# Heliyon

# Prediction of atrial fibrillation and stroke using machine learning models in UK Biobank
## --Manuscript Draft--

| | |
|---|---|
| **Manuscript Number:** | HELIYON-D-23-41056R2 |
| **Article Type:** | Original Research Article |
| **Section/Category:** | Medical Sciences |
| **Keywords:** | Atrial Fibrillation;  stroke;  Risk prediction;  Machine Learning |
| **Manuscript Classifications:** | 10.130: Statistics; 20.140.100: Artificial Intelligence; 20.140.100.140: Machine Learning; 110.170: Bioinformatics; 110.170.100: Biocomputational Method; 110.210: Genetics; 130.100.130: Epidemiology; 130.110: Cardiology; 110.430: Biostatistics; 130.100.160: Chronic Diseases |
| **Corresponding Author:** | Panos Deloukas, PhD<br>Queen Mary University of London<br>London, UNITED KINGDOM |
| **First Author:** | Areti Papadopoulou |
| **Order of Authors:** | Areti Papadopoulou |
| | Daniel Harding |
| | Greg Slabaugh |
| | Eirini Marouli |
| | Panos Deloukas |

| | |
|---|---|
| **Abstract:** | Background: Atrial fibrillation (AF) is the most common cardiac arrythmia; 12.1 million people are expected to be affected by 2030. Importantly, AF is associated with increased risk for ischemic stroke, which is underestimated as AF can be asymptomatic. Methods: To develop ML models for prediction of 1) AF in the general population and 2) ischemic stroke in patients with AF we constructed XGBoost, LightGBM, Random Forest, Deep Neural Network, Support Vector Machine and Lasso penalised logistic regression models using UK-Biobank's extensive real-world clinical data, questionnaires, as well as biochemical and genetic data, and their predictive performances were compared. Ranking and contribution of the different features was assessed by SHapley Additive exPlanations (SHAP) analysis. The clinical tool CHA2DS2-VASc for prediction of ischemic stroke among AF patients, was used for comparison to the best performing ML model. Findings: The best performing model for AF prediction was LightGBM, with an area-under-the-roc-curve (AUROC) of 0.729 (95% confidence intervals (CI): 0.719, 0.738). The best performing model for ischemic stroke prediction in AF patients was XGBoost with AUROC of 0.631 (95% CI: 0.604, 0.657). The improved AUROC in the XGBoost model compared to CHA2DS2-VASc was statistically significant based on DeLong's test (pvalue=2.20E-06). In addition, the SHAP analysis showed that several peripheral blood biomarkers (e.g. creatinine, glycated haemoglobin, monocytes) were associated with ischemic stroke, which are not considered by CHA2DS2-VASc. Low levels of albumin and increased levels of alkaline phosphatase were associated with increased risk of ischemic stroke also in European descent subjects and not only in East Asians as previously reported. Interpretation: The best performing ML models presented have the potential for clinical use, but further validation in independent studies is required. Our results endorse the incorporation of some routinely measured blood biomarkers for ischemic stroke prediction in AF patients. Funding: This work was funded from the National Institute of Health Research (NIHR) Barts Biomedical Research Centre. |

| | |
|---|---|
| **Opposed Reviewers:** | |
| **Additional Information:** | |
| **Question** | **Response** |
| **Publication ethics** | I confirm |

Please confirm that you have reviewed our guidelines for Ethics in Publishing as well as Heliyon's Ethics and Editorial Policies

CellPress

# CELL PRESS DECLARATION OF INTERESTS POLICY

Transparency is essential for a reader's trust in the scientific process and for the credibility of published articles. At Cell Press, we feel that disclosure of competing interests is a critical aspect of transparency. Therefore, we require a "declaration of interests" section in which all authors disclose any financial or other interests related to the submitted work that (1) could affect or have the perception of affecting the author's objectivity or (2) could influence or have the perception of influencing the content of the article.

### *What types of articles does this apply to?*
We require that you disclose competing interests for all submitted content by completing and submitting the form below. We also require that you include a "declaration of interests" section in the text of all articles even if there are no interests to declare.

### *What should I disclose?*
We require that you and all authors disclose any personal financial interests (e.g., stocks or shares in companies with interests related to the submitted work or consulting fees from companies that could have interests related to the work), professional affiliations, advisory positions, board memberships (including membership on a journal's advisory board when publishing in that journal), or patent holdings that are related to the subject matter of the contribution. As a guideline, you need to declare an interest for (1) any affiliation associated with a payment or financial benefit exceeding $10,000 p.a. or 5% ownership of a company or (2) research funding by a company with related interests. You do not need to disclose diversified mutual funds, 401ks, or investment trusts.

Authors should also disclose relevant financial interests of immediate family members. Cell Press uses the Public Health Service definition of "immediate family member," which includes spouse and dependent children.

### *Where do I declare competing interests?*
Competing interests should be disclosed on this form as well as in a "declaration of interests" section in the manuscript. This section should include financial or other competing interests as well as affiliations that are not included in the author list. Examples of "declaration of interests" language include:

> "AUTHOR is an employee and shareholder of COMPANY."
> "AUTHOR is a founder of COMPANY and a member of its scientific advisory board."

*NOTE*: Primary affiliations should be included with the author list and do not need to be included in the "declaration of interests" section. Funding sources should be included in the "acknowledgments" section and also do not need to be included in the "declaration of interests" section. (A small number of front-matter article types do not include an "acknowledgments" section. For these articles, reporting of funding sources is not required.)

### *What if there are no competing interests to declare?*
If you have no competing interests to declare, please note that in the "declaration of interests" section with the following wording:

> "The authors declare no competing interests."

*July 14, 2023*

![CellPress logo]

# CELL PRESS DECLARATION OF INTERESTS FORM

If submitting materials via Editorial Manager, please complete this form and upload with your initial submission. Otherwise, please email as an attachment to the editor handling your manuscript.

***Please complete each section of the form and insert any necessary "declaration of interests" statement in the text box at the end of the form. A matching statement should be included in a "declaration of interests" section in the manuscript.***

### *Institutional affiliations*
We require that you list the current institutional affiliations of all authors, including academic, corporate, and industrial, on the title page of the manuscript. ***Please select one of the following:***

☒ All affiliations are listed on the title page of the manuscript.

☐ I or other authors have additional affiliations that we have noted in the "declaration of interests" section of the manuscript and on this form below.

### *Funding sources*
We require that you disclose all funding sources for the research described in this work. ***Please confirm the following:***

☒ All funding sources for this study are listed in the "acknowledgments" section of the manuscript.*

*A small number of front-matter article types do not include an "acknowledgments" section. For these, reporting funding sources is not required.

### *Competing financial interests*
We require that authors disclose any financial interests and any such interests of immediate family members, including financial holdings, professional affiliations, advisory positions, board memberships, receipt of consulting fees, etc., that:

(1) could affect or have the perception of affecting the author's objectivity, *or*
(2) could influence or have the perception of influencing the content of the article.

***Please select one of the following:***

☒ We, the authors and our immediate family members, have no financial interests to declare.

☐ We, the authors, have noted any financial interests in the "declaration of interests" section of the manuscript and on this form below, and we have noted interests of our immediate family members.

*July 14, 2023*

_**Advisory/management and consulting positions**_
We require that authors disclose any position, be it a member of a board or advisory committee or a paid consultant, that they have been involved with that is related to this study. We also require that members of our journal advisory boards disclose their position when publishing in that journal. **_Please select one of the following:_**

☒ We, the authors and our immediate family members, have no positions to declare and are not members of the journal's advisory board.

☐ The authors and/or their immediate family members have management/advisory or consulting relationships noted in the "declaration of interests" section of the manuscript and on this form below.

_**Patents**_
We require that you disclose any patents related to this work by any of the authors or their institutions. **_Please select one of the following:_**

☒ We, the authors and our immediate family members, have no related patents to declare.

☐ We, the authors, have a patent related to this work, which is noted in the "declaration of interests" section of the manuscript and on this form below, and we have noted the patents of immediate family members.

**_Please insert any "declaration of interests" statements in this space._** This exact text should also be included in the "declaration of interests" section of the manuscript. If no authors have a competing interest, please insert the text, "The authors declare no competing interests."

The authors declare no competing interests

☒ **On behalf of all authors, I declare that I have disclosed all competing interests related to this work. If any exist, they have been included in the "declaration of interests" section of the manuscript.**

Manuscript number
(if available):

ISCIENCE-D-23-01236

_July 14, 2023_

**Declaration of interests**

☒The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

The declaration of competing interests that you upload should be in a standard and editable format.  Please select the suitable option and upload it with this submission. You can download the standard declaration of competing interests form from the following link.
https://declarations.elsevier.com/

This has now been addressed.

Please reference all numbered figures in text. Currently, numbered figures [1, 4] in the manuscript have not been cited in the text.

This was implemented in lines 104, 136, 39, 243, and 262

Please reference all numbered tables in text. Currently, numbered tables [1] in the manuscript have not been cited in text.

Table 1 is cited in line 216.

If not already included, please add a 'Data availability statement' in your manuscript along with the name of the repository and the accession number if applicable. The statement should be located right before the 'References' section.

addressed in lines 391-393.

Furthermore, please ensure that the references in your manuscript are in a numbered format, if not already. To avoid unnecessary delays in the publication of your manuscript, please do not make any additional changes apart from those requested during this revision.

**Prediction of atrial fibrillation and stroke using machine learning models in UK Biobank.**

A. Papadopoulou[1], D. Harding[1], G. Slabaugh[2,3], E. Marouli[1,3,5], P. Deloukas[1,4,5]

1 William Harvey Research Institute, Barts and The London School of Medicine and Dentistry, Queen Mary University of London, London, UK.
2 School of Electronic Engineering and Computer Science, Queen Mary University of London, London, UK.
3 Digital Environment Research Institute, Queen Mary University of London, London, UK.
4 Princess Al-Jawhara Al-Brahim Centre of Excellence in Research of Hereditary Disorders (PACER-HD), King Abdulaziz University, Jeddah, Saudi Arabia
5 These authors contributed equally to this work; corresponding authors

## Abstract

**Objective:** Atrial fibrillation (AF) is the most common cardiac arrythmia, and it is associated with increased risk for ischemic stroke, which is underestimated, as AF can be asymptomatic. The aim of this study was to develop optimal ML models for prediction of AF in the population, and secondly for ischemic stroke in AF patients.

**Methods:** To develop ML models for prediction of 1) AF in the general population and 2) ischemic stroke in patients with AF we constructed XGBoost, LightGBM, Random Forest, Deep Neural Network, Support Vector Machine and Lasso penalised logistic regression models using UK-Biobank's extensive real-world clinical data, questionnaires, as well as biochemical and genetic data, and their predictive performances were compared. Ranking and contribution of the different features was assessed by SHapley Additive exPlanations (SHAP) analysis. The clinical tool $CHA_2DS_2$-VASc for prediction of ischemic stroke among AF patients, was used for comparison to the best performing ML model.

**Findings:** The best performing model for AF prediction was LightGBM, with an area-under-the-roc-curve (AUROC) of 0.729 (95% confidence intervals (CI): 0.719, 0.738). The best performing model for ischemic stroke prediction in AF patients was XGBoost with AUROC of 0.631 (95% CI: 0.604, 0.657). The improved AUROC in the XGBoost model compared to $CHA_2DS_2$-VASc was statistically significant based on DeLong's test (pvalue=2.20E-06). In addition, the SHAP analysis showed that several peripheral blood biomarkers (e.g. creatinine, glycated haemoglobin, monocytes) were associated with ischemic stroke, which are not considered by $CHA_2DS_2$-VASc.

**Implications:** The best performing ML models presented have the potential for clinical use, but further validation in independent studies is required. Our results endorse the incorporation of some routinely measured blood biomarkers for ischemic stroke prediction in AF patients.

## Introduction

Atrial fibrillation (AF) is the most common cardiac arrythmia, which is characterised by a rapid and irregular heartbeat [1, 2]. The incidence of AF is increasing rapidly with 12.1 million people expected to be affected by 2030. This is mainly attributed to the ageing of the population, along with changes in lifestyle. AF, besides doubling the risk of cardiovascular mortality, is associated with increased risk of stroke, ischemic heart disease, heart failure and cognitive dysfunction. More specifically, AF quintuple the risk for ischemic stroke, independent of age. However, AF is sometimes asymptomatic, and thus remains undetected, and subsequently the ischemic stroke risk attributed to AF is under-estimated [1, 2].

Machine learning (ML) algorithms are promising to revolutionise disease prediction, classification of medical images and diagnosis revealing new features, which would have not been discovered using traditional statistical models [3]. ML models use a hypothesis-free approach with no prior assumptions either among the input features or between the features and the outcome. ML methods with varying degree of accuracy have been reported for the prediction of circulatory diseases. However, they have been limited from access to large-scale cohorts with integrated clinical, biochemical and genetic data [3, 4].

There have been several studies that employed ML methods for prediction of circulatory diseases. A recent study in Geisinger's clinical MUSE database with no history of AF, within 1-year of an ECG, employed deep neural networks and reported an area under the receiver operating characteristic (AUROC) of 0.85 for AF prediction [3]. They also reported that 62% of patients who had a stroke caused by AF within 3 years of an ECG, with no prior AF diagnosis, would have been identified by their prediction tool before the stroke occurred [3]. Another study employed four ML models to predict modified Rankin Scale (mRS) at hospital discharge and in-hospital deterioration for acute ischemic stroke patients enrolled on the Stroke Registry in Chang Gung Healthcare System (SRICHS) [4]. Random forest performed well in both outcomes; the AUROC was 0.83 for discharge mRS and 0.71 for in-hospital deterioration [4]. There have also been several studies using ML methods for the prediction of ischemic stroke in AF-patients. In the Korean National Health Insurance (KNHIS) dataset, the authors aimed to predict ischemic stroke occurrence in AF patients using ML models such as DNN, XGBoost and RF, for more than 150,000 AF patients. The best performing model was DNN with an AUROC of 0.727, outperforming $CHA_2DS_2$-VASc with AUROC of 0.651 [5]. Another study using the Fushimi AF registry, showed that CatBoost ML method outperformed $CHA_2DS_2$-VASc, having AUROC 0.72 (95%CI, 0.66-0.79) and 0.62 (95%CI, 0.54-0.70) respectively [6]. Using the Korean Atrial Fibrillation Evaluation Registry in Ischemic Stroke Patients (K-ATTENTION), the authors showed that LightGBM performed the best, with AUROC of 0.772 (95% CI 0.715-0.829), for the prediction of early neurological deterioration (END) among AF-related stroke patients [7]. The studies mentioned above underlined the importance of ML methods, since besides the improved prediction performance that they display in contrast to current clinical tools, they exhibit the potential to unravel new and diverse risk factors associated with the disease.

The aim of this study was to develop optimal ML models for prediction of: 1) AF in the population and 2) ischemic stroke in AF patients. We constructed ML models with six different algorithms in UK-Biobank (500,000 participants with extensive questionnaires, clinical, biochemical and genetic data – Tables S1-S3) and assessed their predictive performances. For ranking of feature importance and contribution to the prediction outcome we used SHapley Additive exPlanations (SHAP) [8].

**Methods**

*Overview of the research framework*

We included clinical data, phenotypes, lifestyle, and medications from UK-Biobank. We imputed the missing values and employed a feature selection process, described in more detail at *Data pre-processing*, to reduce the number of features employed to the ones relative to the outcome. Six ML models were used to create predictive models as described at the *ML methods* below. Each model's hyperparameters were optimised using 10-fold cross validation at the training dataset. The ML models were validated on the test dataset and their performances were compared. Lastly, we employed the SHAP explanations to reveal the features' contributions to the prediction.

*Phenotype and participant selection*

*Data pre-processing*

We examined the UK-Biobank, a prospective cohort of 502,492 participants, aged 37-73 years old, recruited between 2006 and 2010. The dataset includes blood measurements, clinical assessments, anthropometry, cognitive function, hearing, arterial stiffness, hand grip strength, sociodemographic factors, lifestyle, family history, psychosocial factors and dietary intake [9]. Related individuals were removed, and the remaining dataset for analysis included 454,118 participants. Furthermore, we incorporated medications as features, derived from field 20003 (https://biobank.ndph.ox.ac.uk/showcase/field.cgi?id=20003). Additionally, clinical data were employed, coded in ICD10, derived from the Hospital Episodes Statistics (HES), which are linked to the UK-Biobank. From these, we constructed phenotype codes or "phecodes", using a phecode map [10], which are aggregated ICD10 codes defining specific diseases or traits. We employed only the umbrella phecode categories. Detailed list of all the features that we examined can be found at *Table_S1, Table_S2, Table_S3*. Moreover, we created two polygenic scores (PGS) which were included as features for the prediction of ischemic stroke in people with AF. The first one is the AF score, based on 94 genome-wide variants derived from the Roseli et al. [11] genome-wide association study (GWAS) for AF. The second is the Ischemic STROKE score, based on 28 genome-wide variants derived from the Malik et al. [12] GWAS for ischemic stroke. The AF SCORE was also employed as a feature both for the prediction of AF and for the ischemic stroke in AF patients.

The investigator phenotypes dataset from UK-Biobank includes 2,199 fields for 454,118 participants. We set answers "Do not know" and "Prefer not to answer" as NA and removed features that had more than 25% missingness, resulting in 390 investigator phenotypes. Afterwards, we imputed the missing values using a multivariate imputer that estimates each feature from all the others, using *IterativeImputer* from Python [13]. Then, we added 419 phecodes, available for 278,177 participants, derived from HES in UK-Biobank. Lastly, we added the medications from field 20003 (https://biobank.ndph.ox.ac.uk/showcase/field.cgi?id=20003), after applying one-hot-encoding, resulting in 1,289 medications for 294,698 participants (Figure 1).
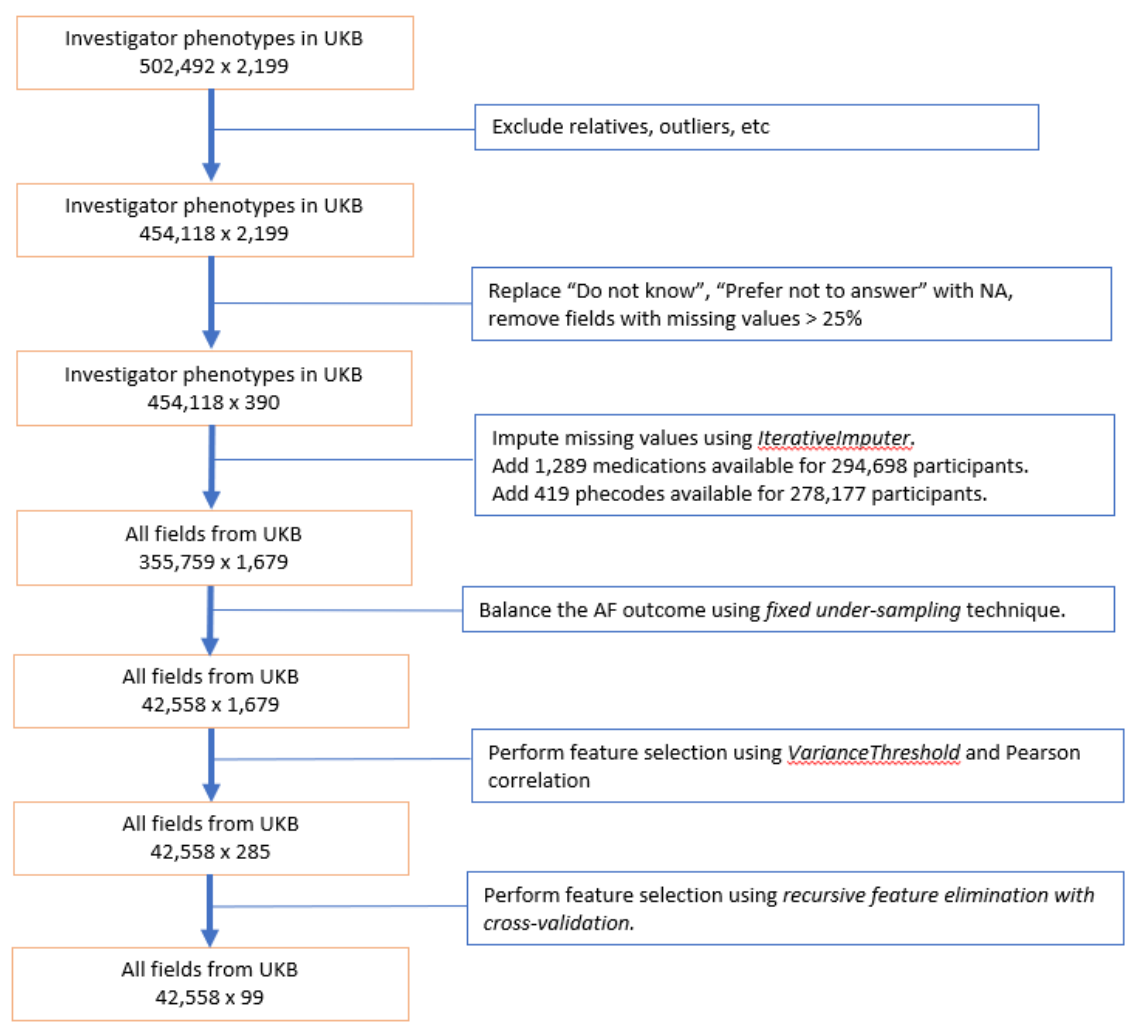
Next, we decided to balance the outcome sample size, since imbalanced data has a negative impact on ML procedures [14]. The classification algorithms have the tendency to get biased estimates towards the majority class, ignoring the minority class. This happens because most of the classifying methods aim to maximize the accuracy rate, meaning the number of correctly classified observations [15, 16]. Therefore, we employed the *fixed under-sampling* technique from Python [17], which is a process for reducing the number of samples in the majority class; the control group in this case. The algorithm randomly selects samples from the control group, in order to have equal representation of both classes. After balancing the outcome, we used *VarianceThreshold* from Python [13], which eliminates all features whose variance does not meet a threshold of 90%. Additionally, we removed the continuous correlated fields using Pearson correlation, at a 0.8 threshold; features strongly correlated with the outcome were maintained. Next, we performed feature selection in order to reduce the computational cost via dimensionality reduction, achieve higher classification accuracy by eliminating the noise, and include the most relevant features for the disease prediction [18]. A recent paper by Ramos-Pérez et al. [19], suggests that the best practice is for the fixed under-sampling technique to precede the feature selection. Therefore, we filtered all the remaining features using recursive feature elimination with cross-validation from Python [13] in order to find the optimal number of features to include in the ML models.

### ***Create the AF outcome***

We removed participants from the UK-Biobank that had cardiac dysrhythmias before the time of enrolment, with one or more of the following codes: non-cancer illness code, self-reported (1471, 1483); operation code (1524); diagnoses – main/secondary ICD10 (I44, I44.1-I44.7, I45, I45.0-I45.6, I45.8-I45.9, I46, I46.2, I46.8-I46.9, I47, I47.0-I47.2, I47.9, I48, I48.0-4, I48.9, I49, I49.0-I49.5, I49.8-I49.9, R00.0, R00.1, R00.2, R94.3, Z86.7, Z95.0, Z95.8-Z95.9); underlying (primary/secondary) cause of death: ICD10 (I44, I44.1-I44.7, I45, I45.0-I45.6, I45.8-I45.9, I46, I46.2, I46.8-I46.9, I47, I47.0-I47.2, I47.9, I48, I48.0-4, I48.9, I49, I49.0-I49.5, I49.8-I49.9, I60-I61, I63-I64 (NOT I63.6), R00.0, R00.1, R00.2, R94.3, Z86.7, Z95.0, Z95.8-Z95.9); diagnoses – main/secondary ICD9

(4273, 430, 431, 4339, 4340, 4341, 4349, 436); operative procedures – main/secondary OPCS (K57.1, K62.1-4). In total, 20,584 participants were excluded, having at least one of the above conditions, before enrolment in the UK-Biobank.

AF cases were defined when having one or more of the following codes: non-cancer illness code, self-reported (1471, 1483); operation code (1524); diagnoses – main/secondary ICD10 (I48, I48.0-4, I48.9); underlying (primary/secondary) cause of death: ICD10 (I48, I48.0-4, I48.9); operative procedures – main/secondary OPCS (K57.1, K62.1-4). In total, 21,279 people developed one of the conditions described above, after enrolment in UK-Biobank (Figure 1).



*Figure 1: Diagram depicting the data curation and feature selection process for the prediction of atrial fibrillation.*

### *Create the AF & Stroke outcome*

Cases were defined as participants who developed ischemic stroke after AF diagnosis in UK-Biobank with one or more of the following codes: diagnoses – main/secondary ICD10 (I63, I63.0-9, I64); diagnoses – main/secondary ICD9 (434, 436); underlying (primary/secondary) cause of death: ICD10 (I63, I63.0-9, I64). Thus, 3,150 people developed ischemic stroke after AF diagnosis and were included as cases, and the controls were people diagnosed with AF and did not develop stroke, as far as the data allow us to know. Based on the selection criteria for AF patients with and without ischemic stroke (Supplementary figure 1), 3,150 prospective

147 cases who developed ischemic stroke after AF diagnosis and equal number of controls, along with 129
148 features, were included in the ML models *(Table_S8)*.

149 *ML models*

150 **XGBoost**
151 In more detail, XGBoost uses regression trees in a sequential learning process as weak learners into a single
152 strong model, where each tree attempts to correct the residuals in the predictions made by previous trees.
153 Regression trees include a continuous score on each leaf, which is the last node once the tree has grown. For
154 a specific observation, the algorithm uses decision rules in the trees to classify it into the leaves. The sum of
155 the scores on each leaf is the final prediction [20].

156 **LightGBM**
157 Machine learning methods relying on Gradient Boosting Decision Tree (GBDT) scan all the data instances, for
158 all the features, to calculate the information gain for each possible split. As a result, the computational time
159 and complexity will increase as the features accumulate. To this end, there are two techniques incorporated
160 at LightGBM algorithm that contribute towards a faster implementation. Firstly, in the Gradient-based One-
161 Side Sampling (GOSS) technique, instances that have larger gradients contribute more to the information gain,
162 and the instances with smaller gradients are randomly sampled to provide accurate and fast estimation.
163 Secondly, the Exclusive Feature Bundling (EFB) technique reduces the number of effective features. For
164 datasets that are sparse, many features are mutually exclusive; they will rarely take nonzero values at the
165 same time, therefore such features are tied into one [21].

166 **Deep Neural Networks (DNN)**
167 Deep learning is a subdomain of ML attempting to learn many levels of representation using multiple layers.
168 These layers transform the data in a non-linear way, and as a result, more complex structure and relationships
169 can be discovered. This method is inspired by the human brain, using a series of connected layers of neurons
170 that constitute a whole network, including at least three layers: input, hidden and output. The input layer
171 consists of multiple neurons, which use as input the original features. The hidden layers can be more than one,
172 depending on the complexity of the dataset. Each layer includes multiple nodes, and each node from the
173 previous layer is connected to each one from the next layer, constituting a fully connected network. Lastly,
174 the output layer, using a sigmoid activation function, concludes in a number between 0 and 1, which
175 represents the probability belonging to one of the two classes [22].

176 **Support Vector Machine (SVM)**
177 SVM is a high accuracy ML model, which can deal with non-linear spaces. It maps the input data into a higher
178 dimension feature space, using a kernel function. Then, a linear decision surface (hyperplane), is created to
179 classify the outcome, with properties that satisfy the generalisation of the algorithm. The optimal hyperplane
180 classifies the data by using its maximal margin, employing a small percentage of the training data, which are
181 named support vectors. The authors support that if the optimal hyperplane is created from a few support
182 vectors, then the algorithm can be generalised, even in a space with infinite dimensions [23].

183 Cross-validation

184 The ML models aim to optimise the general model performance on datasets different from the ones used to
185 train them. Therefore, evaluating the generalisation of ML methods requires the data to be split in three non-
186 overlapping sets of training/validation/test, combined with stratified 10-fold cross-validation (CV),
187 maintaining the same proportion of cases and controls in each fold. Grid search is performed using 9 sets for
188 the parameter tuning, and the 1 remaining set is used for validation. This process is repeated 10 times, until
189 every set is used once for training and once for validation. The best parameters for the model correspond to

the highest score, which is calculated by averaging the results from all repetitions. The test dataset is used to check for overfitting and unbiased evaluation of the final model [13].

*SHAP*

ML models, although accurate and capable of capturing the non-linear relationships, are complex to interpret. A more widespread method for interpretation is SHAP, employed to understand each feature's contribution to the prediction, using cooperative game theoretic tools. The SHAP values are in theory the best solution up to now, however time-consuming, since all possible combinations need to be calculated. TreeExplainer is an expansion of SHAP, employing tree nodes instead of linear models for the estimation of Shapley values. The Shapley values of a tree-based algorithm are calculated as the weighted average of the Shapley values corresponding to individual trees. Thus, it is commonly used to explain tree-based machine learning models, reducing tremendously the computation time. In parallel, SHAP values seem to overcome the interpretability issue by employing both global and local interpretation. Global explanation relies on the effect of input features on the whole model, and local interpretation depicts the effect of input features on single predictions [8].

For the methods described above, Python computer language was employed [24]. The code and libraries that were employed are described in Table_S5.

## Results

Machine learning models can enhance prediction accuracy by utilising extensive datasets and incorporating potential predictors. In our present study, we demonstrated the improvement in prediction accuracy for ischemic stroke among AF patients, compared to current approaches, by employing machine learning modelling. The findings suggest inclusion of commonly measured blood biomarkers for prediction, while advocating for the incorporation of a genetic score for AF prediction. The approaches and modelling introduced in this study hold promise for clinical implementations.

*AF*

We examined 21,279 prospective AF cases and an equal number of controls in UK-Biobank. Baseline characteristics, along with comorbidities and medication, both overall and according to AF cases versus controls, are provided in **Error! Reference source not found.**.

*Table 1:Baseline characteristics for the 21,279 prospective AF cases and equal number of controls.*

|  | Total | AF cases | AF controls | Pvalue* |
|---|---|---|---|---|
| **Sex** |  |  |  |  |
| Females | 20231 (47.5%) | 8122 (38.2%) | 12109 (56.9%) | < 2.2E-16 |
| Males | 22327 (52.5%) | 13157 (61.8%) | 9170 (43.1%) |  |
| **Age (mean, sd)** | 59 (8) | 62 (6) | 57 (8) | < 2.2E-16 |
| **Ethnicity** |  |  |  |  |
| EUR | 41042 (96.9%) | 20791 (97.7%) | 20251 (95.0%) | 5E-03 |
| AFR | 535 (1.2%) | 154 (0.7%) | 381 (1.8%) |  |
| EAS | 127 (0.3%) | 31 (0.2%) | 96 (0.5%) |  |
| SAS | 854 (1.6%) | 303 (1.4%) | 551 (2.7%) |  |
| **Comorbidities** |  |  |  |  |
| Diabetes | 6434 (15.1%) | 4423 (20.8%) | 2011 (9.5%) | < 2.2E-16 |
| Hypertension | 22019 (51.7%) | 14810 (69.6%) | 7209 (33.9%) | < 2.2E-16 |
| Smoking |  |  |  |  |
| Never | 23273 (54.7%) | 11627 (54.6%) | 11646 (54.7%) | 0.8804 |

| | | | | |
|---|---|---|---|---|
| Previous | 14791 (34.8%) | 7389 (34.7%) | 7402 (34.8%) | |
| Current | 4494 (10.6%) | 2263 (10.6%) | 2231 (10.5%) | |
| **Cholesterol lowering medication** | 7459 (17.5%) | 3712 (17.4%) | 3747 (17.6%) | 0.4799 |
| **History of heart diseases** | 21102 (49.6%) | 11233 (52.8%) | 9869 (46.4%) | < 2.2E-16 |
| **History of stroke** | 12317 (28.9%) | 6581 (30.9%) | 5736 (26.9%) | < 2.2E-16 |

*Note. * P-values refer to chi-square test for dichotomous variables and to Mann-Whitney test for continuous data with non-parametric distribution.*

In total, 99 features (*Table_S4*) were employed, using five ML models to predict AF. The results presented in this section correspond to the optimal hyperparameters, derived after 10-fold cross-validation from the examined values included in *Table_S6*. SVM did not converge after running 10 days and utilising 16 cores in Queen Mary's Apocrita HPC facility[1].

The best AUROC value was achieved with LightGBM (Table 2) albeit De-Long's test (Table 3) showed that there is no evidence for significant difference in the AUROCs between LightGBM and XGBoost, DNN, or RF. In contrast, DeLong's test showed that there was statistically significant difference in the AUROCs between LightGBM and penalised LR (pvalue=1.38E-02), after considering multiple correction. The AUROC of penalised LR differed from the AUROC of all other examined ML models based on DeLong's test and this was statistically significant. The AUROC curves for the five models in the test dataset are shown in Figure 2.

*Table 2: Performance of the ML models for AF prediction, on the test dataset, under various metrics.*

| **Models** | **AUROC (95% CI)** | **Accuracy** | **Precision** | **Recall** | **F1 score** |
|---|---|---|---|---|---|
| **LightGBM** | 0.729 (0.719-0.738) | 0.73 | 0.72 | 0.74 | 0.73 |
| **XGBoost** | 0.728 (0.718-0.737) | 0.73 | 0.74 | 0.73 | 0.73 |
| **DNN** | 0.716 (0.706-0.725) | 0.72 | 0.71 | 0.73 | 0.72 |
| **RF** | 0.715 (0.706-0.725) | 0.72 | 0.71 | 0.74 | 0.72 |
| **LR (L1 penalty)** | 0.622 (0.612-0.633) | 0.62 | 0.63 | 0.60 | 0.61 |

AUROC, the area under a receiver operating characteristic curve; Accuracy = (TP + TN) / (TP + TN + FP + FN); Precision = TP / (TP + FP), Recall = TP / (TP+FN) where TP stands for true positive, TN for true negative, FP for false positive, and FN for false negative; F1 score =2 (precision*recall) / (precision + recall).

*Table 3: DeLong's test for the ML model comparisons for AF prediction.*

| **Models** | **LightGBM** | **XGBoost** | **DNN** | **RF** |
|---|---|---|---|---|
| **LightGBM** | - | | | |
| **XGBoost** | 8.28E-01 | - | | |
| **DNN** | 3.67E-02 | 5.78E-02 | - | |

---

[1] This research utilised Queen Mary's Apocrita HPC facility, supported by QMUL Research-IT. http://doi.org/10.5281/zenodo.438045

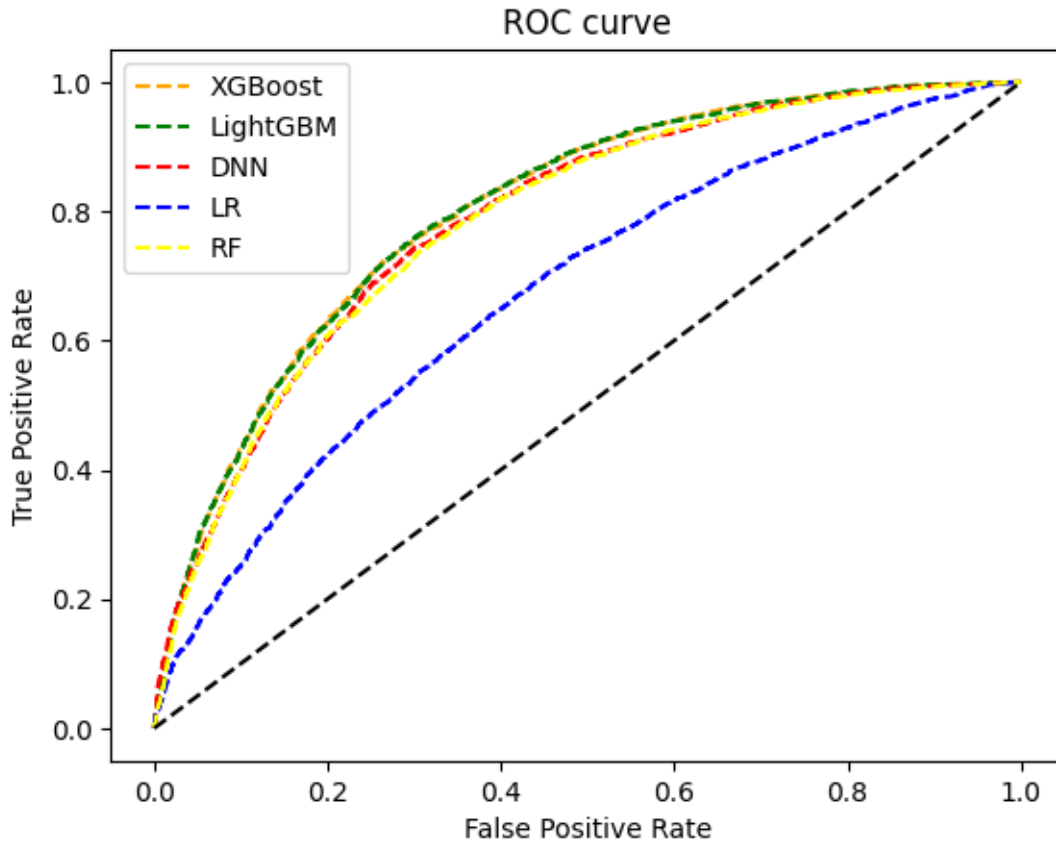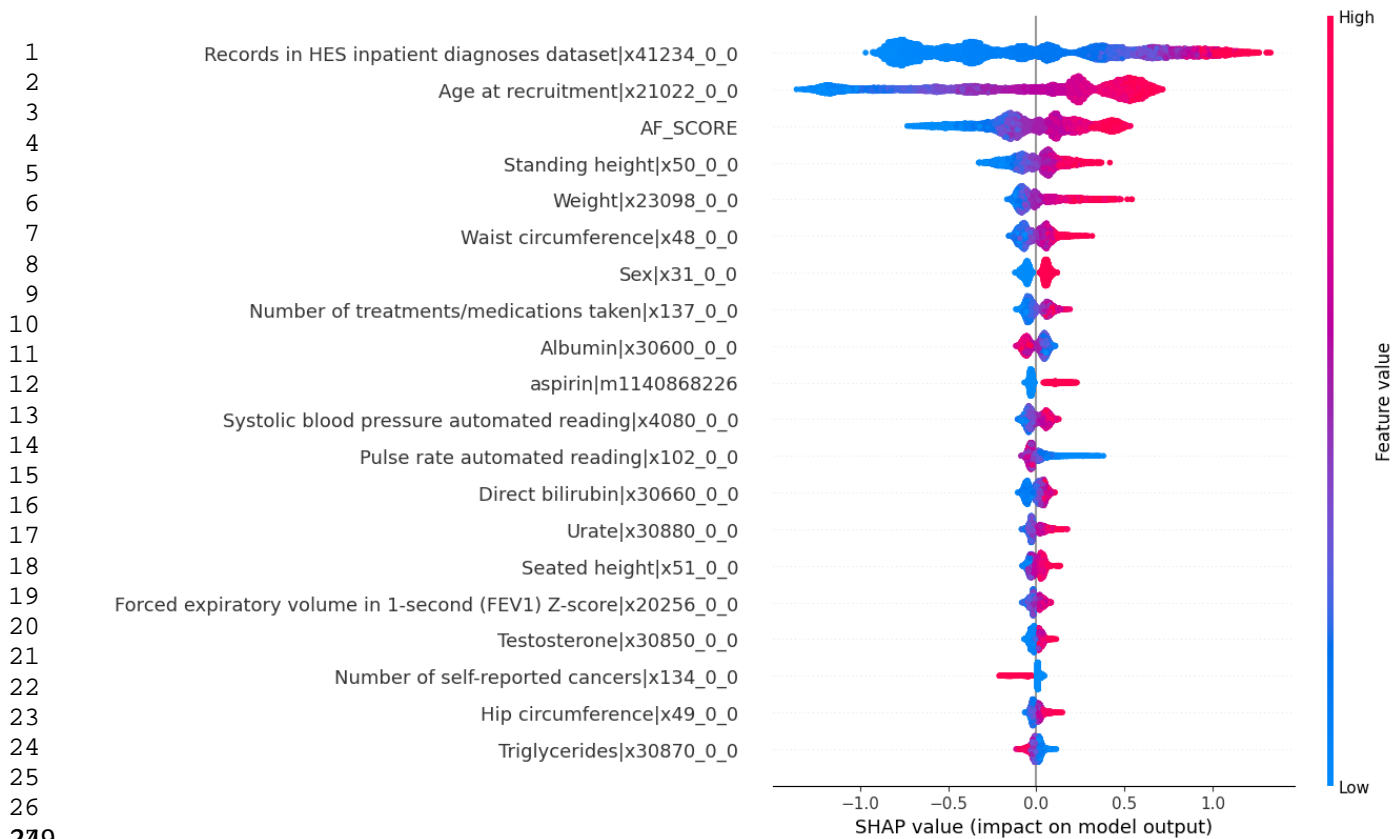| | | | | |
|---|---|---|---|---|
| **RF** | 1.17E-02 | 2.44E-02 | 9.91E-01 | - |
| **LR (L1 penalty)** | 1.38E-32 | 8.82E-32 | 2.41E-24 | 5.73E-27 |



*Figure 2: AUROC for each ML model for AF prediction in the test dataset.*

To estimate the contribution of each feature in each of the five models assessed for prediction of AF, we employed SHAP analysis, which is accurate, fast and stable. Figure 3 displays the top 20 features, ranked according to their SHAP value, for the LightGBM model; features are listed in descending order, starting with the most significant for AF prediction. SHAP values depict the distribution of the effect of each feature on the model output.

Based on Figure 3, SHAP analysis reveals that the top 3 most important variables contributing to the model were "Records in HES inpatient diagnoses dataset" which is the number of times an individual has been hospitalised (fieldID 41234), "Age at recruitment" (fieldID 21022) and "AF SCORE", using the unweighted sum of increasing alleles from Roseli et al. [11]. All the features' contributions, based on SHAP analysis, can be found in *Table_S7.*

*Figure 3: Summary plot of the SHAP values (x-axis) for the top 20 features (y-axis), in descending order, showing the distribution of the impact that each feature has for the AF prediction on the test dataset, employing LightGBM model. Each dot represents a participant. The red dots represent a high feature value and blue dots represent a low feature value for each participant. For example, the AF SCORE had a positive impact on the model output, i.e., a higher AF SCORE increased AF risk.*

### AF & Stroke

We examined 3,150 prospective cases who developed ischemic stroke after being diagnosed with AF, and an equal number of controls in UK-Biobank including 129 features (*Table_S8)* and using six models to predict ischemic stroke in AF cases. As indicated previously, results correspond to the optimal hyperparameters (*Table_S9)*.

The best AUROC value was achieved for XGBoost (Table 4). DeLong's test (Table 5) showed that there is no evidence for significant difference in the AUROCs between XGBoost and all other examined ML models but the penalised LR model (pvalue=2.00E-02) (Figure 4).

*Table 4: Performance of the ML models for the prediction of ischemic stroke in AF patients, on the test dataset, under various metrics.*

| Models | AUROC (95% CI) | Accuracy | Precision | Recall | F1 score |
|---|---|---|---|---|---|
| **XGBoost** | 0.631 (0.604-0.657) | 0.63 | 0.63 | 0.63 | 0.63 |
| **LightGBM** | 0.620 (0.593-0.647) | 0.62 | 0.62 | 0.61 | 0.62 |
| **RF** | 0.599 (0.573-0.625) | 0.60 | 0.61 | 0.56 | 0.58 |

| | | | | |
|---|---|---|---|---|
| SVM | 0.599 (0.572-0.624) | 0.60 | 0.63 | 0.50 | 0.55 |
| DNN | 0.589 (0.562-0.615) | 0.59 | 0.59 | 0.60 | 0.59 |
| LR (L1 penalty) | 0.563 (0.536-0.591) | 0.56 | 0.56 | 0.56 | 0.56 |

AUROC, the area under a receiver operating characteristic curve; Accuracy = (TP + TN) / (TP + TN + FP + FN); Precision = TP / (TP + FP), Recall = TP / (TP+FN) where TP stands for true positive, TN for true negative, FP for false positive, and FN for false negative; F1 score =2 (precision*recall) / (precision + recall).

*Table 5: DeLong's test for the ML model comparisons for ischemic stroke prediction in AF patients.*

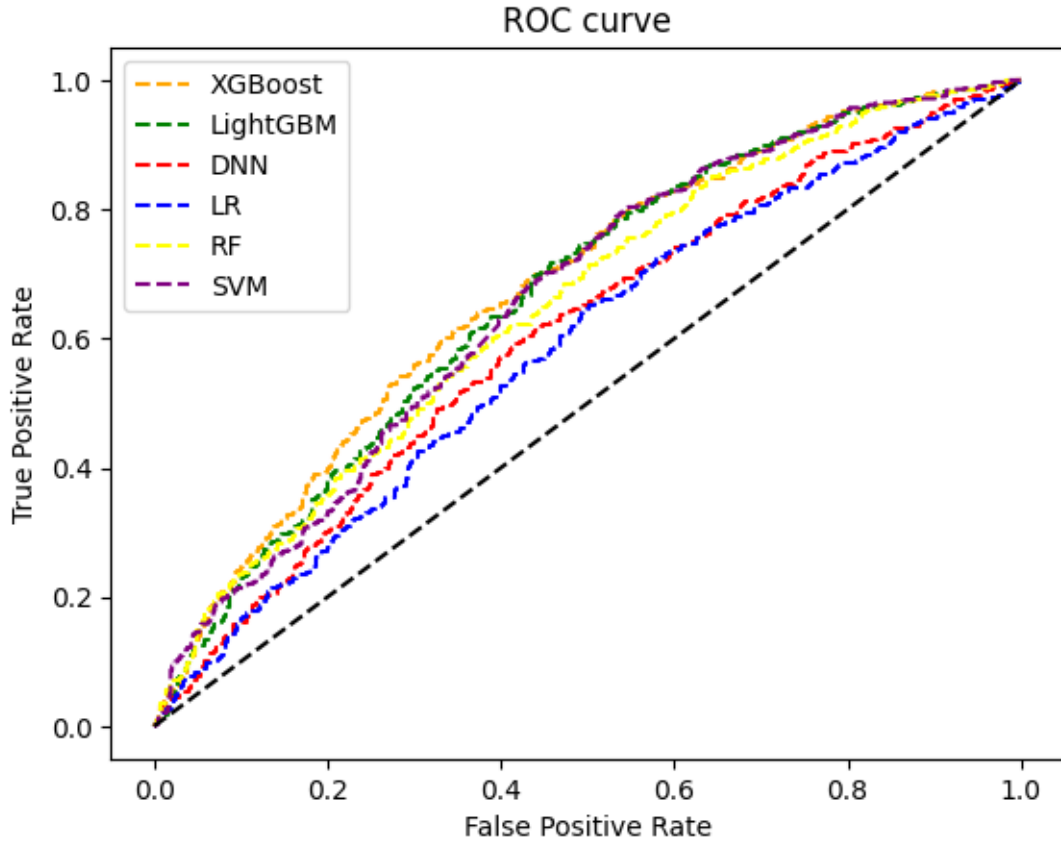| Models | XGBoost | LightGBM | RF | SVM | DNN |
|---|---|---|---|---|---|
| **XGBoost** | - | | | | |
| **LightGBM** | 5.65E-01 | - | | | |
| **RF** | 1.33E-01 | 3.45E-01 | - | | |
| **SVM** | 1.71E-01 | 3.75E-01 | 9.80E-01 | - | |
| **DNN** | 1.34E-01 | 2.89E-01 | 7.54E-01 | 7.45E-01 | - |
| **LR (L1 penalty)** | 2.00E-02 | 5.70E-02 | 2.56E-01 | 4.50E-01 | 2.54E-01 |

*Figure 4: AUROC for each ML model for predicting the development of ischemic stroke in AF patients, on the test dataset.*

As shown in Figure 5, SHAP analysis revealed that the 3 most important variables contributing to prediction of ischemic stroke in AF cases in the model were "Records in HES inpatient diagnoses dataset" which is the number of times an individual has been hospitalised (fieldID 41234), "Age at recruitment" (fieldID 21022), and "Glycated haemoglobin (HbA1c)" which is a blood biochemistry measurement (fieldID 30750). *Table_S10* lists the contribution of each of the 129 features in the model based on SHAP analysis.
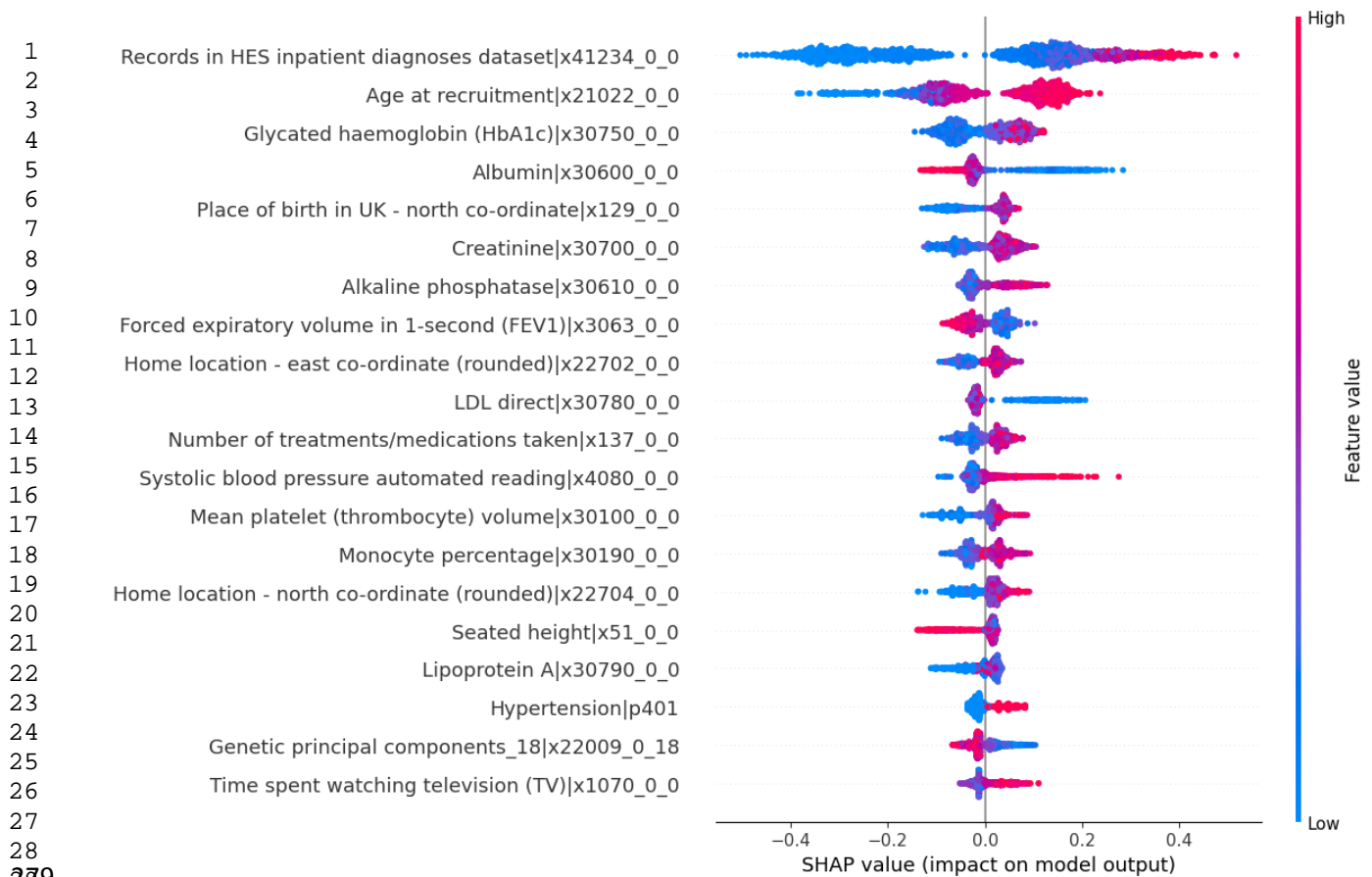
*Figure 5: Summary plot of the SHAP values (x-axis) for the top 20 features (y-axis), in descending order, showing the distribution of the impact that each feature has for the development of ischemic stroke in AF patients, on the test dataset, employing XGBoost model. Each dot represents a participant. The red dots represent a high feature value and blue dots represent a low feature value for each participant.*

### Comparison with CHA$_2$DS$_2$-VASc

The current tool used for prediction of ischemic stroke occurrence among AF patients is CHA$_2$DS$_2$-VASc which considers multiple risk factors; age, sex, heart failure, hypertension, stroke, vascular disease, diabetes [25]. Thus, we decided to compare the performance of the best ML model, XGBoost (Table 4), with CHA$_2$DS$_2$-VASc in UK-Biobank. To construct the CHA$_2$DS$_2$-VASc we employed the codes described in *Table_S11*. The AUROC and 95% CI for CHA$_2$DS$_2$-VASc and XGBoost was 0.611 (0.585 − 0.638) and 0.631 (0.604 − 0.657) in the test set, respectively. The improved AUROC in the XGBoost model compared to CHA$_2$DS$_2$-VASc was statistically significant based on DeLong's test (pvalue=2.20E-06). Furthermore, the SHAP analysis for the XGBoost model (Figure 5), shows that there is a significant number of peripheral blood markers associated with ischemic stroke, which are overlooked from CHA$_2$DS$_2$-VASc.

### Discussion

### Comparison of the performance of ML models for prediction of AF or ischemic stroke in patients with AF

We assessed six ML models in total for prediction of AF (XGBoost, LightGBM, RF, DNN, LR) or ischemic stroke in AF patients (XGBoost, LightGBM, RF, DNN, SVM, LR) and employed SHAP analysis to rank features for predictive importance. SHAP analysis was successful in the visualisation of non-linear relationships between the features used for prediction and the outcome. Additionally, the direction of the SHAP values for the top 20 features agrees with what has been reported so far in the literature. We found that the ensemble learning

models LightGBM (best for AF prediction) and XGBoost (best for prediction of ischemic stroke in patients with AF) achieved higher AUROCs compared to the other examined models, suggesting that these models have better generalisation. DeLong's test showed that penalised LR model had a lower AUROC compared to all other models and these differences were statistically significant (Table 3), indicating that ML models capture useful information by modeling non-linear associations, leading to the discovery of new features.

*AF results*

Advancing age has been shown to be one of the most important risk factors for AF [26], which is corroborated by the present study and ranked as the second most important feature. The third most important feature in the model was the AF SCORE, a set of 94 genome-wide variants associated with AF and explaining 42% of the heritability in Europeans [11], which as expected had a positive impact on the model output, i.e. the higher the AF score the higher the risk of developing AF. Thus, the present results endorse the likely clinical utility of an AF score in disease prediction. However, an optimised AF score for prediction in multi-ethnic populations such as the UK population will be required prior to considering clinical use. Interestingly, standing height was ranked as the fourth most significant feature in LightGBM, which was the best performing model for AF prediction. Greater height has been identified as a risk factor for AF in several studies and in both males and females [27], and it is in agreement with the present analysis. Some studies report that taller people have greater heart chamber size [27], meaning a larger left atrial size, which may be potential explanation albeit not a very robust one as AF is driven by left atrial stretch and fibrosis. Two other anthropometric traits, weight and waist circumference, ranked just below standing height. Obesity is associated with increased risk of left atrial enlargement, atrial fibrosis, electrical derangements of the atria, impaired diastolic function, inflammation and accumulation of pericardial fat, which are all key mechanisms in the pathogenesis of AF [28], and it is supported by the present analysis. The ranking of sex as the seventh most significant feature in the model is also in agreement with epidemiological studies reporting sex differences in AF; males are at higher risk which is in agreement with the results, along with the electrophysiologic properties of the atria and structural remodelling [29]. The analysis presented here also found that participants with lower albumin levels had an increased risk of AF. This is in agreement with a meta-analysis revealing that an increase in albumin level decreased the risk of AF [30]. However, low albumin levels are associated with poor health overall and therefore we cannot exclude confounding. Among the remaining 20 most significant features in the model it is worth noting that (i) direct bilirubin has been reported as an important independent risk factor for AF development in both thyrotoxic patients [31] and a study in postoperative cardiac surgery [32], (ii) urate has been reported to increase the risk of AF and be causally associated to AF through MR analysis in Koreans [33], and (iii) the positive effect of increased testosterone on risk of AF has been reported in males but not in females in the ARIC study [34]; the present study corroborates these results. Finally, only two of the 20 top features have some conflicting data in the literature. FEV-1 levels have an increased risk of AF as shown in other studies [35], and it is corroborated by the present analysis, but the Korean National Health and Nutritional Examination Survey reported an adverse association between FEV-1 and AF development [36]. Decreased levels of triglycerides contribute to increased risk of AF, but a study in Chinese participants contradicts the present analysis, showing no evidence of association between triglycerides and incidence of AF [37].

*AF & Ischemic stroke results*

In the present study, XGBoost model was the best in predicting ischemic stroke in AF patients and showed that it performs better than $CHA_2DS_2$-VASc, albeit marginal this result was statistically significant. Consistent with a recent French study for prediction of incident AF in a post-stroke population [38], the best performing ML model was DNN with a C index of 0.77 (95% CI 0.76-0.78) on the test set, performed better than $CHA_2DS_2$-VASc. In this study, XGBoost was identified as the best ML model for prediction of ischemic stroke in AF patients, with AUROC 0.631 (95% CI 0.604-0.657), in contrast to another two US studies that use more than

347     3.4 [39] and 6.4 [40] million participants, and reported c-index above 0.8. The lower performance of the ML

348     model could be attributed to the fact that we used 6,300 participants in contrast to the million that were used

349     in the US studies [39, 40], thus leading to less power.

350     Unexpectedly, the genetic risk score for ischemic stroke, based on 28 genome-wide variants, was not among

351     the top 20 features of the model, although ischemic stroke is highly heritable [41]. In the top 20 most

352     significant features, medium to high feature values of HbA1c ranked third after sex and was associated with

353     increased risk of stroke in AF patients. This agrees with the Clalit Health Services electronic medical records

354     Israelian database, where participants with diabetes and AF were found to have an increased risk of stroke

355     when their HbA1C levels were ranging from medium to high [42]. The fourth most significant feature was

356     albumin which ranked ninth in the AF prediction model, suggesting a stronger relationship with ischemic

357     stroke in AF patients than AF per se. This is corroborated by a Japanese study, which reported that lower

358     albumin levels were associated with an increased risk of ischemic stroke in both sexes independently of AF

359     status [43]. Four other blood biomarkers, creatinine, alkaline phosphatase, LDL cholesterol, and Lipoprotein A

360     (Lp(a)) ranked among the top 20 features. These results are in agreement with the China National Stroke

361     Registry reporting an association between high levels of alkaline phosphatase with recurrent stroke [44] and

362     the Copenhagen General Population Study showing that high levels of Lp(a) were associated with increased

363     risk of ischemic stroke [45]. It is worth noting that the latter although true for all examined ancestries it varies

364     in strength e.g. higher in African than European Americans [46]. Interestingly, the use of creatinine as marker

365     for increased risk of ischemic stroke in AF patients has not been previously reported and will merit further

366     investigation. Lastly, the twentieth feature identified from the SHAP analysis – time spent watching television

367     – could be considered as a surrogate marker for luck of sleep and physical inactivity; a recent study showed

368     that physical inactivity increases the risk of stroke risk [47].

## Conclusion

370     To conclude. there is a plethora of studies using ML methodology to predict circulatory diseases such as AF

371     [3], cardiovascular disease [48], stroke [4, 5], however none of them has the breadth and richness of electronic

372     health record data that UK Biobank offers, including disease diagnosis, medications and laboratory tests. The

373     strength of the present study is that makes use of the UK Biobank dataset, including up to 2,199 variables. The

374     present study supports the incorporation of a few routinely measured blood biomarkers, whereas the results

375     endorse the inclusion of a genetic score only in the model for AF prediction. The standardization of big data,

376     along with the wide application of machine and deep learning methodologies, enables the identification of

377     previously unknown risk factors for disease prediction. In the current study, the use of creatinine as marker

378     for increased risk of ischemic stroke in AF patients has not been previously reported, however it requires

379     further investigation. Machine learning models that employ large datasets, including potential predictors, can

380     improve prediction accuracy, as presented in the current study, for the prediction ischemic stroke in AF

381     patients using ML models in comparison to $CHA_2DS_2$-VASc, and provide graphical interpretation of the results

382     using SHAP analysis. The models presented here have the potential for clinical use, but validation in further

383     independent studies is required, since the models were developed and assessed in the UK Biobank and might

384     not reflect other datasets with respect to age, sex, socio-economic status [49]. The models would need to be

385     validated across all ancestries as some features vary by ethnicity e.g., Lp(a) and AF genetic score.

## Declaration of interests

387     Nothing to declare.

## Funding source

## Data availability

Individual level data could be accessed upon request and approval from UK Biobank. All the results discussed in this manuscript are available in the Supplementary Material.

**References**

1. Benjamin, E.J., et al., *Heart Disease and Stroke Statistics-2019 Update: A Report From the American Heart Association.* Circulation, 2019. **139**(10): p. e56-e528.
2. Khurshid, S., et al., *Performance of Atrial Fibrillation Risk Prediction Models in Over 4 Million Individuals.* Circ Arrhythm Electrophysiol, 2021. **14**(1): p. e008997.
3. Raghunath, S., et al., *Deep Neural Networks Can Predict New-Onset Atrial Fibrillation From the 12-Lead ECG and Help Identify Those at Risk of Atrial Fibrillation-Related Stroke.* Circulation, 2021. **143**(13): p. 1287-1298.
4. Su, P.Y., et al., *Machine Learning Models for Predicting Influential Factors of Early Outcomes in Acute Ischemic Stroke: Registry-Based Study.* JMIR Med Inform, 2022. **10**(3): p. e32508.
5. Jung, S., et al., *Predicting Ischemic Stroke in Patients with Atrial Fibrillation Using Machine Learning.* Front Biosci (Landmark Ed), 2022. **27**(3): p. 80.
6. Nishi, H., et al., *Predicting cerebral infarction in patients with atrial fibrillation using machine learning: The Fushimi AF registry.* J Cereb Blood Flow Metab, 2022. **42**(5): p. 746-756.
7. Kim, S.H., et al., *Interpretable machine learning for early neurological deterioration prediction in atrial fibrillation-related stroke.* Sci Rep, 2021. **11**(1): p. 20610.
8. Lundberg, S.M. and S.-I. Lee. *A unified approach to interpreting model predictions*. in *Proceedings of the 31st international conference on neural information processing systems*. 2017.
9. Millard, L.A.C., et al., *Searching for the causal effects of body mass index in over 300 000 participants in UK Biobank, using Mendelian randomization.* PLoS Genet, 2019. **15**(2): p. e1007951.
10. Wu, P., et al., *Mapping ICD-10 and ICD-10-CM Codes to Phecodes: Workflow Development and Initial Evaluation.* JMIR Med Inform, 2019. **7**(4): p. e14325.
11. Roselli, C., et al., *Multi-ethnic genome-wide association study for atrial fibrillation.* Nat Genet, 2018. **50**(9): p. 1225-1233.
12. Malik, R., et al., *Multiancestry genome-wide association study of 520,000 subjects identifies 32 loci associated with stroke and stroke subtypes.* Nat Genet, 2018. **50**(4): p. 524-537.
13. Pedregosa, F., et al., *Scikit-learn: Machine learning in Python.* the Journal of machine Learning research, 2011. **12**: p. 2825-2830.
14. Lemaître, G., F. Nogueira, and C.K. Aridas, *Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning.* The Journal of Machine Learning Research, 2017. **18**(1): p. 559-563.
15. Krawczyk, B., et al., *Evolutionary undersampling boosting for imbalanced classification of breast cancer malignancy.* Applied Soft Computing, 2016. **38**: p. 714-726.
16. AlJame, M., et al., *Ensemble learning model for diagnosing COVID-19 from routine blood tests.* Inform Med Unlocked, 2020. **21**: p. 100449.
17. Aridas, G.L.F.N.C.K., *Imbalanced-learn: A Python Toolbox to Tackle the Curse of Imbalanced Datasets in Machine Learning.* Journal of Machine Learning Research, 2017. **18(17)**: p. 1-5.
18. Berisha, V., et al., *Digital medicine and the curse of dimensionality.* NPJ Digit Med, 2021. **4**(1): p. 153.
19. Ismael, R.-P., et al., *When is resampling beneficial for feature selection with imbalanced wide data?* Expert Systems with Applications, 2022. **188**: p. 116015.
20. Chen, T. and C. Guestrin, *XGBoost: A Scalable Tree Boosting System*, in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2016, Association for Computing Machinery: San Francisco, California, USA. p. 785–794.
21. Ke, G., et al., *LightGBM: A Highly Efficient Gradient Boosting Decision Tree.* 2017.
22. LeCun, Y., Y. Bengio, and G. Hinton, *Deep learning.* Nature, 2015. **521**(7553): p. 436-44.
23. Cortes, C. and V. Vapnik, *Support-vector networks.* Machine learning, 1995. **20**(3): p. 273-297.
24. Van Rossum, G. and F.L. Drake, *The python language reference manual*. 2011: Network Theory Ltd.
25. Lip, G.Y., et al., *Refining clinical risk stratification for predicting stroke and thromboembolism in atrial fibrillation using a novel risk factor-based approach: the euro heart survey on atrial fibrillation.* Chest, 2010. **137**(2): p. 263-72.

444  26. Chung, M.K., et al., *Lifestyle and Risk Factor Modification for Reduction of Atrial Fibrillation: A*
445  *Scientific Statement From the American Heart Association.* Circulation, 2020. **141**(16): p. e750-e772.
446  27. Johansson, C., et al., *Weight, height, weight change, and risk of incident atrial fibrillation in middle-*
447  *aged men and women.* J Arrhythm, 2020. **36**(6): p. 974-981.
448  28. Feng, T., et al., *Weight and weight change and risk of atrial fibrillation: the HUNT study.* Eur Heart J,
449  2019. **40**(34): p. 2859-2866.
450  29. Westerman, S. and N. Wenger, *Gender Differences in Atrial Fibrillation: A Review of Epidemiology,*
451  *Management, and Outcomes.* Curr Cardiol Rev, 2019. **15**(2): p. 136-144.
452  30. Wang, Y., et al., *Relationship Between Serum Albumin and Risk of Atrial Fibrillation: A Dose-Response*
453  *Meta-Analysis.* Front Nutr, 2021. **8**: p. 728353.
454  31. Sun, D., et al., *Direct bilirubin level is an independent risk factor for atrial fibrillation in thyrotoxic*
455  *patients receiving radioactive iodine therapy.* Nucl Med Commun, 2019. **40**(12): p. 1289-1294.
456  32. Turkkolu, S.T., E. Selcuk, and C. Koksal, *Biochemical predictors of postoperative atrial fibrillation*
457  *following cardiac surgery.* BMC Cardiovasc Disord, 2021. **21**(1): p. 167.
458  33. Hong, M., et al., *A mendelian randomization analysis: The causal association between serum uric*
459  *acid and atrial fibrillation.* Eur J Clin Invest, 2020. **50**(10): p. e13300.
460  34. Berger, D., et al., *Plasma total testosterone and risk of incident atrial fibrillation: The Atherosclerosis*
461  *Risk in Communities (ARIC) study.* Maturitas, 2019. **125**: p. 5-10.
462  35. Au Yeung, S.L., et al., *Impact of lung function on cardiovascular diseases and cardiovascular risk*
463  *factors: a two sample bidirectional Mendelian randomisation study.* Thorax, 2022. **77**(2): p. 164-171.
464  36. Lee, S.N., et al., *Association between lung function and the risk of atrial fibrillation in a nationwide*
465  *population cohort study.* Sci Rep, 2022. **12**(1): p. 4007.
466  37. Li, X., et al., *Lipid profile and incidence of atrial fibrillation: A prospective cohort study in China.* Clin
467  Cardiol, 2018. **41**(3): p. 314-320.
468  38. Bisson, A., et al., *Prediction of incident atrial fibrillation in post-stroke patients using machine*
469  *learning: a French nationwide study.* Clin Res Cardiol, 2023. **112**(6): p. 815-823.
470  39. Lip, G.Y.H., et al., *Improving Stroke Risk Prediction in the General Population: A Comparative*
471  *Assessment of Common Clinical Rules, a New Multimorbid Index, and Machine-Learning-Based*
472  *Algorithms.* Thromb Haemost, 2022. **122**(1): p. 142-150.
473  40. Lip, G.Y.H., et al., *Improving dynamic stroke risk prediction in non-anticoagulated patients with and*
474  *without atrial fibrillation: comparing common clinical risk scores and machine learning algorithms.*
475  Eur Heart J Qual Care Clin Outcomes, 2022. **8**(5): p. 548-556.
476  41. O'Sullivan, J.W., et al., *Combining Clinical and Polygenic Risk Improves Stroke Prediction Among*
477  *Individuals With Atrial Fibrillation.* Circ Genom Precis Med, 2021. **14**(3): p. e003168.
478  42. Kezerle, L., et al., *Relation of Hemoglobin A1C Levels to Risk of Ischemic Stroke and Mortality in*
479  *Patients With Diabetes Mellitus and Atrial Fibrillation.* Am J Cardiol, 2022. **172**: p. 48-53.
480  43. Li, J., et al., *Serum Albumin and Risks of Stroke and Its Subtypes- The Circulatory Risk in Communities*
481  *Study (CIRCS).* Circ J, 2021. **85**(4): p. 385-392.
482  44. Zong, L., et al., *Alkaline Phosphatase and Outcomes in Patients With Preserved Renal Function:*
483  *Results From China National Stroke Registry.* Stroke, 2018. **49**(5): p. 1176-1182.
484  45. Kamstrup, P.R., *Lipoprotein(a) and Cardiovascular Disease.* Clin Chem, 2021. **67**(1): p. 154-166.
485  46. Kumar, P., et al., *Lipoprotein (a) level as a risk factor for stroke and its subtype: A systematic review*
486  *and meta-analysis.* Sci Rep, 2021. **11**(1): p. 15660.
487  47. Katzmarzyk, P.T., et al., *Physical inactivity and non-communicable disease burden in low-income,*
488  *middle-income and high-income countries.* Br J Sports Med, 2022. **56**(2): p. 101-106.
489  48. Joo, G., et al., *Clinical Implication of Machine Learning in Predicting the Occurrence of Cardiovascular*
490  *Disease Using Big Data (Nationwide Cohort Data in Korea).* IEEE Access, 2020. **8**: p. 157643-157653.
491  49. Fry, A., et al., *Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank*
492  *Participants With Those of the General Population.* Am J Epidemiol, 2017. **186**(9): p. 1026-1034.

493

Click here to access/download
**Supplementary Material**
Supplementary Figures.docx

Supplementary Tables

Click here to access/download
**Supplementary Material**
Supplementary_material.xlsx