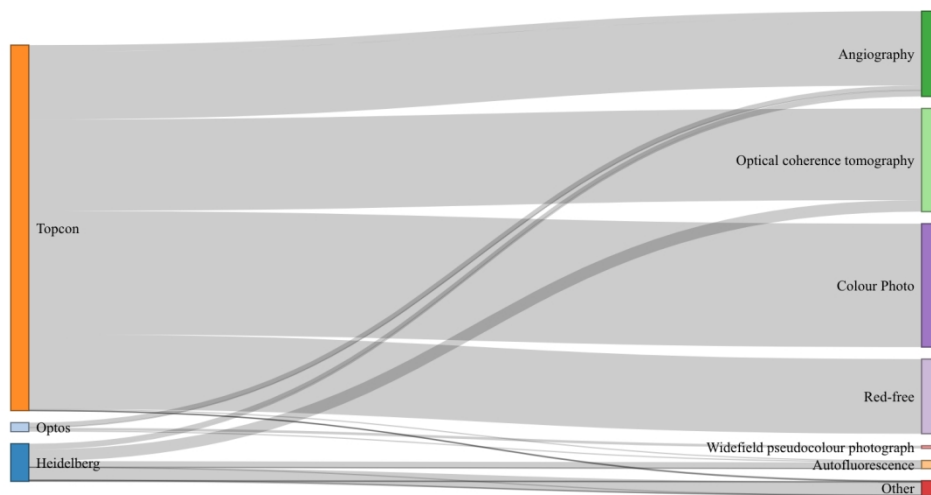


Figure 5: Parallel sets diagram illustrating the imaging modality across vendors within AlzEye. The majority of images were acquired on the Topcon system and the most frequent modalities were colour photography and optical coherence tomography. Designed using the networkD3 package.

584x324mm (72 x 72 DPI)



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

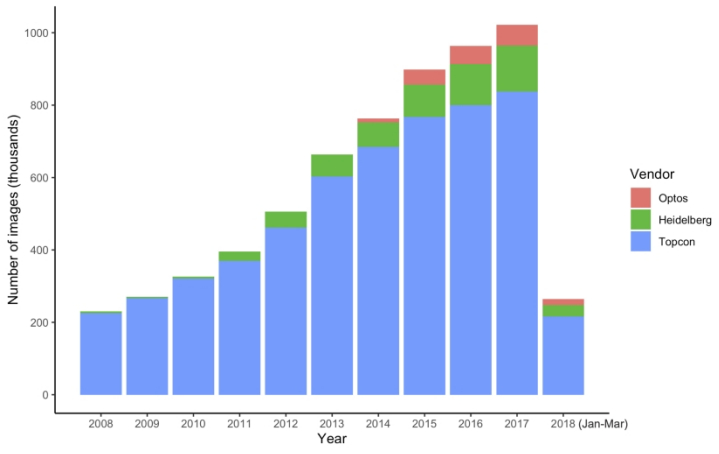


Figure 6: Stacked bar chart of the annual number of images acquired during the study period for the three leading device vendors at Moorfields Eye Hospital. Data for 2018 represents 3 months only prior to the study end-date.

1411x793mm (72 x 72 DPI)

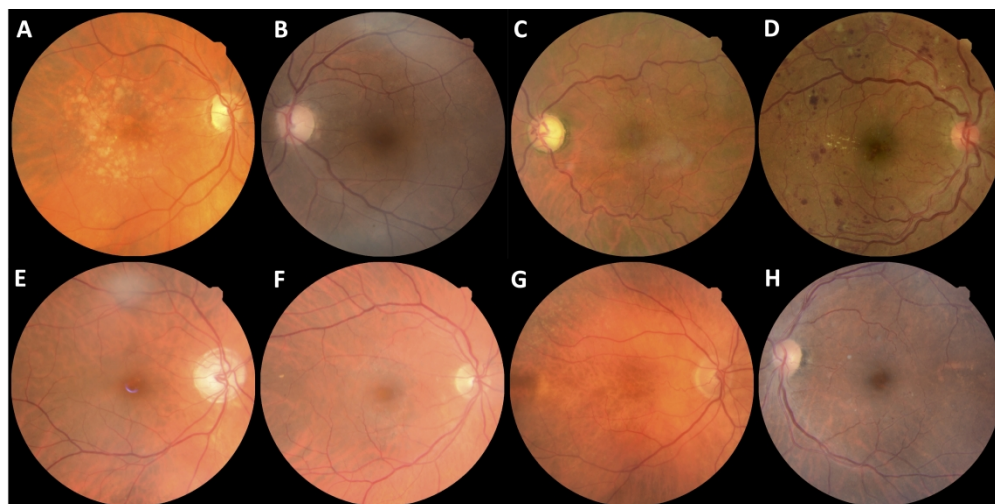


Figure 7: Example colour retinal photographs of ophthalmic and systemic diseases within AlzEye. A: age-related macular degeneration, B: cataract, C: glaucoma, D: proliferative diabetic retinopathy E: prevalent Alzheimer's disease, F: incident ischaemic stroke, G: incident myocardial infarction, H: prevalent vascular dementia.

1800x899mm (72 x 72 DPI)