# Validation, bias assessment, and optimization of the UNAFIED 2-year risk prediction model for undiagnosed atrial fibrillation using national electronic health data

Mohammad Ateya, PharmD, MS,[1] Danai Aristeridou, MSc,[1] George H. Sands, MD,[1] Jessica Zielinski, BS,[1] Randall W. Grout, MD, MS,[2,3] A. Carmine Colavecchia, PharmD, PhD,[1] Oussama Wazni, MD, FHRS,[4] Saira N. Haque, PhD, MHSA[1]

*From the [1]Pfizer Inc, New York, New York, [2]Regenstrief Institute, Indianapolis, Indiana, [3]Indiana University School of Medicine, Indianapolis, Indiana, and [4]Cleveland Clinic, Cleveland, Ohio.*

**BACKGROUND** Prediction models for atrial fibrillation (AF) may enable earlier detection and guideline-directed treatment decisions. However, model bias may lead to inaccurate predictions and unintended consequences.

**OBJECTIVE** The purpose of this study was to validate, assess bias, and improve generalizability of "UNAFIED-10," a 2-year, 10-variable predictive model of undiagnosed AF in a national data set (originally developed using the Indiana Network for Patient Care regional data).

**METHODS** UNAFIED-10 was validated and optimized using Optum de-identified electronic health record data set. AF diagnoses were recorded in the January 2018–December 2019 period (outcome period), with January 2016–December 2017 as the baseline period. Validation cohorts (patients with AF and non-AF controls, aged ≥40 years) comprised the full imbalanced and randomly sampled balanced data sets. Model performance and bias in patient subpopulations based on sex, insurance, race, and region were evaluated.

**RESULTS** Of the 6,058,657 eligible patients (mean age 60 ± 12 years), 4.1% (n = 246,975) had their first AF diagnosis within the outcome period. The validated UNAFIED-10 model achieved a higher C-statistic (0.85 [95% confidence interval 0.85–0.86] vs 0.81 [0.80–0.81]) and sensitivity (86% vs 74%) but lower specificity (66% vs 74%) than the original UNAFIED-10 model. During retraining and optimization, the variables insurance, shock, and albumin were excluded to address bias and improve generalizability. This generated an 8-variable model (UNAFIED-8) with consistent performance.

**CONCLUSION** UNAFIED-10, developed using regional patient data, displayed consistent performance in a large national data set. UNAFIED-8 is more parsimonious and generalizable for using advanced analytics for AF detection. Future directions include validation on additional data sets.

**KEYWORDS** Atrial fibrillation; Screening; Predictive model; Machine learning; Electronic health record

## Introduction

Atrial fibrillation (AF) is a major public health concern that is growing in importance because of its association with significant morbidity, mortality, and economic burden.[1,2] Patients with AF have an elevated risk of cardiovascular (CV) events, including ∼5 times higher risk of AF-related stroke, systemic embolism, and ischemic heart disease, than does the general population.[1–3] Although an aging US population partially explains the continued increasing prevalence of AF, previous studies have indicated that ∼11%–23% of the total AF prevalence in 2009 (3.9 million patients) and 2015 (5.6 million patients) had undiagnosed AF.[4–10] Projections indicate that the prevalence of AF in the United States will reach ∼12.1 million patients by 2030 and ∼17.9 million patients (aged >55 years) in Europe by 2060.[4–7]

AF can be asymptomatic and transient, which makes it difficult to detect.[11,12] Many individuals remain undiagnosed (∼16%–22%) until they have a transient ischemic attack or stroke.[12–14] Several predictive models to detect undiagnosed AF have been developed using artificial intelligence (AI) or machine learning (ML).[15–20] If successful in identifying patients at higher risk of AF, earlier access to guideline-recommended oral anticoagulant treatment may minimize the incidence of AF-related stroke and other related CV events.[15–19] However, most models have been developed using local or region-specific populations and include features that are not readily accessible in electronic health records (EHRs) and clinical databases, which limits the ability to operationalize them in the real world. Furthermore, the models have a 5- to 10-year prediction range and potentially less actionable than a

**Address reprint requests and correspondence:** Dr Mohammad Ateya, Pfizer Inc., 66 Hudson Blvd E, New York, NY 10001. E-mail address: Mohammad.Ateya@pfizer.com.

## KEY FINDINGS

- Atrial fibrillation (AF) is a major public health concern associated with significant morbidity and mortality. Prediction models for AF may enable earlier detection and guideline-directed treatment decisions. However, model bias may lead to inaccurate predictions and unintended consequences.

- The 2-year, 10-variable predictive model of undiagnosed AF using common electronic health data (UNAFIED, referred to as UNAFIED-10 in this study) was originally developed using regional data from the Indiana Network for Patient Care. It includes 10 predictor variables: age, sex, body mass index, combined chronic obstructive pulmonary disease (COPD) and obstructive sleep apnea (OSA), heart failure, acute heart disease, kidney disease, shock, albumin, and insurance. This study externally validated and assessed bias of UNAFIED-10 using the nationwide US Optum® electronic health record (EHR) database.

- Findings from the external validation of UNAFIED-10 demonstrated better predictive performance compared with the published UNAFIED-10 model development study (C-statistic 0.850 vs 0.806). These support the model's potential applicability beyond the regional data set used for the original model development. We further report model retraining and optimization undertaken to generate a more generalizable and parsimonious 8-variable predictive model called UNAFIED-8.

- UNAFIED-8 includes 8 predictor variables—age, sex, body mass index, COPD, OSA, kidney disease, heart disease, and heart failure—that are routinely collected and easily extractable from EHR data. It showed consistent good performance in identifying patients at higher risk of AF in the 2 years before their eventual AF diagnosis (C-statistic 0.845; specificity 79%; accuracy 78%; sensitivity 76%). Further, performance disparity analyses showed that UNAFIED-10 and UNAFIED-8 tended to have higher sensitivity in patients with Medicare insurance than in those with other insurance categories, higher precision in White patients than in Asian patients, and higher specificity in patients with other insurance categories than in Medicare beneficiaries.

- UNAFIED-8 is more parsimonious and generalizable for using advanced analytics for AF detection. Future directions include validating the model on data from other US and non-US clinical databases, integrating bias mitigation strategies to ascertain optimal performance and fairness in patient subpopulations, and studying the implementation of the model in clinical settings.

model predicting more imminent risk (eg, in 1–2 years). These predictive models also commonly lack external validation in alternative patient populations or health care systems (limiting transferability) and were not evaluated with bias assessment tools.[21–24]

The clinical applications of AI/ML models may improve CV health outcomes, but several studies have raised concerns about possible model biases leading to inaccurate decisions and harmful results.[14,25–27] The accuracy and fairness of AI/ML predictive models are highly dependent on the data and algorithm used to develop, train, and test them.[28] Model bias includes the likelihood of a model to favor one demographic group over another, which might contribute to unfairness.[29] Algorithmic and population-specific biases may perpetuate sex, racial, or socioeconomic disparities and restrict model generalizability.[30]

The recently published 2-year predictive model of undiagnosed AF using common electronic health data (UNAFIED), referred to as UNAFIED-10 here, was developed using local data from the Indiana Network for Patient Care (INPC), a large health information exchange in Indiana.[20] UNAFIED-10 includes 10 predictor variables commonly available in EHRs: age, sex, body mass index, combined chronic obstructive pulmonary disease (COPD) and obstructive sleep apnea (OSA), heart failure, acute heart disease, kidney disease, shock, albumin, and insurance. It achieved a C-statistic of 0.81 (95% confidence interval [CI] 0.80–0.81) during validation using INPC data.[20] A successful noninterventional proof-of-concept pilot deployment of UNAFIED-10 was launched in the Epic EHR system across all settings of Eskenazi Health.[20] During the proof-of-concept pilot, UNAFIED-10 identified 35.5% (n=7916) of 22,272 patients (aged ≥40 years) at higher risk of AF.

The $CHA_2DS_2$-VASc score is used to calculate stroke risk for patients with AF. Of the 7916 patients with higher AF risk identified during the proof-of-concept deployment of UNAFIED-10, 70% (n=5582) had a $CHA_2DS_2$-VASc score of ≥2.[20] This suggests that they may benefit from guideline-recommended anticoagulant therapy to reduce their risk of stroke if diagnosed with AF.[20,31,32] A 9-month clinical pilot was then conducted using UNAFIED-10 to automatically identify patients with an elevated 2-year risk of AF in a cardiology clinic.[33]

UNAFIED-10 was developed using data mostly from a single US state. The predictors included variables (eg, insurance) that may limit the model's generalizability beyond the regional setting in which it was developed. This study reports the external validation and AI bias assessment of UNAFIED-10 using the nationwide Optum® de-identified electronic health record (Optum EHR) data set. We further report model retraining and optimization undertaken to generate a more generalizable and parsimonious 8-variable predictive model called UNAFIED-8.

## Methods
### Study design, patients, and compliance
This retrospective nested case-control study used de-identified clinical data from the US Optum EHR repository to validate
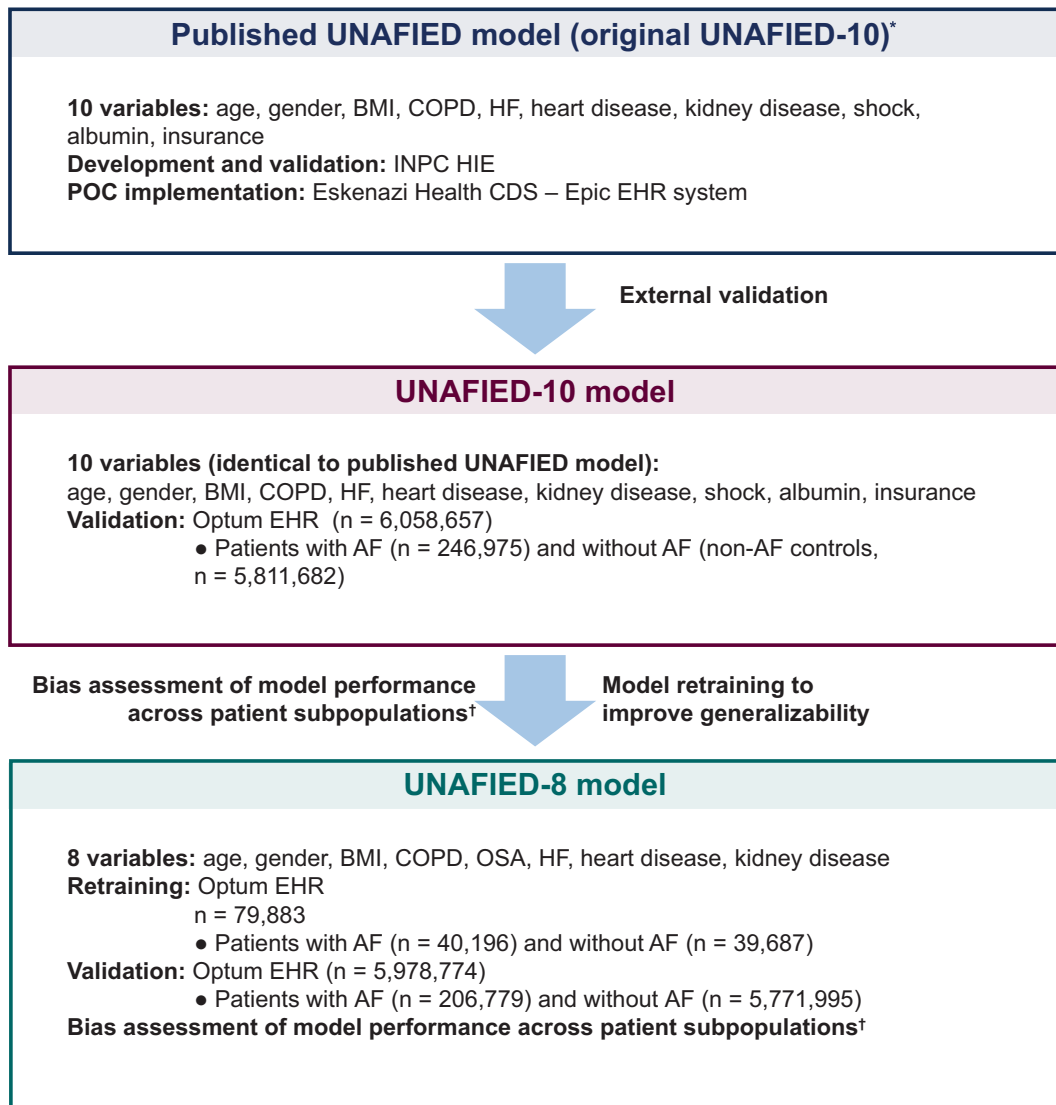
## Published UNAFIED model (original UNAFIED-10)[*]

**10 variables:** age, gender, BMI, COPD, HF, heart disease, kidney disease, shock, albumin, insurance
**Development and validation:** INPC HIE
**POC implementation:** Eskenazi Health CDS – Epic EHR system

**External validation**

## UNAFIED-10 model

**10 variables (identical to published UNAFIED model):**
age, gender, BMI, COPD, HF, heart disease, kidney disease, shock, albumin, insurance
**Validation:** Optum EHR  (n = 6,058,657)
● Patients with AF (n = 246,975) and without AF (non-AF controls, n = 5,811,682)

**Bias assessment of model performance across patient subpopulations[†]**    **Model retraining to improve generalizability**

## UNAFIED-8 model

**8 variables:** age, gender, BMI, COPD, OSA, HF, heart disease, kidney disease
**Retraining:** Optum EHR
n = 79,883
● Patients with AF (n = 40,196) and without AF (n = 39,687)
**Validation:** Optum EHR (n = 5,978,774)
● Patients with AF (n = 206,779) and without AF (n = 5,771,995)
**Bias assessment of model performance across patient subpopulations[†]**

**Figure 1**    Model validation and optimization flowchart from the original UNAFIED-10 model to the new UNAFIED-8 model. AF = atrial fibrillation/atrial flutter; BMI = body mass index; CDS = clinical decision support; COPD = chronic obstructive pulmonary disease; EHR = electronic health record; HF = heart failure; INPC HIE = Indiana Network for Patient Care health information exchange; OSA = obstructive sleep apnea; POC = proof-of-concept; UNAFIED = undiagnosed atrial fibrillation prediction using electronic health data. *Refer to Grout et al.[20] †Patient subpopulations were based on sex, type of health insurance, race, and region.

and optimize the UNAFIED-10 logistic regression model (Figure 1).[20] Optum's longitudinal EHR repository is derived from dozens of 'health care provider organizations in the United States', encompassing all types of health care providers, institutions, or organizations included in the Optum EHR database and contains data from >100 million patients. The data are certified as de-identified by an independent statistical expert following the US Health Insurance Portability and Accountability Act statistical de-identification rules and managed according to Optum customer data use agreements. The original UNAFIED-10 model was reproduced using Optum EHR data for external validation on a national scale followed by bias assessment. The model was further retrained and modified to address bias and improve generalizability. This resulted in a more parsimonious 8-variable predictive model called UNAFIED-8.

This study followed the Transparent Reporting of a Multivariable Prediction Model for Individual Prognosis or Diagnosis (TRIPOD) guidelines (Online Supplemental Table S1).[34]

### Validation of UNAFIED-10 using Optum EHR data

Patients 40 years and older during the outcome period (January 1, 2018–December 31, 2019) were eligible for inclusion if they were active within the Optum EHR database, had ≥2 baseline in-person clinic visits, and no AF diagnoses before January 1, 2018 (start of the outcome period) (Online Supplemental Figure S1). AF was determined whether the patient had an *International Classification of Diseases, Tenth Revision* code of I48.0, I48.1, I48.2, I48.3, I48.4, I48.91, I48.92, I48.11, I48.19, I48.20, or I48.21, an *International Classification of Diseases, Ninth Revision* code of 427.31 or 427.32, or a report of AF or atrial flutter in their medical record (Online Supplemental Table S2). In both UNAFIED-10 and UNAFIED-8 models,
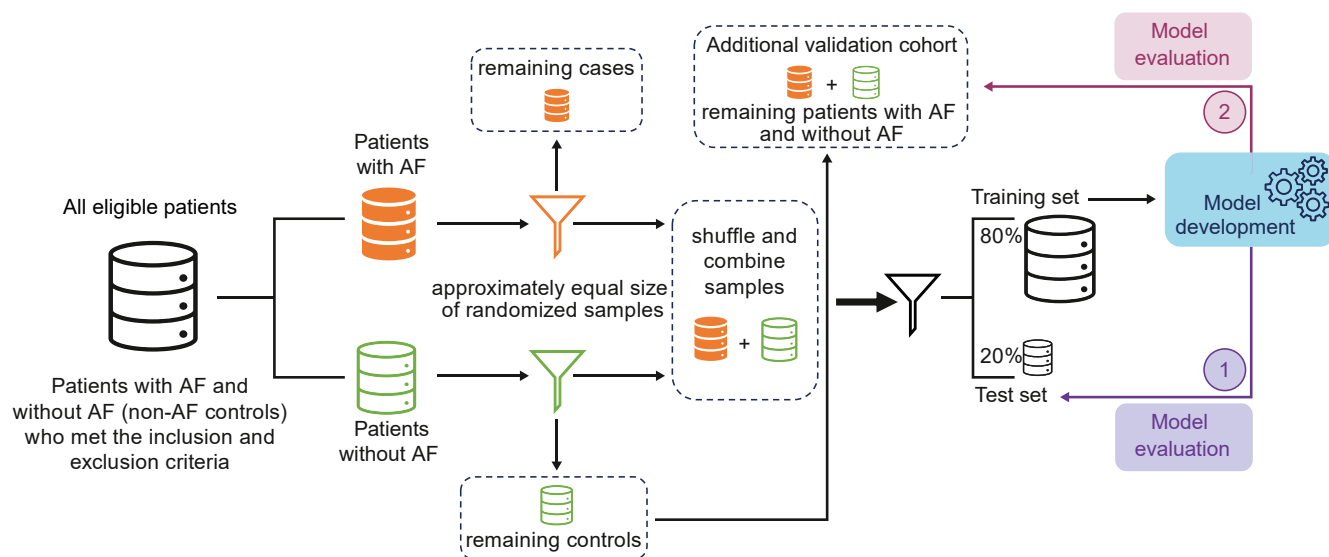
**Figure 2**    Process for training, testing, and validation of the UNAFIED-8 model. AF = atrial fibrillation/atrial flutter; UNAFIED = undiagnosed atrial fibrillation prediction using electronic health data.

the baseline period (January 1, 2016–December 31, 2017) for patients with AF was defined as 2 years before the date of first AF diagnosis (index date). Patients without AF (non-AF controls) had a fixed index date of January 1, 2018.

Parameter estimates for the externally validated UNAFIED-10 model are the same as the original UNAFIED model (Online Supplemental Table S3). Features of UNAFIED-10 were recreated for external validation using Optum EHR data (Online Supplemental Table S4). Two sources—discrete and natural language processing (NLP)–enriched data—were created and randomly sampled to form the final validation cohort, consisting of the full imbalanced data set and a balanced data set from each data source (Online Supplemental Figure S2). Discrete data were derived from structured Optum EHR data (eg, *International Classification of Diseases, Ninth Revision* and *International Classification of Diseases, Tenth Revision* diagnosis codes, demographics, and laboratory values). NLP-enriched data consisted of features from discrete data and NLP records. NLP variables were obtained from Optum proprietary preconfigured NLP concepts, created on broad topics (eg, medications, signs, disease and symptoms, measurements, and observations). Data were harvested from the note fields within the EHR provided to Optum from >50 large health care systems throughout the United States. Variable definitions for the discrete and NLP-enriched data sets are provided in Online Supplemental Table S4. All eligible patients were included in the full imbalanced data set. Random sampling of the full imbalanced data set was done to generate the balanced data sample including ~50% of patients with AF and ~50% of patients without AF.

## Statistical analysis
Selection, coding, transformation of variables, and identification of patients with AF and without AF were based on (or modified from) UNAFIED-10 (Online Supplemental Tables S2 and S4).

Model performance was evaluated using the following metrics: C-statistic (area under the receiver operating characteristic [ROC] curve), accuracy, sensitivity (true positive rate), specificity (true negative rate), negative predictive value, and precision (positive predictive value). AF risk scores were calculated for each patient on the basis of the parameters of the published UNAFIED-10 model (Online Supplemental Table S3), and the same classification cutoff calculated using the Youden's index (0.591) was used to classify patients into higher vs standard risk of undiagnosed AF.[20,35] Analyses were done using Dataiku Data Science Studio Version 11.4.4 (Dataiku, New York, NY) and Python 3.6 (Python Software Foundation, Fredericksburg, VA).

## Bias assessment of model performance across patient subpopulations
UNAFIED models were evaluated for performance bias using the Python library Aequitas audit toolkit version 0.42.0 (University of Chicago, Chicago, IL) and odds ratios.[36] Disparities in performance metrics (sensitivity, precision, and specificity) were analyzed in patient subpopulations based on sex, insurance, race, and region (male, Medicare, White, and US Midwest were used as reference categories). Significant disparity was determined using the 80% rule; it was considered significant if the selected bias disparity ratio is <0.8 or ≥1.2.[37]

## UNAFIED-8 model development and validation
On the basis of external validation and bias assessment findings of UNAFIED-10, the model was retrained, tested, and validated (Figure 2), which subsequently generated the updated UNAFIED-8 model. To create a balanced data set during model development, random samples of patients with AF (n=40,196) and without AF (n=39,687) were drawn from

the full imbalanced data set of patients with AF (n=246,975) and without AF (n=5,811,682). Among the balanced cohorts of patients with AF and without AF, 80% were used for model training whereas the remaining 20% represented a holdout sample for testing. Additional validation was performed on the imbalanced data set after excluding individuals in the balanced data set used for training and testing the model. A total of 24 models, consisting of 8 different versions of UNAFIED-10, with or without 4 features (insurance, shock, albumin, or laboratory values), were tested with 3 different strategies for coding the age variable (Online Supplemental Table S5). Variance inflation factor (VIF) was executed to alleviate multicollinearity among variables (VIF = 5, moderate threshold). Variables with VIF > 5 were dropped. Statistically insignificant variables were also excluded one at a time. The Youden's index for determining the optimal cutoff threshold was calculated for UNAFIED-8. Subsequently, bias assessment of model performance across patient subpopulations was conducted as described earlier. Analyses were performed using Dataiku Data Science Studio Version 11.4.4 and Python 3.6.

## Results

### UNAFIED-10 model validation using Optum EHR data

Among all eligible patients (n=6,058,657; 61% were female), 4.1% had their first AF diagnosis recorded within the outcome period (patients with AF, n=246,975; mean age 72 ± 11 years) whereas 95.9% (patients without AF, n=5,811,682; mean age 60 ± 12 years) did not. The balanced validation cohort included a random sample of 9983 patients: 50.2% in the AF group (mean age 72 ± 11 years; 48.6% were female) and 49.8% in the non-AF control group (mean age 60 ± 12 years; 61.2% were female). Patient characteristics for the validation cohorts using either discrete or NLP-enriched Optum® EHR data are summarized in Table 1 and Online Supplemental Table S6, respectively.

The C-statistic (95% CI) was similar (0.850 [0.845–0.857] vs 0.851 [0.849– 0.851]) (Table 2) between the complete and balanced validation cohorts from Optum EHR discrete data, which was not improved with NLP data enrichment (0.847 vs 0.842) (Online Supplemental Table S7). The ROC curve of the validation cohorts is shown in Figure 3 and Online Supplemental Figure S3. Using the same AF risk threshold score (0.591) as in the original UNAFIED model, validation using Optum EHR data achieved comparable sensitivity (87% vs 86%) and specificity (66% vs 65%) but differences were observed in the accuracy (67% vs 76%) and negative predictive value (99% vs 83%) between the complete and balanced validation data sets (Table 2). Since the positive predictive value is dependent on prevalence, the imbalanced complete validation cohort had a lower positive predictive value than the balanced data set (10% vs 72%). In comparison to the original UNAFIED-10 model validated using INPC data, the externally validated UNAFIED-10 model

showed a higher C-statistic (0.850 [95% CI 0.845–0.857] vs 0.806 [95% CI 0.802–0.810]) (Table 2) at the same classification threshold of 0.591.

### Bias assessment of UNAFIED-10 across patient subpopulations

Disparity analysis for UNAFIED-10 showed some differences in performance, depending on the selected performance metric and patient sex, insurance type, race, or region (Figure 4). Sensitivity (true positive rate) was lowest in patients with unknown insurance compared with Medicare beneficiaries (76% vs 99%). Similarly, disparity in model precision was seen in patients with unknown insurance compared with Medicare beneficiaries (55% vs 74%). Precision was observed to be lowest in patients of Asian vs White (42% vs 66%) race, with an overall lower representation of Asian patients. Specificity was higher in patients of Asian vs White (61% vs 51%) race and lower in those with Medicare (12%) vs all other insurance categories (Medicaid [62%], commercial [56%], or unknown [63%] insurance). Overall, higher positive predictive values were observed in Medicare beneficiaries than in those with all other insurance types and in White patients than in those of all other races.

### UNAFIED-8 development and validation

After the external validation and bias assessment of UNAFIED-10, 8 different versions across 3 different coding strategies for the age variable were tested to reduce the number of features for greater generalizability and to address bias (Online Supplemental Table S5). Of all the 24 models, a final model with 5-year age buckets and comprising 8 variables, called UNAFIED-8, was selected. The insurance variable demonstrated unwanted bias and restricted model use to regions inside the United States; the shock variable had a low prevalence in patients with AF (7%) and without AF (2%); and the albumin variable was not measured for most patients. Furthermore, 5-year age buckets were selected to provide greater granularity. Other modifications in UNAFIED-8 consisted of excluding tachycardia-induced cardiomyopathy *International Classification of Diseases* codes from the definition of the heart failure feature and separating COPD and OSA into 2 features. The new, more generalized, and parsimonious UNAFIED-8 model has 8 predictor variables: age, sex, body mass index, COPD, OSA, kidney disease, heart disease, and heart failure. The variables and their corresponding parameter estimates are provided in Table 3.

The performance of UNAFIED-8 was assessed using 2 different data sets consisting of a test set and a validation cohort (Figure 2). The balanced test data set included 15,977 patients (mean age 66 ± 13 years; 55% were female) with 8016 (50%) and 7961 (50%) patients in the AF and non-AF cohorts, respectively. The imbalanced validation cohort consisted of 5,978,774 patients (61% were female), including 206,779 (3.5%; mean age 72 ± 11 years) who developed AF during the outcome period and 5,771,995 (96.5%; mean age 60 ± 12 years) who did not. The balanced validation cohort

**Table 1** Patient characteristics in UNAFIED-10 and UNAFIED-8 validation cohorts using discrete data from the Optum® EHR

| Characteristic | UNAFIED-10 validation | | | | UNAFIED-8 | | | |
|---|---|---|---|---|---|---|---|---|
| | All eligible patients (N = 6,058,657) | | Balanced data set (n = 9983) | | Balanced data set* (n=79,883) | | Validation cohort† (n = 5,978,774) | |
| | Patients with AF | Patients without AF | Patients with AF | Patients without AF | Patients with AF | Patients without AF | Patients with AF | Patients without AF |
| Number | 246,975 (4.1) | 5,811,682 (95.9) | 5010 (50.2) | 4973 (49.8) | 40,196 (50.3) | 39,687 (49.7) | 206,779 (3.5) | 5,771,995 (96.5) |
| Age (y) | 72.4 ± 11.4 | 59.91 ± 11.8 | 72.3 ± 11.3 | 60.1 ± 11.8 | 72.4 ± 11.4 | 60.0 ± 11.8 | 72.4 ± 11.4 | 59.9 ± 11.8 |
| Sex | | | | | | | | |
| Female | 117,885 (47.7) | 3,549,049 (61.1) | 2433 (48.6) | 3042 (61.2) | 19,215 (47.8) | 24,277 (61.2) | 98,670 (47.7) | 3,524,772 (61.1) |
| Male | 129,090 (52.3) | 2,262,633 (38.9) | 2577 (51.4) | 1931 (38.8) | 20,981 (52.2) | 15,410 (38.8) | 108,109 (52.3) | 2,247,223 (38.9) |
| Body mass index | | | | | | | | |
| Underweight | 5,059 (2.0) | 49,818 (0.9) | 92 (1.8) | 42 (0.8) | 848 (2.1) | 327 (0.8) | 4,211 (2.0) | 49,491 (0.9) |
| Normal | 50,520 (20.5) | 1,070,627 (18.4) | 1043 (20.8) | 928 (18.7) | 8,138 (20.2) | 7,281 (18.3) | 42,382 (20.5) | 1,063,346 (18.4) |
| Overweight | 71,059 (28.8) | 1,636,455 (28.2) | 1432 (28.6) | 1402 (28.2) | 11,642 (29.0) | 11,129 (28.0) | 59,417 (28.7) | 1,625,326 (28.2) |
| Obese | 104,565 (42.3) | 2,255,535 (38.8) | 2148 (42.9) | 1931 (38.8) | 17,005 (42.3) | 15,501 (39.1) | 87,560 (42.3) | 2,240,034 (38.8) |
| Missing | 15,772 (6.4) | 799,247 (13.8) | 295 (5.9) | 670 (13.5) | 2,563 (6.4) | 5,449 13.7) | 13,209 (6.4) | 793,798 (13.8) |
| Heart failure | 73,468 (29.7) | 168,554 (2.9) | 1515 (30.2) | 161 (3.2) | 8,632 (21.5) | 685 (1.7) | 27,961 (13.5) | 103,138 (1.8) |
| Heart disease | 135,967 (55.1) | 713,082 (12.3) | 2769 (55.3) | 619 (12.4) | 22,009 (54.8) | 4,887 (12.3) | 113,958 (55.1) | 708,195 (12.3) |
| Kidney disease | 111,386 (45.1) | 860,065 (14.8) | 2257 (45.1) | 755 (15.2) | 17,985 (44.7) | 5,794 (14.6) | 93,401 (45.2) | 854,271 (14.8) |
| COPD‡ | 80,227 (32.5) | 809,875 (13.9)* | 1603 (32.0) | 700 (14.1) | 7,176 (17.9) | 2,595 (6.5) | 30,259 (14.6) | 370,849 (6.4) |
| OSA‡ | | | | | 5,444 (13.5) | 2,595 (6.5) | 21,533 (10.4) | 373,096 (6.5) |
| Shock§ | 17,859 (7.2) | 121,795 (2.1) | 332 (6.6) | 96 (1.9) | Excluded | Excluded | Excluded | Excluded |
| Albumin: low§ǁ | 47,909 (19.4) | 250,174 (4.3) | 971 (19.4) | 224 (4.5) | Excluded | Excluded | Excluded | Excluded |
| Insurance type§ | | | | | | | | |
| Commercial | 149,652 (60.6) | 4,263,104 (73.4) | 3065 (61.2) | 3727 (74.9) | | | | |
| Medicaid | 12,110 (4.9) | 281,145 (4.8) | 253 (5.0) | 267 (5.4) | | | | |
| Medicare | 71,661 (29.0) | 649,935 (11.2) | 1444 (28.8) | 566 (11.4) | | | | |
| Unknown | 13,552 (5.5) | 617,498 (10.6) | 248 (5.0) | 413 (8.3) | | | | |

Values are presented as mean ± SD or n (%).

AF = atrial fibrillation/atrial flutter; COPD = chronic obstructive pulmonary disease; EHR = electronic health record; OSA = obstructive sleep apnea; UNAFIED = undiagnosed atrial fibrillation prediction using electronic health data.

*The balanced data set is the data set from which we extracted the training and test subsets.

†Additional UNAFIED-8 validation was done on the imbalanced data set excluding the balanced case and control cohorts used in model training and testing.

‡COPD and OSA were combined and considered as 1 variable in the UNAFIED-10 model but were separated into 2 variables in the UNAFIED-8 model.

§In the UNAFIED-8 model, the variables shock, albumin, and insurance were excluded.

ǁIn this study, low albumin is <3.5 g/dL.

**Table 2**    Performance of the different UNAFIED models

| Metric | UNAFIED: INPC HIE (n = 44,772)* | UNAFIED-10: Optum EHR validation | | UNAFIED-8: Optum EHR | |
|---|---|---|---|---|---|
| | | Balanced data set (n = 9983) | All eligible patients (N = 605,8657) | Test set (n = 15,977) | Validation cohort (n = 5,978,774) |
| C-statistic (95% CI)* | 0.806 (0.802–0.810) | 0.851 (0.845–0.857) | 0.850 (0.849–0.851) | 0.853 (0.845–0.857) | 0.845 (0.844–0.846) |
| Sensitivity (%)[‡] | 74.0 | 86.6 | 86.3 | 79.8 | 75.5 |
| Specificity (%)[‡] | 74.0 | 65.0 | 65.7 | 75.3 | 78.6 |
| PPV (%) | Not reported | 71.6 | 9.7 | 77.0 | 11.2 |
| NPV (%) | Not reported | 82.7 | 99.1 | 78.3 | 98.9 |
| Accuracy (%) | Not reported | 75.9 | 66.5 | 77.6 | 78.4 |
| Cutoff for AF diagnosis[‡] | 0.591 | 0.591 | 0.591 | 0.445 | 0.476 |

AF = atrial fibrillation; CI = confidence interval; EHR = electronic health record; HIE = health information exchange; INPC = Indiana Network for Patient Care; NPV = negative predictive value; PPV = positive predictive value; UNAFIED = undiagnosed atrial fibrillation prediction using electronic health data.
*Values are from the original UNAFIED model validation phase.[20]
[†]The C-statistic (95% CI) in the development phase of the original UNAFIED model was 0.796 (0.792–0.799; n = 53,552), which achieved 74% sensitivy and 74% specificity.[20]
[‡]The optimal cutoff value was determined using the Youden's index (sensitivity + specificity – 1).

included 79,883 patients (54% were female; 50% of patients with AF and 50% of patients without AF). Patient characteristics are summarized in Table 1.

Using the test set, the optimal cutoff score was 0.445, with a Youden's index of 0.5508 and a C-statistic of 0.853 (Figure 3 and Online Supplemental Figure S4), which resulted in 80% sensitivity, 75% specificity, 77% positive predictive value, 78% negative predictive value, and 78% accuracy (Table 2). Another version of the model was created by excluding laboratory values from the definitions of UNAFIED-8 features to assess the usability of the model by organizations with no access to laboratory values. The parameter estimates and ROC curve for UNAFIED-8 without laboratory values are presented in Online Supplemental Table S8 and Online Supplemental Figure S5.

For the UNAFIED-8 imbalanced complete validation cohort, the optimal cutoff score was 0.476, with a Youden's index of 0.5404 (Table 2 and Online Supplemental Figure S5). The C-statistic was 0.845, which resulted in 76% sensitivity, 79% specificity, 78% accuracy, 11% positive predictive value, and 99% negative predictive value.
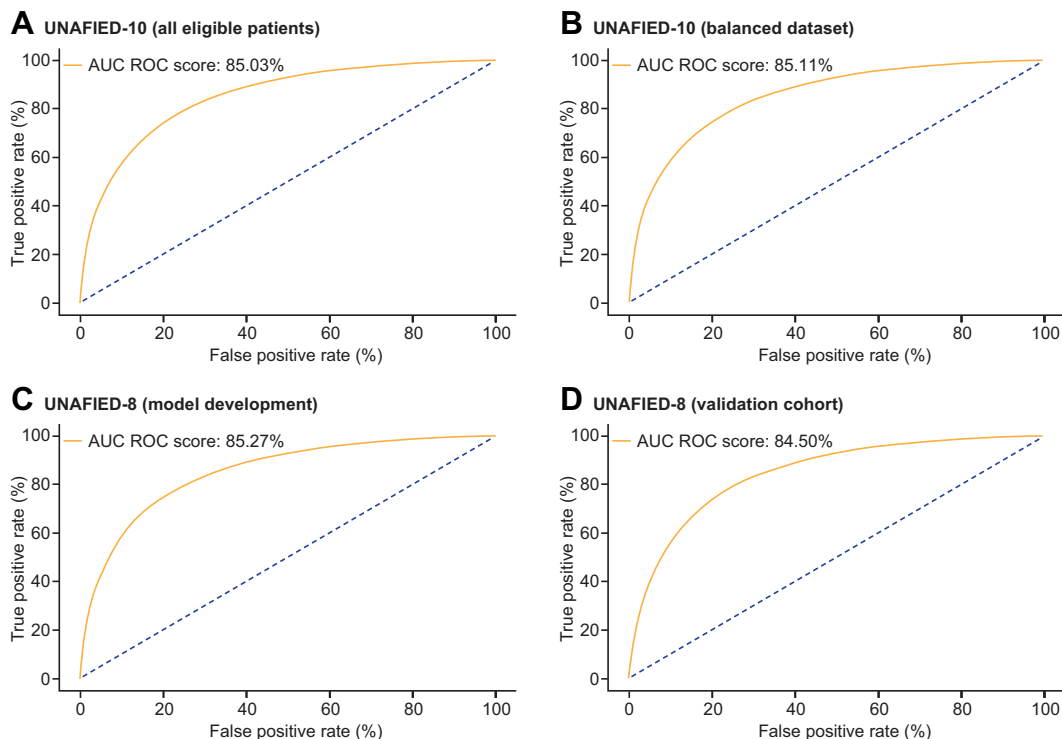


**Figure 3**    ROC curves for the UNAFIED-10 model validation cohort: (**A**) complete and (**B**) balanced discrete data sets. ROC curves for the UNAFIED-8 model: (**C**) development and (**D**) validation cohorts. AUC ROC = area under the receiver operating characteristic curve; UNAFIED = undiagnosed atrial fibrillation prediction using electronic health data.
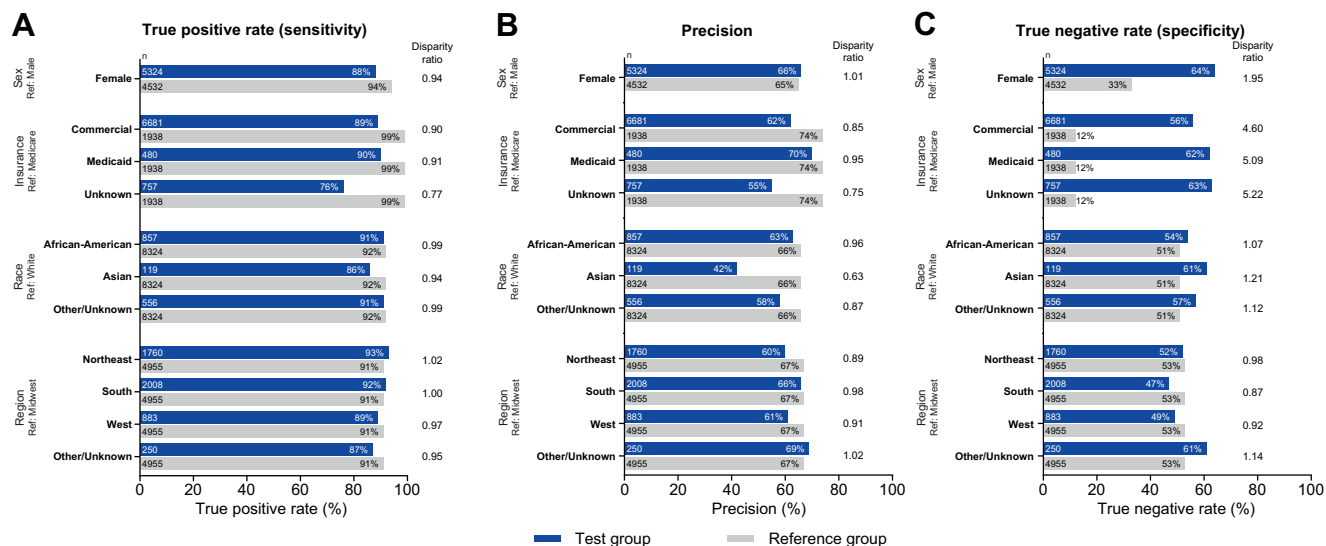
**Figure 4** Bias assessment of UNAFIED-10: (**A**) sensitivity, (**B**) precision, and (**C**) specificity. Ref = reference group; UNAFIED = undiagnosed atrial fibrillation prediction using electronic health data.

Overall, UNAFIED-8 had a higher C-statistic (95% CI) compared with the original UNAFIED-10 model (0.845 [0.844–0.846] vs 0.806 [0.802–0.810]) and consistent with UNAFIED-10 that was validated using Optum EHR data (0.845 vs 0.850). Moreover, the UNAFIED-8 validation cohort showed higher specificity (79% vs 66%) and accuracy (78% vs 67%) but lower sensitivity (76% vs 86%) compared with the UNAFIED-10 Optum EHR complete validation cohort (Table 2). Performance of the

UNAFIED-8 models with or without laboratory values were similar (C-statistic 0.845 vs 0.847) (Table 2 and Online Supplemental Table S9).

## Bias assessment of UNAFIED-8 across patient subpopulations

Findings from the UNAFIED-8 model performance disparity analysis showed differences in sensitivity between patients with Medicaid (65%), commercial (69%), or an unknown

**Table 3** UNAFIED-8 model variables and parameter estimates

| Parameter*† | UNAFIED-8 model development | | |
| | Description | Estimate | Odds ratio (95% CI) |
| --- | --- | --- | --- |
| Intercept | | −1.8785 | |
| Age | ≥40 to <45 y | −0.3460 | 0.65 (0.58–0.72) |
| | ≥50 to <55 y | 0.3375 | 1.57 (1.44–1.71) |
| | ≥55 to <60 y | 0.6239 | 2.35 (2.17–2.55) |
| | ≥60 to <65 y | 0.9870 | 3.64 (3.37–3.94) |
| | ≥65 to <70 y | 1.3666 | 5.55 (5.13–6.00) |
| | ≥70 to <75 y | 1.7449 | 8.70 (8.04–9.41) |
| | ≥75 to <80 y | 2.1015 | 12.85 (11.83–13.95) |
| | ≥80 to <85 y | 2.4137 | 18.42 (16.82–20.16) |
| | ≥85 y | 2.9047 | 30.43 (27.69–33.44) |
| Heart disease | Present | 1.3388 | 8.62 (8.31–8.93) |
| Body mass index | Missing | −0.3607 | 0.45 (0.43–0.47) |
| | <18.5 kg/m² | 0.4660 | 2.48 (2.18–2.82) |
| | >30 kg/m² | 0.2515 | 1.05 (1.01–1.09) |
| Sex | Female | −0.4060 | 0.58 (0.57–0.60) |
| Heart failure* | Present | 1.5826 | 15.57 (14.39–16.85) |
| Kidney disease | Present | 0.5493 | 4.74 (4.58–4.90) |
| OSA* | Present | 0.3954 | 2.24 (2.13–2.35) |
| COPD* | Present | 0.4065 | 3.11 (2.96–3.26) |

CI = confidence interval; COPD = chronic obstructive pulmonary disorder; *ICD-9* = International Classification of Diseases, Ninth Revision; ICD-10 = International Classification of Diseases, Tenth Revision; OSA = obstructive sleep apnea; UNAFIED = undiagnosed atrial fibrillation prediction using electronic health data.
*The UNAFIED-8 model differed from the original UNAFIED-10 model in defining heart failure and COPD. In UNAFIED-8, the tachycardia codes (*ICD-9* code 425.X; *ICD-10* code I42.X) were removed from heart failure, and OSA (*ICD-9* codes 32723, 32720, 32729, 78051, 78053, 78057; *ICD-10* codes G4733, G4730, G4739) was separated from COPD.[20]
†Reference parameters: age ≥45 to <50 y, body mass index 18.5–24.9 kg/m², male, no diagnosis of heart disease, COPD, heart failure, or kidney disease.
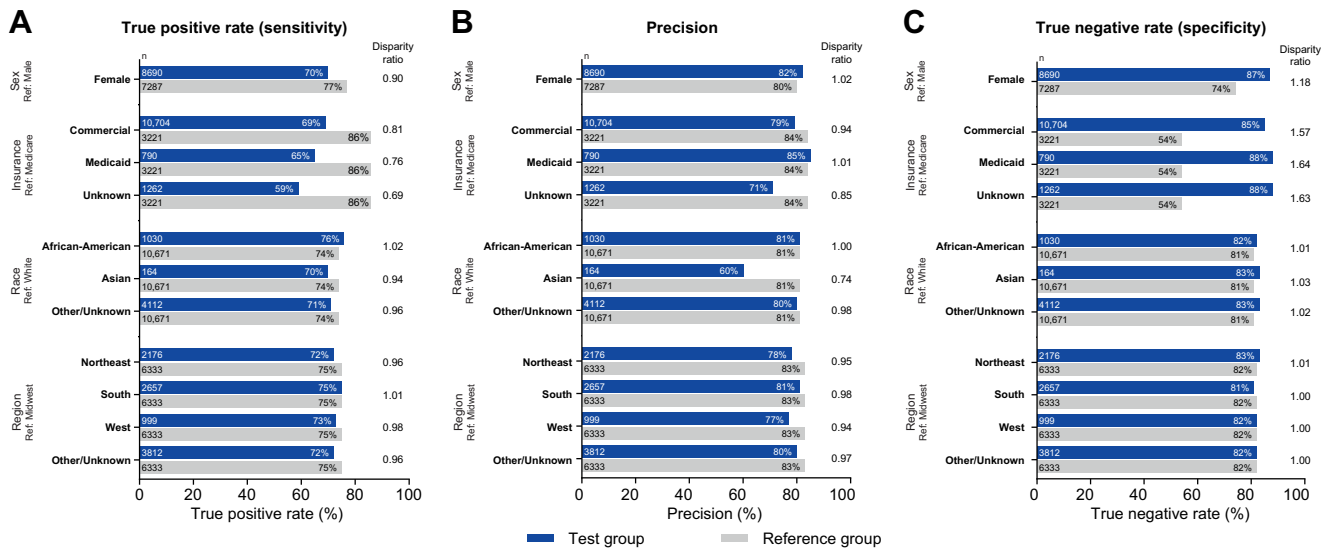
**Figure 5** Bias assessment of UNAFIED-8: (**A**) sensitivity, (**B**) precision, and (**C**) specificity. Ref = reference group; UNAFIED = undiagnosed atrial fibrillation prediction using electronic health data.

(59%) insurance compared with Medicare beneficiaries (86%) (Figure 5). Similar to the disparity assessment results for the externally validated UNAFIED-10 model, precision was lowest in patients of Asian (60%) compared with White (81%) race, and there was a lower representation of Asian patients. Furthermore, specificity was lower in Medicare beneficiaries (54%) than in recipients of all other insurance types (commercial [85%], Medicaid [88%], and unknown [88%]).

## Discussion

This study used a large longitudinal national data set with >100 million patients to externally validate and assess bias of the previously published 10-variable clinical prediction model for undiagnosed AF (UNAFIED-10), which was developed using a diverse patient population mostly from a single US state.[20] UNAFIED-10 was updated to a more parsimonious and generalizable 8-variable model (UNAFIED-8).

External validation is important to assess a model's reproducibility and transportability in new settings. However, most models are rarely externally validated or used in clinical practice.[38] The external validation of UNAFIED-10 using a national data set showed better predictive performance compared with the published UNAFIED-10 model development study (C-statistic 0.850 vs 0.806). Findings from external validation suggest that UNAFIED-10 has acceptable predictive performance in settings beyond the regional INPC data set used in its original model development.

Bias in medical AI may affect patient care and could arise from several factors (eg, disparities in institutional or health care practices, sampling, demographic representation in data sets, and algorithm bias).[21,22,39] However, bias assessment is not necessarily a standard part of model development. Findings from performance disparity analyses showed that both UNAFIED-10 and UNAFIED-8 models tended to have higher sensitivity in patients with Medicare insurance than in those

with other insurance categories, higher precision in White patients than in Asian patients, and higher specificity in patients with other insurance categories than in Medicare beneficiaries. Disparities in model performance may be attributed to Medicare beneficiaries being generally older (aged ≥65 years), which is associated with higher AF prevalence, and to underrepresentation of Asian patients in the data sets. By assessing and identifying biases in model performance within patient subpopulations, mitigation strategies can be implemented to ascertain optimal performance and fairness. Further postprocessing to recalibrate or tune the model to a specific population may reduce bias and improve fairness.[29]

The updated UNAFIED-8 model has 8 predictor variables (age, sex, body mass index, COPD, OSA, kidney disease, heart disease, and heart failure) that are routinely collected and easily extractable from EHR data (Figure 1 and Table 3). Overall, UNAFIED-8 demonstrated consistent good performance in identifying patients at higher risk of AF in the 2 years before their eventual AF diagnosis (C-statistic 0.845; specificity 79%; accuracy 78%; sensitivity 76%). Variables (insurance, shock, albumin level) that restricted the generalizability of UNAFIED-10 were eliminated in UNAFIED-8. We also presented a version of the UNAFIED-8 model without laboratory values in the feature definitions and with similar predictive performance to allow its potential implementation at health care organizations without full access to laboratory data (eg, payers). UNAFIED-8 may be used alone or with other modalities in the future. When the model is implemented, it would be dynamic, and the risk scores will be automatically updated as new relevant diagnoses or laboratory results are documented. Furthermore, when appropriately integrated into workflows, UNAFIED-8 may be a useful tool for health care professionals to better identify patients for AF screening, potentially leading to earlier clinical interventions and reducing the prevalence of AF-related stroke.

Clinical practice guidelines provide varying endorsements of population-based AF screening. The European Society of Cardiology recommends opportunistic screening for AF by pulse taking or electrocardiogram (ECG) rhythm strip in individuals 65 years and older.[40] It also recommends systematic ECG screening to detect AF in individuals 75 years and older or at high risk of stroke.[40] However, the US Preventive Services Task Force indicated, in their 2022 publication, that current evidence is insufficient to assess the balance of benefits and harms of screening for AF in asymptomatic adults 50 years and older.[41] Multiple randomized and nonrandomized clinical trials on population-based AF screening showed increased rates for detecting AF in previously undiagnosed patients, but most trials did not demonstrate reduction in stroke, systemic embolism, or death; and a number of these trials were not powered to detect such differences.[41,42] Moreover, the STROKESTOP trial in patients aged 75–76 years in Sweden demonstrated a small net benefit of AF screening compared with standard of care (no screening) for the combined primary end point of ischemic or hemorrhagic stroke, systemic embolism, bleeding leading to hospitalization, and all-cause death after a median follow-up of 6.9 years (hazard ratio 0.96 [95% CI 0.92–1.00]; $P = .045$).[43] This study indicated that screening is safe and beneficial in older populations.[43] However, the potential benefits and harms of screening and earlier treatment must be balanced. Screening may contribute to the reduction in the risk of AF-related stroke or systemic embolism and reduce AF-related morbidity and mortality through early identification of patients at high AF risk, thereby prompting timely treatment decisions. However, false-positive diagnoses from screening may lead to unnecessary oral anticoagulation therapy, potentially increasing bleeding risk. Abnormal screening results also require additional confirmatory testing, potentially contributing to patient anxiety and higher costs.[24,40]

There is growing interest in understanding the impact of using clinical prediction scores and AI-based models on targeting individuals at highest AF risk and facilitating more efficient AF screening. Different predictive approaches for evaluating AF risk have been proposed, including clinical risk scores (eg, Cohorts for Heart and Aging Research in Genomic Epidemiology-AF, Atherosclerosis Risk in Communities, EHR-AF, and C2HEST), AI models using clinical variables to predict AF, and ML models using raw 12-lead ECGs or a combination of raw 12-lead ECGs and clinical variables.[17,42,44,45] Compared with these existing models, UNAFIED may provide a simpler option to estimate AF risk with acceptable predictive performance.

In the present study, our analyses were limited to patients with available clinical data in the Optum EHR. One of the advantages of using EHR data is the ability to include patients regardless of their insurance coverage status or insurance type. We used a single instance of documentation for the outcome and predictor variables as proxy for the presence of these variables. EHR data are collected for clinical care and generally may have inherent limitations (eg, the possibility of entry and coding errors, missing information, and inconsistencies in data collection and reporting).[46] Moreover, model validation was conducted with a retrospective design using patient cohorts with known outcomes and the results may be different when the model is used to prospectively predict AF risk in patients with unknown status. Also, we are unable to recommend specific interventions for patients identified as having a higher risk of AF, and additional research is needed to understand the best approach of managing these patients.

## Conclusion

UNAFIED-10, developed using US regional patient data, displayed consistent performance in a large national data set (Optum EHR). The model's ability to prospectively predict 2-year AF occurrence may guide earlier AF detection and guideline-recommended treatment decisions. UNAFIED-8 is more parsimonious and generalizable for using advanced analytics for AF detection. Future directions include UNAFIED-8 model validation on data from other US and non-US clinical databases and studying the implementation of the model in clinical settings.

## Acknowledgments

## References

1. Odutayo A, Wong CX, Hsiao AJ, Hopewell S, Altman DG, Emdin CA. Atrial fibrillation and risks of cardiovascular disease, renal disease, and death: systematic review and meta-analysis. BMJ 2016;354:i4482.
2. Sanoski CA. Clinical, economic, and quality of life impact of atrial fibrillation. J Manag Care Pharm 2009;15:S4–S9.

3. Wolf PA, Abbott RD, Kannel WB. Atrial fibrillation as an independent risk factor for stroke: the Framingham Study. Stroke 1991;22:983–988.

4. Krijthe BP, Kunst A, Benjamin EJ, et al. Projections on the number of individuals with atrial fibrillation in the European Union, from 2000 to 2060. Eur Heart J 2013;34:2746–2751.

5. Miyasaka Y, Barnes ME, Gersh BJ, et al. Secular trends in incidence of atrial fibrillation in Olmsted County, Minnesota, 1980 to 2000, and implications on the projections for future prevalence. Circulation 2006;114:119–125.

6. Lippi G, Sanchis-Gomar F, Cervellin G. Global epidemiology of atrial fibrillation: an increasing epidemic and public health challenge. Int J Stroke 2021; 16:217–221.

7. Colilla S, Crow A, Petkun W, Singer DE, Simon T, Liu X. Estimates of current and future incidence and prevalence of atrial fibrillation in the U.S. adult population. Am J Cardiol 2013;112:1142–1147.

8. Caplan Z. U.S. older population grew from 2010 to 2020 at fastest rate since 1880 to 1890. https://www.census.gov/library/stories/2023/05/2020-census-united-states-older-population-grew.html. Accessed November 14, 2023.

9. Turakhia MP, Shafrin J, Bognar K, et al. Estimated prevalence of undiagnosed atrial fibrillation in the United States. PLoS One 2018;13:e0195088.

10. Turakhia MP, Guo JD, Keshishian A, et al. Contemporary prevalence estimates of undiagnosed and diagnosed atrial fibrillation in the United States. Clin Cardiol 2023;46:484–493.

11. Boriani G, Laroche C, Diemberger I, et al. Asymptomatic atrial fibrillation: clinical correlates, management, and outcomes in the EORP-AF Pilot General Registry. Am J Med 2015;128:509–518.e502.

12. Jaakkola J, Mustonen P, Kiviniemi T, et al. Stroke as the first manifestation of atrial fibrillation. PLoS One 2016;11:e0168010.

13. Sposato LA, Chaturvedi S, Hsieh CY, Morillo CA, Kamel H. Atrial fibrillation detected after stroke and transient ischemic attack: a novel clinical concept challenging current views. Stroke 2022;53:e94–e103.

14. Garzon-Siatoya WT, Morales-Lara AC, Adedinsewo DA. Artificial intelligence solutions for cardiovascular disease detection and management in women: promise and perils. Cardiovasc Innov Appl 2023;8.

15. Sivanandarajah P, Wu H, Bajaj N, Khan S, Ng FS. Is machine learning the future for atrial fibrillation screening? Cardiovasc Digit Health J 2022;3:136–145.

16. Schnabel RB, Sullivan LM, Levy D, et al. Development of a risk score for atrial fibrillation (Framingham Heart Study): a community-based cohort study. Lancet 2009;373:739–745.

17. Alonso A, Krijthe BP, Aspelund T, et al. Simple risk model predicts incidence of atrial fibrillation in a racially and geographically diverse population: the CHARGE-AF consortium. J Am Heart Assoc 2013;2:e000102.

18. Chamberlain AM, Agarwal SK, Folsom AR, et al. A clinical risk score for atrial fibrillation in a biracial prospective cohort (from the Atherosclerosis Risk in Communities [ARIC] study). Am J Cardiol 2011;107:85–91.

19. Aronson D, Shalev V, Katz R, Chodick G, Mutlak D. Risk score for prediction of 10-year atrial fibrillation: a community-based study. Thromb Haemost 2018; 118:1556–1563.

20. Grout RW, Hui SL, Imler TD, et al. Development, validation, and proof-of-concept implementation of a two-year risk prediction model for undiagnosed atrial fibrillation using common electronic health data (UNAFIED). BMC Med Inform Decis Mak 2021;21:112.

21. Kerasidou A. Ethics of artificial intelligence in global health: explainability, algorithmic bias and trust. J Oral Biol Craniofac Res 2021;11:612–614.

22. Norori N, Hu Q, Aellen FM, Faraci FD, Tzovara A. Addressing bias in big data and AI for health care: a call for open science. Patterns (N Y) 2021;2:100347.

23. Gulati G, Upshaw J, Wessler BS, et al. Generalizability of cardiovascular disease clinical prediction models: 158 independent external validations of 104 unique models. Circ Cardiovasc Qual Outcomes 2022;15:e008487.

24. Khurshid S, Mars N, Haggerty CM, et al. Predictive accuracy of a clinical and genetic risk model for atrial fibrillation. Circ Genom Precis Med 2021;14:e003355.

25. Isaksen JL, Baumert M, Hermans ANL, Maleckar M, Linz D. Artificial intelligence for the detection, prediction, and management of atrial fibrillation. Herzschrittmacherther Elektrophysiol 2022;33:34–41.

26. Siontis KC, Yao X, Pirruccello JP, Philippakis AA, Noseworthy PA. How will machine learning inform the clinical care of atrial fibrillation? Circ Res 2020; 127:155–169.

27. Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm used to manage the health of populations. Science 2019;366:447–453.

28. Castelnovo A, Crupi R, Greco G, Regoli D, Penco IG, Cosentini AC. A clarification of the nuances in the fairness metrics landscape. Sci Rep 2022;12:4209.

29. Fletcher RR, Nakeshimana A, Olubeko O. Addressing fairness, bias, and appropriate use of artificial intelligence and machine learning in global health. Front Artif Intell 2020;3:561802.

30. Yang J, Soltan AAS, Eyre DW, Clifton DA. Algorithmic fairness and bias mitigation for clinical machine learning with deep reinforcement learning. Nat Mach Intell 2023;5:884–894.

31. January CT, Wann LS, Calkins H, et al. 2019 AHA/ACC/HRS focused update of the 2014 AHA/ACC/HRS guideline for the management of patients with atrial fibrillation: a report of the American College of Cardiology/American Heart Association Task Force on Clinical Practice Guidelines and the Heart Rhythm Society in collaboration with the Society of Thoracic Surgeons. Circulation 2019; 140:e125–e151.

32. Lip GY, Nieuwlaat R, Pisters R, Lane DA, Crijns HJ. 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;137:263–272.

33. Rajkumar J, Direnzo B, Walroth T, et al. Identification of patients at higher risk of having atrial fibrillation (AF) using an electronic health record (EHR) predictive model. J Am Coll Cardiol 2023;81:1737.

34. Collins GS, Reitsma JB, Altman DG, Moons KG. Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD). Ann Intern Med 2015;162:735–736.

35. Youden WJ. Index for rating diagnostic tests. Cancer 1950;3:32–35.

36. Saleiro P, Kuester B, Hinkson L, et al. Aequitas: a bias and fairness audit toolkit v2. arXiv published online ahead of print April 29, 2019; https://doi.org/10.48550/arXiv.1811.05577.

37. Feldman MF, Friedler S, Moeller J, Scheidegger C, Venkatasubramanian S. Certifying and removing disparate impact. arXiv published online ahead of print July 16, 2015; https://doi.org/10.48550/arXiv.1412.3756.

38. Ramspek CL, Jager KJ, Dekker FW, Zoccali C, van Diepen M. External validation of prognostic models: what, why, how, when and where? Clin Kidney J 2021;14:49–58.

39. Vokinger KN, Feuerriegel S, Kesselheim AS. Mitigating bias in machine learning for medicine. Commun Med (Lond) 2021;1:25.

40. Hindricks G, Potpara T, Dagres N, et al. 2020 ESC Guidelines for the diagnosis and management of atrial fibrillation developed in collaboration with the European Association for Cardio-Thoracic Surgery (EACTS): the Task Force for the diagnosis and management of atrial fibrillation of the European Society of Cardiology (ESC) developed with the special contribution of the European Heart Rhythm Association (EHRA) of the ESC. Eur Heart J 2021;42:373–498.

41. Davidson KW, Mangione C, Ogedegbe G. US Preventive Services Task Force recommendation statement on screening for atrial fibrillation—reply. JAMA 2022;327:2022.

42. Pipilas D, Friedman SF, Khurshid S. The use of artificial intelligence to predict the development of atrial fibrillation. Curr Cardiol Rep 2023;25:381–389.

43. Svennberg E, Friberg L, Frykman V, Al-Khalili F, Engdahl J, Rosenqvist M. Clinical outcomes in systematic screening for atrial fibrillation (STROKESTOP): a multicentre, parallel group, unmasked, randomised controlled trial. Lancet 2021;398:1498–1506.

44. Raghunath S, Pfeifer JM, Ulloa-Cerna AE, 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;143:1287–1298.

45. Attia ZI, Noseworthy PA, Lopez-Jimenez F, et al. An artificial intelligence-enabled ECG algorithm for the identification of patients with atrial fibrillation during sinus rhythm: a retrospective analysis of outcome prediction. Lancet 2019;394:861–867.

46. Bowman S. Impact of electronic health record systems on information integrity: quality and safety implications. Perspect Health Inf Manag 2013;10:1c.