Machine learning is more accurate and biased than risk scoring tools in the prediction of postoperative atrial fibrillation after cardiac surgery

Joyce C Ho *¹, Shalmali Joshi ⁺², Eduardo Valverde ^{‡3}, Kathryn Wood ^{§4}, Kendra Grubb ^{¶5},
 Miguel Leal ^{\105}, and Vicki Stover Hertzberg **⁴

¹Department of Computer Science, Emory University
 ²Department of Biomedical Informatics, Columbia University
 ³Department of Computer Science, Georgia Institute of Technology
 ⁴School of Nursing, Emory University
 ⁵School of Medicine, Emory University

11

Abstract

Incidence of postoperative atrial fibrillation (POAF) after cardiac surgery remains high and is 12 associated with adverse patient outcomes. Risk scoring tools have been developed to predict 13 POAF, yet discrimination performance remains moderate. Machine learning (ML) models can 14 achieve better performance but may exhibit performance heterogeneity across race and sex 15 subpopulations. We evaluate 8 risk scoring tools and 6 ML models on a heterogeneous cohort 16 derived from electronic health records. Our results suggest that ML models achieve higher 17 discrimination yet are less fair, especially with respect to race. Our findings highlight the need 18 for building accurate and fair ML models to facilitate consistent and equitable assessment of 19 POAF risk. 20

21

Keywords: POAF, ML, Fairness

*joyce.c.ho@emory.edu

- ⁺sj3261@cumc.columbia.edu
- ^{*}evalverde3@gatech.edu
- [§]kathryn.wood@emory.edu
- [¶]kendra.janel.grubb@emory.edu
- [®]miguel.a.leal@emory.edu
- **vhertzb@emory.edu

22 1 Introduction

Although there have been advancements in cardiac surgery techniques, the incidence of postop-23 erative atrial fibrillation (POAF) following cardiac surgery has not decreased significantly and still 24 ranges from 15% to 50% [1, 2]. Unfortunately, there are short- and long-term adverse outcomes 25 associated with POAF including morbidity, mortality, and longer, more expensive hospitalizations 26 [3, 4, 5, 6, 7]. Early identification of patients at risk for developing POAF has long been desired 27 to guide preventative and treatment strategies. To this end, more than a dozen POAF risk scoring 28 algorithms have been introduced encompassing a variety of risk factors including patient demo-29 graphics and clinical characteristics as well as surgical characteristics. Yet a recent review found 30 only patient age had no conflicting evidence across existing studies [8]. Moreover, these scoring 31 systems offer moderate discrimination with area under the receiver operating characteristic curve 32 (AUROC) scores ranging between 0.55 and 0.87 and may not generalize broadly as the performance 33 is assessed on relatively small, homogeneous patient populations. 34

Machine learning (ML) has been proposed as an alternative to achieve better predictive perfor-35 mance [9]. A recent scoping review found that support vector machines (SVM), gradient boosting 36 machines (GBM), and random forests (RF) using clinical characteristics can predict POAF risk more 37 accurately than existing risk scores with promising specificity, sensitivity, and AUROC scores [9]. 38 Three existing works compared multiple ML algorithms with Lu et al [10] and Parise et al. [11] 39 concluding that SVM achieved the best performance while GBM performed the best in Karri et al. 40 [12]. Despite their promise, indiscriminate application of ML models can exacerbate existing health 41 disparities if they are not trained on a representative sample [13]. 42

⁴³ Unfortunately, significant race and sex disparities exist as the number of patients undergoing car-⁴⁴ diac surgery procedures and the outcomes for these patients [14]. Incidence of POAF after coro-⁴⁵ nary artery bypass graft (CABG) surgery is higher in White patients [15]. It has also been suggested ⁴⁶ males are more likely to experience POAF following CABG [16, 17] although there exists conflicting ⁴⁷ evidence [18]. However, only 2 studies utilizing ML report the ethnicity composition of the under-⁴⁸ lying dataset and both studies assessed the performance in populations with less than 4% Black ⁴⁹ patients [12, 19]. Thus, a crucial unanswered question is whether the better performance of ML ⁵⁰ algorithms may exacerbate existing disparities.

The objective of this study is to assess both the predictive performance and fairness of existing POAF risk scoring tools with popular ML algorithms on a heterogeneous population, with more than 20% of the patients identifying as Black. We assess the fairness of the predictive models in both race and sex subpopulations. We also restrict our evaluation to common structured data found within electronic health records (EHRs) as such algorithms can provide quicker (and hopefully more accurate) management strategies [9].

57 2 Methods

58 2.1 Data Source

Our study was conducted using de-identified EHRs from the Emory Healthcare clinical data ware-59 house. Secondary data analysis was approved by the Emory University Institutional Review Board. 60 Adult patients who received cardiac surgery in the outpatient or inpatient setting between Jan-61 uary 1, 2013 and December 31, 2017 were included. Cardiac surgery was defined using the Current 62 Procedural Terminology (CPT) codes as either venous grafting for CABG or surgical procedures on 63 cardiac valves (see Supplemental Material for full list). For security purposes, patient identifiers 64 were omitted and certain records were excluded based on the date shifting logic. Patients who 65 had a prior history of atrial fibrillation (AF), defined by the International Classification of Diseases 66 codes of '427.31' for the 9th revision (ICD-9) or 'I48.XX' for the 10th revision (ICD-10) were ex-67 cluded from the study. We used the presence of the AF ICD-9 or ICD-10 code following the cardiac 68 surgery procedure date to identify cases of POAF. The value of 0 was assigned to patients that did 69 not experience POAF and had at least 1 encounter after the cardiac surgery. 70

All the clinical variables including age, sex, race, height, weight, and blood pressure were extracted from the EHR. We used the most recent value collected within the 1 year prior to the cardiac surgery date. The presence of clinical comorbidities for the risk scoring systems was determined using diagnostic (ICD-9 or ICD-10), procedural (CPT), and medication codes. For the ML clinical variables, we grouped the diagnostic codes using the single-level Clinical Classifications Software (CCS) system and medication codes using Anatomical Therapeutic Chemical (ATC) Level 3 classification

77 codes.

78 2.2 Risk Scores

We evaluated POAF risk scoring systems and incident AF risk scoring systems that utilize com-79 monly collected measures in structured EHR data. Although a recent review identified at least 80 12 distinct POAF scoring systems, [8] several used echocardiographic measurements such as left 81 atrial dilation, left atrial diameter, and left ventricular ejection fraction which are often captured 82 in unstructured text and are not easily accessible broadly. As such, we focused on the following 8 83 risk scores: (1) CHADS₂,[20] (2) CHA₂DS₂-VASc,[20] (3) HATCH, [21] (4) COM-AF, [22] (5) C₂HEST, 84 [23] (6) mC₂HEST, [24] (7) AFRI [25], and (8) CHARGE-AF [26]. The Python code for the scor-85 ing systems is openly available as a GitHub repository (https://github.com/joyceho/afib). The 86 predictor variables for each model can be found in Supplemental Table 3. 87

88 2.3 Machine Learning Models

Six commonly used ML algorithms were explored that have been previously benchmarked from 89 previous existing studies: (1) logistic regression (LR), (2) decision tree (DT), (3) SVM, (4) RF, (5) GBM, 90 and (6) multi-layer perceptron (MLP). The ML models were constructed using the popular Python 91 open-source software library, scikit-learn version 1.5.0 [27]. The ML models were supplied with 92 age, race, gender, CCS, and ATC codes. CCS and ATC codes that were not present in at least 93 20% of the patients were excluded. A total of 71 variables were supplied as input to the models. 94 Exhaustive hyperparameter optimization was performed using 5-fold cross validation on the training 95 dataset (see Supplemental Table 5 for the parameter search space for each model). The optimal 96 hyperparameter for each model was identified using AUROC. 97

98 2.4 Training and Evaluation

⁹⁹ We used stratified Monte Carlo cross validation to randomly split the data into 70-30% train-test. ¹⁰⁰ This process was repeated 10 times to assess model performance. Data imputation was required ¹⁰¹ for age, height, weight, and blood pressure. Mean imputation from the training data was used. ¹⁰² Predictive performance was measured using AUROC and area under the precision recall curve

(AUPRC) on the test set. AUROC and AUPRC were calculated using the scikit-learn package on 103 the test data. We also assessed the fairness of the models on the sex (female and male) and race 104 (White and Black/Other) subgroups. Two popular group fairness metrics were used, demographic 105 parity ratio (DPR) and equalized odds ratio (EQR). DPR measures whether the predictive proportion 106 of POAF across the subgroups are equal (i.e., the prediction risk should be independent of sex or 107 race). EQR ensures the true positive rate and false positive rate of predictions are the same across 108 the subgroups. Both DPR and EQR range from 0 to 1 with 1 indicates fairness across the subgroups. 109 DPR and EQR were computed using the fairlearn package version 0.10.0.[28] All analyses were 110 performed using Python version 3.9.7. 111

112 **3 Results**

3.1 Patient Characteristics

Out of the final study population of 4961 patients, 1953 (39.4%) experienced POAF following cardiac surgery with an average onset of xxx. Baseline characteristics of the overall study population and the 2 outcome groups (no POAF and POAF) are reported in Table 1. The incidence of POAF experienced in males (40.1%) and Whites (42.5%) was statistically higher than in females (37.9%) and Blacks (32.3%), respectively.

119 3.2 Performance Comparison

Table 2 summarizes the discrimination and fairness performance of the 14 models (8 risk scoring algorithms and 6 ML models). For each performance metric, the value represents the mean across the 10 test splits. Statistical significance in discrimination performance between any 2 models was assessed using a one-tailed paired t-test that the difference is greater than 0 (i.e., one model consistently outperforms the other).

The ML model that achieved the best discrimination was RF with AUROC and AUPRC of 0.671 and 0.558, respectively. Only GBM yielded a p-value above 0.001 for AUROC (0.03) and AUPRC (0.06) during the one-tailed paired t-test between RF and the other 5 ML models. Among the risk scoring systems, CHARGE-AF achieved the best performance with AUROC and AUPRC of 0.585 and 0.449,

	Missing	Overall	No POAF	POAF	P-Value
		(n=4961)	(n=3008)	(n=1953)	
Age, mean (SD)	19	64.6 (12.4)	63.3 (13.1)	66.5 (10.9)	<0.001
Male, n (%)	0	3300 (66.5)	1976 (65.7)	1324 (67.8)	0.133
Ethnicity, n (%)	0				<0.001
White		3426 (69.1)	1970 (65.5)	1456 (74.6)	
Black		1085 (21.9)	735 (24.4)	350 (17.9)	
Other		450 (9.1)	303 (10.1)	147 (7.5)	
CHF, n (%)	0	2066 (41.6)	1226 (40.8)	840 (43.0)	0.123
HTN, n (%)	0	2999 (60.5)	1888 (62.8)	1111 (56.9)	<0.001
DM, n (%)	0	1382 (27.9)	864 (28.7)	518 (26.5)	0.098
Stroke, n (%)	0	959 (19.3)	560 (18.6)	399 (20.4)	0.123
Vascular Disease, n (%)	0	2004 (40.4)	1253 (41.7)	751 (38.5)	0.027
COPD, n (%)	0	509 (10.3)	307 (10.2)	202 (10.3)	0.914
CAD, n (%)	0	4041 (81.5)	2503 (83.2)	1538 (78.8)	<0.001
SHF, n (%)	0	379 (7.6)	228 (7.6)	151 (7.7)	0.887
Hyperthyroid, n (%)	0	18 (0.4)	14 (0.5)	4 (0.2)	0.211
Valvular, n (%)	0	3099 (62.5)	1783 (59.3)	1316 (67.4)	<0.001
PVD, n (%)	0	1265 (25.5)	723 (24.0)	542 (27.8)	0.004
Obesity, n (%)	0	420 (8.5)	255 (8.5)	165 (8.4)	1.000
Height (m), mean (SD)	65	67.7 (4.6)	67.4 (4.6)	68.0 (4.5)	<0.001
Weight (kg) , mean (SD)	41	192.0 (45.8)	189.9 (44.5)	195.3 (47.5)	<0.001
SBP, mean (SD)	40	135.3 (22.3)	135.1 (22.2)	135.5 (22.4)	0.559
DBP, mean (SD)	41	74.4 (12.9)	74.8 (12.9)	73.8 (12.8)	0.015

Table 1: Baseline characteristics in patients with and without POAF.

Table 2: Average discrimination and fairness performance of the prediction models across 10 Monte Carlo cross-validation splits.

			Ra	Race		ex	
Model	AUROC	AUPRC	DPR	EQR	DPR	EQR	
RF	0.671	0.558	0.824	0.816	0.967	0.946	
GBM	0.666	0.553	0.837	0.835	0.975	0.948	
SVM	0.658	0.540	0.828	0.827	0.976	0.958	
LR	0.645	0.521	0.659	0.618	0.949	0.903	
DT	0.630	0.534	0.970	0.950	0.962	0.942	
MLP	0.622	0.494	0.789	0.772	0.953	0.923	
CHARGE-AF	0.578	0.443	0.896	0.895	0.941	0.926	
AFRI	0.568	0.474	1.000	1.000	1.000	1.000	
COM-AF	0.550	0.416	0.993	0.980	0.884	0.858	
HAVOC	0.527	0.401	0.999	0.999	0.999	0.998	
HATCH	0.527	0.400	1.000	1.000	1.000	1.000	
mC_2HEST	0.526	0.408	1.000	1.000	1.000	1.000	
CHA_2DS_2 -VASc	0.525	0.408	1.000	1.000	1.000	1.000	
$CHADS_2$	0.517	0.401	1.000	1.000	1.000	1.000	
C_2 HEST	0.509	0.401	1.000	1.000	1.000	1.000	

respectively. Notably, all 6 ML models outperformed CHARGE-AF and the other risk scoring tools at statistically significant levels (p-value < 0.001) for both discrimination metrics.

The risk scoring systems generally yielded the best group fairness concerning race as all but CHARGE-AF, COM-AF, and HAVOC resulted in both DPR and EQR of 1. Notably, CHARGE-AF is the only risk scoring system incorporating race as a variable (see Supplemental Table 3), yet achieves the worst group fairness. In contrast, all the ML models except DT perform worse in terms of DPR and EQR to CHARGE-AF (DPR = 0.896 and EQR = 0.895).

A similar group fairness trend is observed for sex in terms of risk scoring systems again as CHARGE AF, COM-AF, and HAVOC do not yield DPR and EQR of 1. However, the ML models achieve slightly
 better performance than CHARGE-AF and COM-AF in terms of DPR and EQR for race. Surprisingly,
 COM-AF which incorporates sex as a variable yields the worst DPR and EQR performance with
 values of 0.884 and 0.858, respectively.

141 **4** Discussion

In this study, we evaluated the performance of 6 ML models and 8 risk scoring algorithms to predict 142 POAF. We demonstrated that RF outperformed the other ML methods and all of the risk scores 143 considered in terms of AUROC and AUPRC. Furthermore, there were statistical differences between 144 the discrimination performance of RF and the other models except for the GBM algorithm. The 145 AUROC and AUPRC of these 14 models were all under 0.671 and 0.558 respectively. Compared 146 to the existing ML studies, the discrimination performance is lower as they achieved an AUROC 147 of at least 0.72. However, these models used indicators related to cardiac surgery which are not 148 commonly available in the structured EHR data. 149

In contrast to the discrimination performance, six of the risk scores outperformed all of the ML methods and three of the risk scores with respect to metrics of fairness. In fact, the results indicate the ML models exacerbated race and sex differentials when used for POAF prediction, which is consistent with existing evidence for other outcomes such as cardiovascular risk [29], dermatology [30], and population health [31]. Thus, better discrimination performance may not always be desired as it might exacerbate existing race and sex disparities. This suggests further investigation is

- ¹⁵⁶ necessary to holistically assess the efficacy of ML algorithms for POAF prognostication in real
- ¹⁵⁷ clinical contexts, [32] and whether bias mitigation mechanisms should be adopted to minimize
- ¹⁵⁸ disparities in outcomes and interventions.

Abbreviations

- ¹⁶⁰ **POAF**: post-operative atrial fibrillation; **AUROC**: area under the receiver operating characteristic
- ¹⁶¹ curve; **ML**: machine learning; **SVM**: support vector machines; **GBM**: gradient boosting machines;
- ¹⁶² **RF**: random forests; **CABG**: coronary artery bypass graft

5 Supplemental Information

164 5.1 Risk scoring algorithms

The variables used for each of the 8 risk scores are summarized in Table 3. While other risk scoring systems have been developed using POAF as an event of interest[8], they are not benchmarked in our study as they use variables such as left atrial dilatation, left atrial diameter, left ventricular ejection fraction, and length of stenosis. The scores in Table 3 can be computed from demographic information, diagnosis tables (ICD-9/ICD-10), vital signs commonly collected, and medication tables.

Table 3: Risk scoring systems, the original event of interests, and their associated variables. The events of interest are thromboembolic events in patients with atrial fibrillation (VTE), incident AF (AF), and POAF.

Score	Event	Age	Gender	Race	CHF	HTN	DM	Stroke	Vascular	СОРD	CAD	SHF	Hyperthyroid	PVD	IM	Height	Weight	SBP	DBP	Smoking	HTN Meds
CHADS ₂ [20]	VTE	х			х	х	х	х													
CHA_2DS_2 -VASc [20]	VTE	х	х		Х	х	х	Х	х												
HATCH [21]	AF	х			Х	х				х											
COM-AF [<mark>22</mark>]	POAF	х	х		Х	х	х	Х													
C ₂ HEST [23]	AF	х				х				х	х	Х	х								
mC ₂ HEST [24]	AF	х				Х				Х	х	х	х								
AFRI [<mark>25</mark>]	POAF	х												Х		х	х				
CHARGE-AF [<mark>26</mark>]	AF	Х		х	Х	х	х								Х	х	х	Х	Х	х	Х

171 **5.2 ML for POAF**

¹⁷² A recent scoping review identified 7 papers that used ML for predicting POAF after cardiac surgery.[9]

¹⁷³ Of the 7 studies, 3 relied on electrocardiogram data while the remaining 4 used clinical documenta-

tion, administrative data, or Holter monitoring. The sample size, ethnicity composition, and model

performance of the 4 ML studies using administrative data are summarized in Table 4. As can be

¹⁷⁶ seen, none of the patient populations contains more than 3.4% Black.

Table 5 summarizes the hyperparameter search space for each of the ML models. For each traintest split and ML model, GridSearchCV in scikit-learn was performed using 5-folds on the train

Authors	Dataset	n	White-Black (%)	Best ML (AUROC)
Magee et al. [19]	Cardiopulmonary Research Sci- ence and Technology Institute	19620	90.7-3.2	LR (0.72)
Karri et al. [<mark>12</mark>]	MIMIC-III	6040	74.0-3.4	GBM (0.74)
Lu et al. [10]	Second Affiliated Hospital of Zhejiang University School of Medicine	1400	Unknown	SVM (0.78)
Parise et al. [11]	Maastricht University Medical Center+	394	Unknown	SVM (0.95)

Table 5: Hyperparameter search space for the different ML models.

Model	Parameter Space
LR	C: [0.001, 0.01, 0.1, 1]
DT	Criterion: [gini, entropy]
	Max depth: [3, 4, 5, 6, 7, 8]
RF	Max depth: [10, 15, 20]
	Min leaf samples: [5, 10]
	Num estimators: [50, 100, 200, 300]
GBM	Max depth: [5, 8, 10]
	Learning rate: [0.1, 0.01]
	Min leaf samples: [5, 10]
	Num estimators: [50, 100, 150, 200]
SVM	C: [0.1, 1, 10, 100]
	Kernel: [linear, rbf]
MLP	Hidden layers: [(15, 10), (15,), (25,), (50,), (25, 25,)]

¹⁷⁹ split to find the optimal hyperparameter values. The ML model is then retrained using the optimal

¹⁸⁰ hyperparameter values and the performance is evaluated on the test set.

181 References

- [1] Giovanni Filardo, Ralph J Damiano, Gorav Ailawadi, Vinod H Thourani, Benjamin D Pollock,
 Danielle M Sass, Teresa K Phan, Hoa Nguyen, and Briget Da Graca. Epidemiology of new onset atrial fibrillation following coronary artery bypass graft surgery. *Heart*, 104(12):985–992,
 2018.
- [2] Orlando R Suero, Ahmed K Ali, Lauren R Barron, Matthew W Segar, Marc R Moon, and Subhasis
 Chatterjee. Postoperative atrial fibrillation (poaf) after cardiac surgery: clinical practice review.
 Journal of Thoracic Disease, 16(2), 2024.
- [3] Rachel Eikelboom, Rohan Sanjanwala, Me-Linh Le, Michael H Yamashita, and Rakesh C Arora.
 Postoperative atrial fibrillation after cardiac surgery: a systematic review and meta-analysis.
 The Annals of Thoracic Surgery, 111(2):544–554, 2021.
- [4] Ben O'Brien, Peter S. Burrage, Jennie Yee Ngai, Jordan M. Prutkin, Chuan-Chin Huang, Xinling
 Xu, Sanders H. Chae, Bruce A. Bollen, Jonathan P. Piccini, Nanette M. Schwann, Aman Mahajan,
 Marc Ruel, Simon C. Body, Frank W. Sellke, Joseph Mathew, and J. Daniel Muehlschlegel. Soci ety of cardiovascular anesthesiologists/european association of cardiothoracic anaesthetists
 practice advisory for the management of perioperative atrial fibrillation in patients undergoing
 cardiac surgery. *Journal of Cardiothoracic and Vascular Anesthesia*, 33(1):12–26, 2019.
- [5] Ahmed AlTurki, Mariam Marafi, Riccardo Proietti, Daniela Cardinale, Robert Blackwell, Paul
 Dorian, Amal Bessissow, Lucy Vieira, Isabelle Greiss, Vidal Essebag, Jeff S. Healey, and Thao
 Huynh. Major adverse cardiovascular events associated with postoperative atrial fibrillation
 after noncardiac surgery: a systematic review and meta-analysis. *Circulation: Arrhythmia and Electrophysiology*, 13(1):e007437, 2020.
- [6] Peter S Burrage, Ying H Low, Niall G Campbell, and Ben O'Brien. New-onset atrial fibrillation
 in adult patients after cardiac surgery. *Current anesthesiology reports*, 9:174–193, 2019.
- [7] Michael K. Wang, Pascal B. Meyre, Rachel Heo, P.J. Devereaux, Lauren Birchenough, Richard
 Whitlock, William F. McIntyre, Yu Chiao Peter Chen, Muhammad Zain Ali, Fausto Biancari,
 Jawad Haider Butt, Jeff S. Healey, Emilie P. Belley-Côté, Andre Lamy, and David Conen. Short-

term and long-term risk of stroke in patients with perioperative atrial fibrillation after cardiac
 surgery: systematic review and meta-analysis. *CJC open*, 4(1):85–96, 2022.

- [8] Hugh Fleet, David Pilcher, Rinaldo Bellomo, and Tim G Coulson. Predicting atrial fibrillation
 after cardiac surgery: a scoping review of associated factors and systematic review of existing
 prediction models. *Perfusion*, 38(1):92–108, 2023.
- [9] Adham H El-Sherbini, Aryan Shah, Richard Cheng, Abdelrahman Elsebaie, Ahmed A Harby,
 Damian Redfearn, and Mohammad El-Diasty. Machine learning for predicting postoperative
 atrial fibrillation after cardiac surgery: A scoping review of current literature. *The American Journal of Cardiology*, 209:66–75, 2023.

[10] Yufan Lu, Qingjuan Chen, Hu Zhang, Meijiao Huang, Yu Yao, Yue Ming, Min Yan, Yunxian Yu, and
 Lina Yu. Machine learning models of postoperative atrial fibrillation prediction after cardiac
 surgery. *Journal of Cardiothoracic and Vascular Anesthesia*, 37(3):360–366, 2023.

- [11] Orlando Parise, Gianmarco Parise, Akshayaa Vaidyanathan, Mariaelena Occhipinti, Ali Ghar aviri, Cecilia Tetta, Elham Bidar, Bart Maesen, Jos G. Maessen, Mark La Meir, and Sandro
 Gelsomino. Machine learning to identify patients at risk of developing new-onset atrial fib rillation after coronary artery bypass. *Journal of Cardiovascular Development and Disease*,
 10(2):82, 2023.
- [12] Roshan Karri, Andrew Kawai, Yoke Jia Thong, Dhruvesh M Ramson, Luke A Perry, Reny Segal,
 Julian A Smith, and Jahan C Penny-Dimri. Machine learning outperforms existing clinical
 scoring tools in the prediction of postoperative atrial fibrillation during intensive care unit
 admission after cardiac surgery. *Heart, Lung and Circulation*, 30(12):1929–1937, 2021.

[13] Irene Y Chen, Emma Pierson, Sherri Rose, Shalmali Joshi, Kadija Ferryman, and Marzyeh Ghassemi. Ethical machine learning in healthcare. *Annual review of biomedical data science*, 4:123–144, 2021.

[14] Mohamad Alkhouli, Fahad Alqahtani, David R Holmes, and Chalak Berzingi. Racial disparities
 in the utilization and outcomes of structural heart disease interventions in the united states.
 Journal of the American Heart Association, 8(15):e012125, 2019.

[15] David W Yaffee, Raymond G McKay, Jeffrey Mather, Scott Vella Sorensen, Andrew Kehm, Sean
 McMahon, Trevor Sutton, and Sabet W Hashim. Racial disparities in atrial fibrillation after
 coronary artery bypass: Impact of left atrial volume. *Annals of Thoracic Surgery Short Reports*,
 1(4):631–634, 2023.

[16] Mariana Fragão-Marques, Jennifer Mancio, João Oliveira, Inês Falcão-Pires, and Adelino
 Leite-Moreira. Gender differences in predictors and long-term mortality of new-onset postop erative atrial fibrillation following isolated aortic valve replacement surgery. *Annals of Thoracic and Cardiovascular Surgery*, 26(6):342–351, 2020.

[17] Giovanni Filardo, Gorav Ailawadi, Benjamin D Pollock, Briget da Graca, Teresa K Phan, Vinod
 Thourani, and Ralph J Damiano Jr. Postoperative atrial fibrillation: sex-specific characteristics
 and effect on survival. *The Journal of thoracic and cardiovascular surgery*, 159(4):1419–1425,
 2020.

[18] Giovanni Filardo, Gorav Ailawadi, Benjamin D Pollock, Briget Da Graca, Danielle M Sass,
 Teresa K Phan, Debbie E Montenegro, Vinod Thourani, and Ralph Damiano. Sex differences in
 the epidemiology of new-onset in-hospital post-coronary artery bypass graft surgery atrial
 fibrillation: a large multicenter study. *Circulation: Cardiovascular Quality and Outcomes*,
 9(6):723-730, 2016.

[19] Mitchell J Magee, Morley A Herbert, Todd M Dewey, James R Edgerton, William H Ryan, Syma
 Prince, and Michael J Mack. Atrial fibrillation after coronary artery bypass grafting surgery:
 development of a predictive risk algorithm. *The Annals of thoracic surgery*, 83(5):1707–1712,
 2007.

[20] Gregory YH Lip, Robby Nieuwlaat, Ron Pisters, Deirdre A Lane, and Harry JGM Crijns. Refin ing 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*,
 137(2):263–272, 2010.

[21] Cees B De Vos, Ron Pisters, Robby Nieuwlaat, Martin H Prins, Robert G Tieleman, Robert-Jan S
 Coelen, Antonius C van den Heijkant, Maurits A Allessie, and Harry JGM Crijns. Progression

from paroxysmal to persistent atrial fibrillation: clinical correlates and prognosis. *Journal of* the American College of Cardiology, 55(8):725–731, 2010.

[22] Lucrecia M Burgos, Andreína Gil Ramírez, Leonardo Seoane, Juan F Furmento, Juan P Costa bel, Mirta Diez, and Daniel Navia. New combined risk score to predict atrial fibrillation after
 cardiac surgery: Com-af. *Annals of cardiac anaesthesia*, 24(4):458–463, 2021.

[23] Yan-Guang Li, Daniele Pastori, Alessio Farcomeni, Pil-Sung Yang, Eunsun Jang, Boyoung
 Joung, Yu-Tang Wang, Yu-Tao Guo, and Gregory YH Lip. A simple clinical risk score (c2hest) for
 predicting incident atrial fibrillation in asian subjects: derivation in 471,446 chinese subjects,
 with internal validation and external application in 451,199 korean subjects. *Chest*, 155(3):510–
 518, 2019.

[24] Yan-Guang Li, Jin Bai, Gongbu Zhou, Juan Li, Yi Wei, Lijie Sun, Lingyun Zu, and Shuwang
 Liu. Refining age stratum of the c2hest score for predicting incident atrial fibrillation in a
 hospital-based chinese population. *European journal of internal medicine*, 90:37-42, 2021.

[25] Mikhael F El-Chami, Patrik D Kilgo, K Miriam Elfstrom, Michael Halkos, Vinod Thourani, Omar M
 Lattouf, David B Delurgio, Robert A Guyton, Angel R Leon, and John D Puskas. Prediction of
 new onset atrial fibrillation after cardiac revascularization surgery. *The American journal of cardiology*, 110(5):649–654, 2012.

[26] Alvaro Alonso, Bouwe P Krijthe, Thor Aspelund, Katherine A Stepas, Michael J Pencina, Car-279 lee B Moser, Moritz F Sinner, Nona Sotoodehnia, João D Fontes, A Cecile JW Janssens, 280 Richard A Kronmal, Jared W Magnani, Jacqueline C Witteman, Alanna M Chamberlain, 281 Steven A Lubitz, Renate B Schnabel, Sunil K Agarwal, David D McManus, Patrick T Ellinor, 282 Martin G Larson, Gregory L Burke, Lenore J Launer, Albert Hofman, Daniel Levy, John S Gott-283 diener, Stefan Kääb, David Couper, Tamara B Harris, Elsayed Z Soliman, Bruno H C Stricker, 284 Vilmundur Gudnason, Susan R Heckbert, and Emelia J Benjamin. Simple risk model predicts 285 incidence of atrial fibrillation in a racially and geographically diverse population: the charge-af 286 consortium. Journal of the American Heart Association, 2(2):e000102, 2013. 287

[27] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vander-

plas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard
 Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
 12:2825–2830, 2011.

[28] Hilde Weerts, Miroslav Dudík, Richard Edgar, Adrin Jalali, Roman Lutz, and Michael Madaio.
 Fairlearn: Assessing and improving fairness of ai systems. *Journal of Machine Learning Research*, 24(257):1–8, 2023.

[29] Uri Kartoun, Shaan Khurshid, Bum Chul Kwon, Aniruddh P Patel, Puneet Batra, Anthony Philip pakis, Amit V Khera, Patrick T Ellinor, Steven A Lubitz, and Kenney Ng. Prediction performance
 and fairness heterogeneity in cardiovascular risk models. *Scientific Reports*, 12(1):12542, 2022.

[30] Adewole S Adamson and Avery Smith. Machine learning and health care disparities in der matology. *JAMA dermatology*, 154(11):1247–1248, 2018.

[31] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. Dissecting racial
 bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453,
 2019.

[32] Melissa Mccradden, Oluwadara Odusi, Shalmali Joshi, Ismail Akrout, Kagiso Ndlovu, Ben
 Glocker, Gabriel Maicas, Xiaoxuan Liu, Mjaye Mazwi, Tee Garnett, Lauren Oakden-Rayner, Myrt ede Alfred, Irvine Sihlahla, Oswa Shafei, and Anna Goldenberg. What's fair is... fair? presenting
 justefab, an ethical framework for operationalizing medical ethics and social justice in the in tegration of clinical machine learning: Justefab. In *Proceedings of the 2023 ACM Conference* on Fairness, Accountability, and Transparency, pages 1505–1519, 2023.