

# New artificial intelligence prediction model using serial prothrombin time international normalized ratio measurements in atrial fibrillation patients on vitamin K antagonists: GARFIELD-AF

**Shinichi Goto<sup>1</sup>, Shinya Goto<sup>2\*</sup>, Karen S. Pieper<sup>3</sup>, Jean-Pierre Bassand<sup>3,4</sup>, Alan John Camm<sup>5</sup>, David A. Fitzmaurice<sup>6</sup>, Samuel Z. Goldhaber<sup>7</sup>, Sylvia Haas<sup>8</sup>, Alexander Parkhomenko<sup>9</sup>, Ali Oto<sup>10</sup>, Frank Misselwitz<sup>11</sup>, Alexander G.G. Turpie<sup>12</sup>, Freek W.A. Verheugt<sup>13</sup>, Keith A.A. Fox<sup>14</sup>, Bernard J. Gersh<sup>15</sup>, and Ajay K. Kakkar<sup>3,16</sup>; for the GARFIELD-AF Investigators**

<sup>1</sup>Department of Cardiology, Keio University School of Medicine, Shinanomachi 35, Shinjuku 160-8582, Tokyo, Japan; <sup>2</sup>Department of Medicine (Cardiology), Tokai University School of Medicine, 143 Shimokasuya, Isehara, Kanagawa 259-1193, Japan; <sup>3</sup>Department of Clinical Research, Thrombosis Research Institute, Emmanuel Kaye Building, Manresa Road, Chelsea, London SW3 6LR, UK; <sup>4</sup>Department of Cardiology, University of Besançon Boulevard Fleming, 25000 Besançon, France; <sup>5</sup>Cardiology Clinical Academic Group, Molecular & Clinical Sciences Institute, St. George's University of London, Cranmer Terrace, Tooting, London, UK; <sup>6</sup>Department of Cardio-respiratory Primary Care, Warwick Medical School, University of Warwick, Coventry, CV4 7AL, UK; <sup>7</sup>Department of Medicine, Harvard Medical School, Brigham and Women's Hospital, 75 Francis Street, Boston, MA 02115, USA; <sup>8</sup>Formerly Klinikum rechts der Isar, Technical University of Munich, Normannenstr. 34a, Munich 80333, Germany; <sup>9</sup>National Scientific Center, Strazhesko Institute of Cardiology, 5 Narodnogo Opolcheniya Street, Kiev 03680, Ukraine; <sup>10</sup>Department of Cardiology, Memorial Ankara Hospital, Sihhiye, 06100, Ankara, Turkey; <sup>11</sup>Therapeutic areas Thrombosis & Hematology, Bayer AG, Müllerstraße 178, 13353 Berlin, Germany; <sup>12</sup>Department of Medicine, McMaster University, 237 Barton St E Hamilton, Ontario L8L 2X2, Canada; <sup>13</sup>Department of Cardiology, Onze Lieve Vrouwe Gasthuis (OLVG), Oosterpark 9, NL-1091-AC Amsterdam, Netherlands; <sup>14</sup>Edinburgh Centre for Cardiovascular Science, University of Edinburgh, Queen's Medical Research Institute, 47 Little France Crescent, Edinburgh EH16 4TJ, UK; <sup>15</sup>Department of Cardiovascular Medicine, Mayo Clinic College of Medicine, Rochester, MN, USA; and <sup>16</sup>Department of Surgery, University College London, Gower St, Bloomsbury, London WC1E 6BT, UK

Received 23 October 2019; revised 14 November 2019; editorial decision 28 November 2019; accepted 5 December 2019

## Aims

Most clinical risk stratification models are based on measurement at a single time-point rather than serial measurements. Artificial intelligence (AI) is able to predict one-dimensional outcomes from multi-dimensional datasets. Using data from Global Anticoagulant Registry in the Field (GARFIELD)-AF registry, a new AI model was developed for predicting clinical outcomes in atrial fibrillation (AF) patients up to 1 year based on sequential measures of prothrombin time international normalized ratio (PT-INR) within 30 days of enrolment.

## Methods and results

Patients with newly diagnosed AF who were treated with vitamin K antagonists (VKAs) and had at least three measurements of PT-INR taken over the first 30 days after prescription were analysed. The AI model was constructed with multilayer neural network including long short-term memory and one-dimensional convolution layers. The neural network was trained using PT-INR measurements within days 0–30 after starting treatment and clinical outcomes over days 31–365 in a derivation cohort (cohorts 1–3;  $n = 3185$ ). Accuracy of the AI model at predicting major bleed, stroke/systemic embolism (SE), and death was assessed in a validation cohort (cohorts 4–5;  $n = 1523$ ). The model's c-statistic for predicting major bleed, stroke/SE, and all-cause death was 0.75, 0.70, and 0.61, respectively.

\* Corresponding author. Tel: +81-463-93-1121, Fax: +81-463-93-6679, Email: [sgoto3@mac.com](mailto:sgoto3@mac.com)

© The Author(s) 2019. Published by Oxford University Press on behalf of the European Society of Cardiology.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact [journals.permissions@oup.com](mailto:journals.permissions@oup.com)

## Conclusions

Using serial PT-INR values collected within 1 month after starting VKA, the new AI model performed better than time in therapeutic range at predicting clinical outcomes occurring up to 12 months thereafter. Serial PT-INR values contain important information that can be analysed by computer to help predict adverse clinical outcomes.

## Keywords

Atrial fibrillation • Artificial intelligence • Machine learning

## Introduction

In chronic diseases such as atrial fibrillation (AF) risk stratification using prediction models is useful for clinical decision-making. Several models predict clinical events such as stroke and bleeding.<sup>1–3</sup> The CHA<sub>2</sub>DS<sub>2</sub>-VASc and HAS-BLED scores are widely used to select suitable AF patients for oral anticoagulation (OAC).<sup>4–6</sup> However, some of the variables in these scoring systems are not consistently related to outcomes.<sup>7</sup> Novel machine learning technology has facilitated the development of more accurate models such as the Global Anticoagulant Registry in the Field (GARFIELD)-AF risk model.<sup>8</sup> However, these models incorporate data obtained at a single time-point, baseline. Although computers can process multi-dimensional data such as changes of variables over time, few models have used these inputs to predict future clinical events.<sup>9,10</sup>

Vitamin K antagonists (VKAs) continue to be prescribed for the prevention of stroke in patients with AF, despite the more recent introduction of non-VKA oral anticoagulants (NOACs).<sup>11,12</sup> The VKAs are the only recommended choice of OAC for AF patients with haemodynamically overt mitral stenosis and mechanical heart valve. Clinicians adjust the dose of VKA based on an individual patient's prothrombin time international normalized ratio (PT-INR) at each visit. Time in therapeutic range (TTR) is widely used to standardize the effects of VKA therapy over periods beyond 6 months.<sup>13–17</sup> Various bleeding risk scores feature a TTR component to enhance accuracy,<sup>18</sup> and TTR has predictive power for thrombotic and bleeding events.<sup>19,20</sup> However, information on serial changes in PT-INR during early-phase VKA therapy, which may reflect many occult clinical characteristics of patients such as genotype,<sup>21,22</sup> concomitant medications,<sup>23</sup> and lifestyle,<sup>24</sup> were not included in these TTR-based models.

Advances in artificial intelligence (AI) using recurrent neural networks allow the identification and translation of multi-dimensional data including time-series data directly into meaningful models.<sup>25</sup> Herein, we describe a new AI model for predicting clinical outcomes over 31–365 days after patient enrolment. The model evaluates serially measured PT-INR within the first 30 days of treatment only without other clinical parameters, using data from the largest multinational prospective registry in AF, GARFIELD-AF. The predictive accuracy of the AI model was compared with that of TTR. The working hypothesis was to test whether serially measured PT-INR in early phase can provide information to predict future clinical events.

## Methods

### Design

The AI model was derived from prospective GARFIELD-AF data gathered in adults with newly diagnosed AF.<sup>26</sup> Three independent AI models

were developed with the same composite of neural network structure with multi-dimensional patient-level PT-INR values obtained within the first 30 days after starting treatment. The model tabulated the clinical events of major bleed, ischaemic stroke/systemic embolism (SE), and death occurring within days 31–365.

### Registry population

The GARFIELD-AF is an ongoing, international, prospective registry of newly diagnosed patients with AF at risk of stroke. The study design, baseline characteristics, and main results have been published.<sup>26–29</sup> Eligible patients were adults aged >18 years who had been diagnosed with non-valvular AF within the previous 6 weeks and had at least one risk factor for stroke as judged by the investigator. Risk factors were not pre-specified in the protocol. Any use of antithrombotic agents was shared decision between clinicians and patients only. Patients with a transient reversible cause of AF and those for whom follow-up was not envisaged were excluded. The present analysis was conducted in patients enrolled in GARFIELD-AF cohorts 1–5 between March 2010 and August 2016. Data were extracted from the study database in November 2017.

### Study population

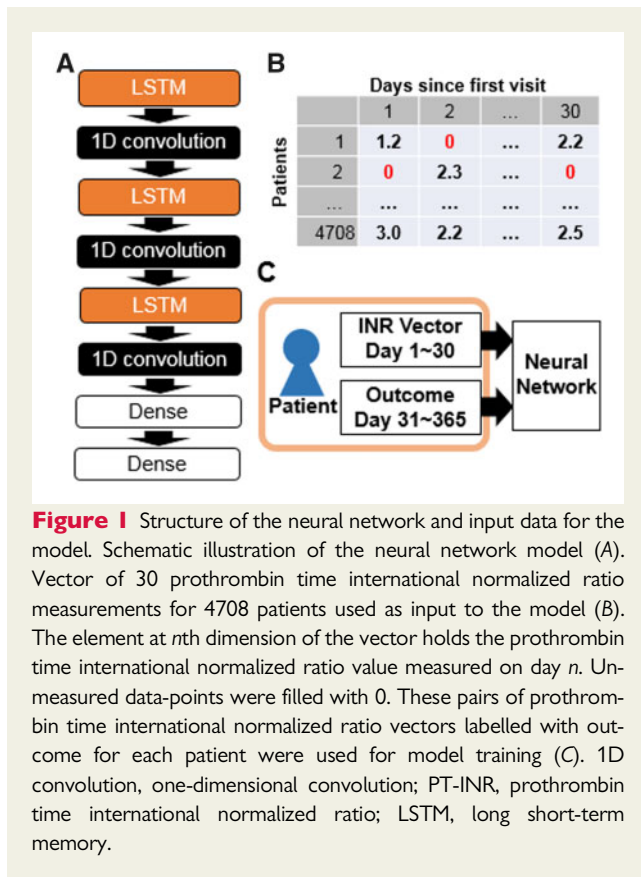
Patients who received anticoagulation therapy with VKA and had three or more PT-INR measurements within the first 30 days after enrolment were included in the model. Patients were excluded if they had experienced any outcome events such as serious bleeding or stroke or died within the first 30 days. In this analysis, day of first visit was set as day 0. Patients from cohorts 1–3 (recruited between March 2010 and October 2014) were included in the derivation cohort whereas those in cohort 4–5 (recruited March 2014 to August 2016) comprised the validation cohort. This study design was considered stringent because each GARFIELD-AF cohort exhibited substantial differences in terms of participating countries, use of anticoagulants, and outcomes.<sup>11</sup>

### Follow-up

Collection of follow-up data occurred at 4 monthly intervals based on medical records and, sometimes, telephone interviews up to 24 months. The incidence of ischaemic stroke, transient ischaemic attack (TIA), SE, acute coronary syndrome, hospitalization, death (cardiovascular and non-cardiovascular), chronic heart failure (CHF; occurrence or worsening), and bleeding (severity and location) was documented. An audit and quality control programme was applied, and data were examined for completeness and accuracy by the co-ordinating centre (TRI, London, UK). By design, 20% of all electronic case report forms in the GARFIELD-AF registry were monitored against source documentation at sites over the 8 years of recruitment and follow-up.

### Outcomes

Outcome measures used in this analysis were major bleeding, stroke/SE, and all-cause death occurring between days 31 and 365. Major bleed was classified by investigators according to International Society on Thrombosis and Haemostasis definition.<sup>30</sup> Stroke/SE was defined as a combined Endpoint of ischaemic stroke, SE, and TIA.



**Figure 1** Structure of the neural network and input data for the model. Schematic illustration of the neural network model (A). Vector of 30 prothrombin time international normalized ratio measurements for 4708 patients used as input to the model (B). The element at  $n$ th dimension of the vector holds the prothrombin time international normalized ratio value measured on day  $n$ . Un-measured data-points were filled with 0. These pairs of prothrombin time international normalized ratio vectors labelled with outcome for each patient were used for model training (C). 1D convolution, one-dimensional convolution; PT-INR, prothrombin time international normalized ratio; LSTM, long short-term memory.

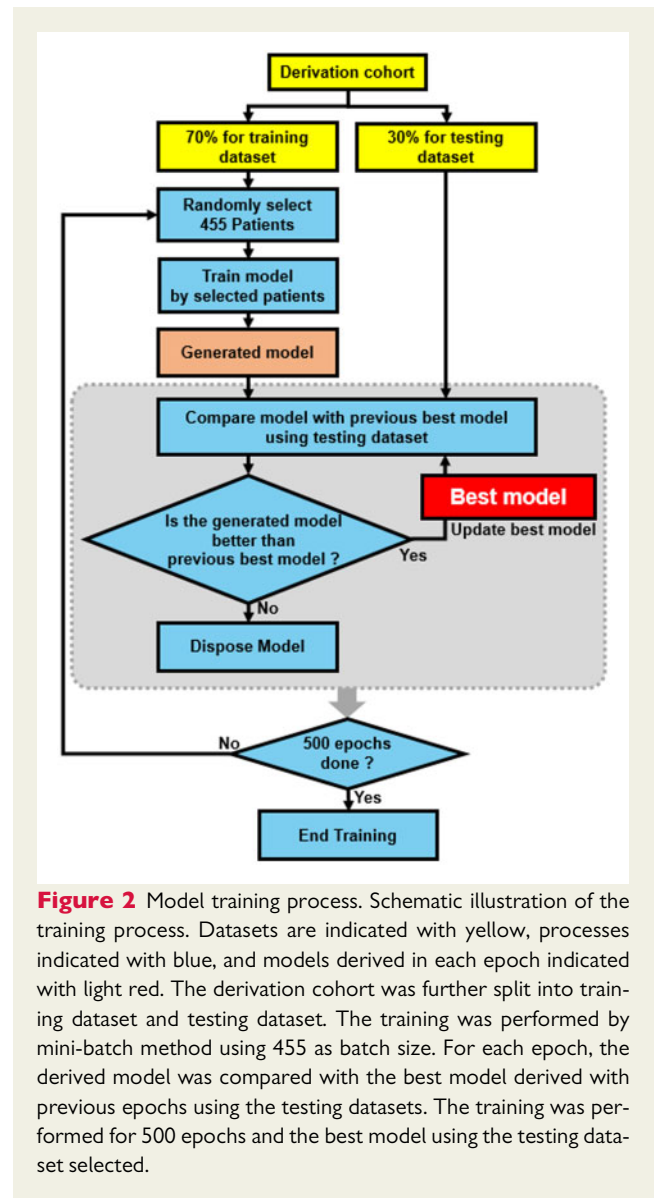
## Artificial intelligence model

The structure of neural networks for the AI model is shown in Figure 1A. To deal with serial data on raw PT-INR measurements, the AI model was constructed by stacking multiple layers of special neurons that can deal with time-dependent data, namely one-dimensional convolution layer and long short-term memory (LSTM) layer. The LSTM layer transfer rectified data to each neighbouring neuron.<sup>31</sup> This structure allows the layer to learn time-dependent data in sequential order.

The neural network model was trained independently for each outcome event. For training, PT-INR measurement patterns for each individual patient were converted to a 30 dimensional PT-INR vector as shown in Figure 1B. All PT-INR measurements obtained within the first 30 days were input to the model. The measured PT-INR value was inserted into  $n$ th element of the 30 dimensional vector, where  $n$  is the number of days after starting VKA. Un-measured data-points were filled with 0. Each vector for patients was labelled with the occurrence of outcome (0 for no event and 1 for event for all three outcome measures) within days 31–365. The neural networks were trained with the multi-dimensional dataset of the PT-INR vector and outcome label as shown in Figure 1C.

## Model training

The process of model training is shown in Figure 2. The training was performed using only patient data from the derivation cohort. The derivation cohort was further split into training (70%) and testing (30%) datasets. The training was performed for 500 epochs and each training epoch included a mini-batch of 455 patients randomly selected from the training dataset. Conceptually, the performance of the model is designed to improve by training with longer epochs. However, this approach can also result in overfitting. To avoid this pitfall and select the model with best



**Figure 2** Model training process. Schematic illustration of the training process. Datasets are indicated with yellow, processes indicated with blue, and models derived in each epoch indicated with light red. The derivation cohort was further split into training dataset and testing dataset. The training was performed by mini-batch method using 455 as batch size. For each epoch, the derived model was compared with the best model derived with previous epochs using the testing datasets. The training was performed for 500 epochs and the best model using the testing dataset selected.

performance, the model was evaluated using the testing dataset at the end of each epoch. The final model was that which performed best with the testing dataset. The performance was measured by calculating the c-statistics of the prediction model for all the data in testing dataset. No data from validation cohort were used for training.

## Model validation

The derived models were validated by inputting the 30 day PT-INR vector and obtaining prediction scores for each outcome. Predicted outcomes were compared with the actual clinical course for each individual patient in the validation cohort. Receiver operating characteristic (ROC) curves were drawn to evaluate the predictive value of the model. The threshold to achieve overall best accuracy for the model was determined and the model's sensitivity and specificity calculated at that threshold. To test the ability of the model to discriminate between high- and low-risk patients for each event, three sets of Kaplan–Meier plots were drawn for event rates stratified as high and low risk with the threshold.





**Table 1** Patients' baseline demographics and clinical characteristics

	≥3 PT-INRs (N = 4708)	0–2 PT-INRs (N = 9630)	≥3 PT-INRs subgroup	
			Derivation (N = 3185)	Validation (N = 1523)
Sex, n (%)				
Female	2085 (44.3)	4330 (45.0)	1420 (44.6)	665 (43.7)
Male	2623 (55.7)	5300 (55.0)	1765 (55.4)	858 (56.3)
Age at dx, years	72.1 (9.9)	70.0 (10.7)	72.2 (9.7)	72.0 (10.2)
BMI, kg/m <sup>2</sup>	28.7 (5.9)	28.1 (5.7)	28.6 (5.7)	29.0 (6.1)
LVEF, %	53.7 (12.9)	55.7 (12.7)	53.2 (13.2)	54.7 (12.2)
Type of AF, n (%)				
New	2409 (51.2)	4087 (42.4)	1706 (53.6)	703 (46.2)
Paroxysmal	798 (16.9)	2207 (22.9)	567 (17.8)	231 (15.2)
Permanent	877 (18.6)	1514 (15.7)	487 (15.3)	390 (25.6)
Persistent	624 (13.3)	1822 (18.9)	425 (13.3)	199 (13.1)
CHF, n (%)	721 (15.3)	2149 (22.3)	466 (14.6)	255 (16.7)
CAD, n (%)	878 (18.6)	1896 (19.7)	511 (16.0)	367 (24.1)
ACS	461 (9.8)	872 (9.1)	292 (9.2)	169 (11.1)
CHA <sub>2</sub> DS <sub>2</sub> -VASc	3.4 (1.5)	3.3 (1.5)	34 (1.5)	33 (1.4)
HAS-BLED	1.4 (0.9)	1.4 (0.9)	15 (0.9)	14 (0.9)

Values are mean (SD) unless specified otherwise.

ACS, acute coronary syndromes; AF, atrial fibrillation; BMI, body mass index; CAD, coronary artery disease; CHF, congestive heart failure; LVEF, left ventricular ejection fraction.

than those with greater values.<sup>36</sup> Within 30 days, TTR has low predictive power because early-phase PT-INR values vary greatly due to a number of influencing factors including genetics,<sup>21,22</sup> choice of commercial thromboplastin and coagulometer device,<sup>37–39</sup> and patients' lifestyles.<sup>40</sup> With the use of AI, we show here the presence of important information in raw PT-INR patterns over first 30 days that can predict clinical events occurring from days 31 to 365.

Multiple useful models exist to predict clinical outcomes in patients with AF.<sup>1–3,8,41,42</sup> However, most use single time-point data. HAS-BLED score, on the other hand, does include time-series data on PT-INR in the guise of labile PT-INR, which is expressed by TTR.<sup>3</sup> Our AI model using time-series PT-INR values has better predictive power than TTR for major clinical events, at least in the early phase of VKA initiation. Even with multi-dimensional data including 31 datasets our AI model output is a prediction score given as a single value. Thus, output of the new model may be included in conventional scoring models by introducing a cut-off, similarly to integration of TTR.<sup>18</sup> The ROC analysis in validation cohort revealed that our AI model has modest predictive power with a best c-statistic 0.78, for major bleeding. However, the model could be usefully incorporated into previous models and thereby improve their accuracy, as has been done with TTR.<sup>18</sup> Our prediction model could expand to automatic prediction of clinical outcomes from multi-dimensional data when incorporated into electrical patient recording systems, for example. Since our models are able to predict clinical outcomes in the early phase of treatment, they may discriminate patients who are unsuitable for VKA therapy and suggest switching them to NOACs, which are associated with lower bleeding risk compared with VKA.

Our novel AI model comprising a neural network can efficiently connect multiple time-dependent measurements to clinical outcomes

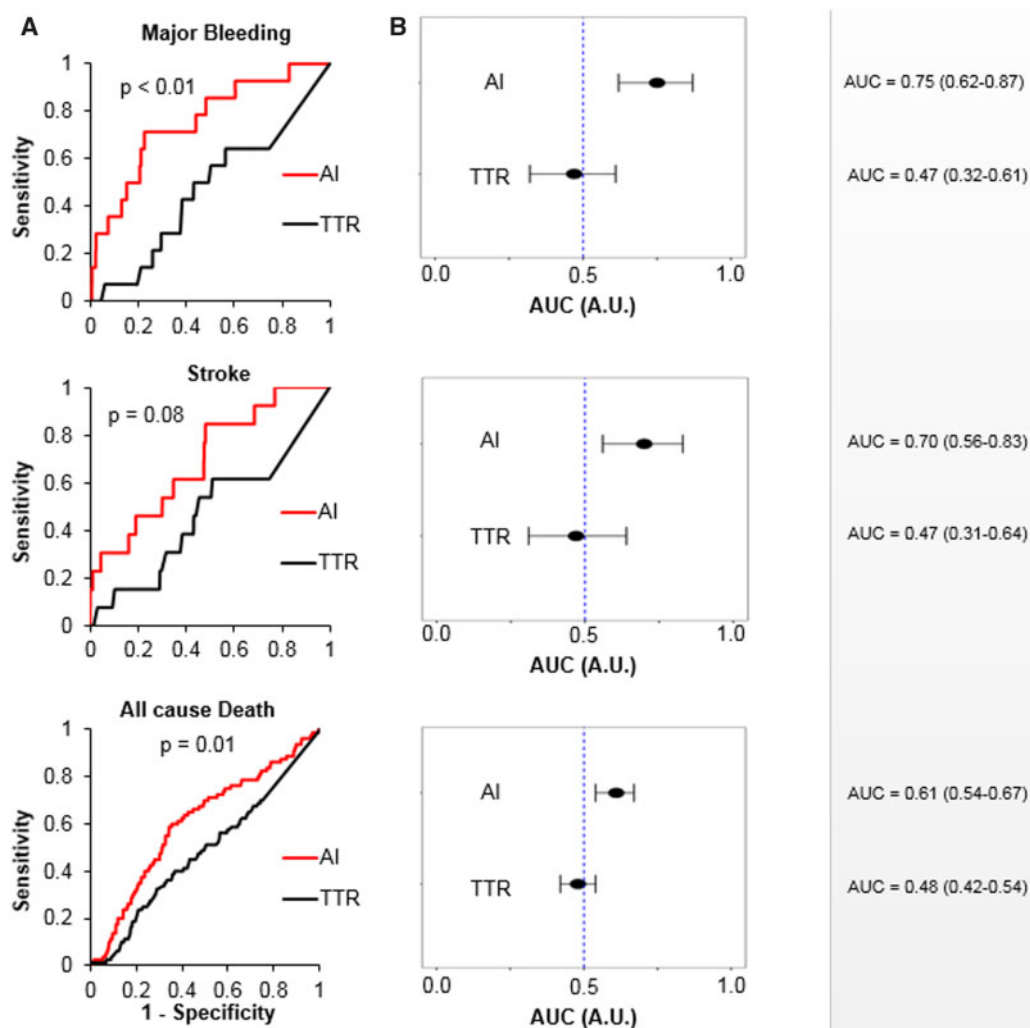
to form a prediction model. Although in this study, the network was used only to learn PT-INR patterns as specific target, the same structure may have the ability to convert other multi-dimensional time-dependent measurements to prediction models. Therefore, the network may provide a new means to incorporate time-dependent data in prediction models.

Output values from our AI model are related to risk of future events but not their probability. Therefore, calibration of the model with typical Hosmer-Lemeshow goodness-of-fit (GOF) test is not feasible.

### Study limitations

Several limitations of this analysis should be noted. First, validation of the AI models was performed using datasets derived from the GARFIELD-AF registry. External validations of the AI model were not conducted. Thus, validity of this model beyond GARFIELD-AF patients is unknown. On the other hand, large dissimilarities between cohorts 1–3 and 4 and 5 were noted, suggesting that our model is sufficiently robust to apply in daily clinical practice. Prothrombin time international normalized ratio within 30 days may be influenced by concomitant dosing with parenteral anticoagulants. However, our model attempted to account for all influencing factors beyond the effects of VKA. We hope that other researchers will test our model's performance in external datasets.

Second, the AI model was trained only with PT-INR data and did not include other information such as sex, age, biomarkers, concomitant drugs, or other serially measured values. Although consecutive patient data were analysed, unrecognized confounders may exist. Many other known risk factors for adverse outcome events were not considered in our models.



**Figure 4** ROC analysis of the artificial intelligence model. Comparison of receiver operating characteristic curves compiled from artificial intelligence model and time in therapeutic range (A). Comparisons were performed using stratified bootstrap method with 2000 bootstrap rounds. Forest plots of area under curve of the receiver operating characteristic curve for each outcome (B). The 95% confidence intervals were calculated by stratified bootstrap method with 2000 bootstrap rounds. AUC, area under curve; CI, confidence interval; ROC, receiver operating characteristic.

**Table 2** Best predictive accuracies and corresponding sensitivities and specificities (95% CIs) for validation cohort

AI	Accuracy	Sensitivity	Specificity
Major bleed	0.78 (0.40–0.92)	0.79 (0.50–1.00)	0.78 (0.39–0.93)
Stroke	0.53 (0.24–0.98)	0.85 (0.31–1.00)	0.53 (0.23–0.99)
All-cause death	0.64 (0.51–0.69)	0.63 (0.50–0.76)	0.65 (0.50–0.70)

AI, artificial intelligence.

Third, by selecting only patients with >3 PT-INR measurements within 30 days, two-thirds of the entire cohort were excluded, which could introduce selection bias. Furthermore, patients do not necessarily remain stable after day 31 and our model cannot capture changes

at time-points later than day 31. Future studies will examine the impact of time periods beyond 30 days in relation to AI risk prediction.

Fourth, although our results suggest the presence of crucial information within the PT-INR measurement pattern to predict patients' clinical course, the nature of that information is unknown. It might be present in the target PT-INR value, PT-INR fluctuations, PT-INR measurement frequency, or elsewhere.

Fifth, the c-statistics, sensitivity, and specificity of our models are far from perfect. Further studies to improve predictive accuracy possibly by adding other clinical characteristics and measurements are necessary.

Sixth, statistical significance was not achieved in either the derivation or validation cohort in comparison with TTR for prediction of all-cause death and stroke. This could be explained by low numbers of events limiting statistical power. Moreover, even though the number of deaths observed was not low, they could have been caused by factors not







37. Poller L, Ibrahim S, Keown M, Pattison A, Jespersen J; European Action on Anticoagulation. The prothrombin time/international normalized ratio (PT-INR) line: derivation of local INR with commercial thromboplastins and coagulometers—two independent studies. *J Thromb Haemost* 2011;**9**:140–148.
38. Christensen TD, Larsen TB. Precision and accuracy of point-of-care testing coagulometers for self-testing and management of oral anticoagulation therapy. *J Thromb Haemost* 2012;**10**:251–260.
39. Hemkens LG, Hilden KM, Hartschen S, Kaiser T, Didjurgeit U, Hansen R, Bender R, Sawicki PT. A randomized trial comparing INR monitoring devices in patients with anticoagulation self-management: evaluation of a novel error-grid approach. *J Thromb Thrombolysis* 2008;**26**:22–30.
40. Custódio das Dôres SM, Booth SL, Martini LA, de Carvalho Gouvêa VH, Padovani CR, de Abreu Maffei FH, Campana AO, Rupp de Paiva SA. Relationship between diet and response to warfarin: a factor analysis. *Eur J Nutr* 2007;**46**:147–154.
41. Fang MC, Go AS, Chang Y, Borowsky LH, Pomernacki NK, Udaltsova N, Singer DE. A new risk scheme to predict warfarin-associated hemorrhage: the ATRIA (Anticoagulation and Risk Factors in Atrial Fibrillation) Study. *J Am Coll Cardiol* 2011;**58**:395–401.
42. Gage BF, Yan Y, Milligan PE, Waterman AD, Culverhouse R, Rich MW, Radford MJ. Clinical classification schemes for predicting hemorrhage: results from the National Registry of Atrial Fibrillation (NRAF). *Am Heart J* 2006;**151**:713–719.