

# A Deep-Learning Algorithm-Enhanced Electrocardiogram Interpretation for Detecting Pulmonary Embolism

Yu-Cheng Chen,<sup>1</sup> Sung-Chiao Tsai,<sup>2</sup> Chin Lin,<sup>3,4,5</sup> Chin-Sheng Lin,<sup>2</sup> Wen-Hui Fang,<sup>6</sup>  
Yu-Sheng Lou,<sup>3,4</sup> Chia-Cheng Lee<sup>7,8</sup> and Pang-Yen Liu<sup>2</sup>

**Background:** The early diagnosis of pulmonary embolism (PE) remains a challenge. Electrocardiograms (ECGs) and D-dimer levels are used to screen potential cases.

**Objective:** To develop a deep learning model (DLM) to detect PE using ECGs and investigate the clinical value of false detections in patients without PE.

**Methods:** Among patients who visited the emergency department between 2011 and 2019, PE cases were identified through a review of medical records. Non-PE ECGs were collected from patients without a diagnostic code for PE. There were 113 PE and 51,456 non-PE ECGs in the training and validation sets for developing the DLM, respectively, and 27 PE and 13,105 non-PE cases in an independent testing set for performance validation. A human-machine competition was conducted from the testing set to compare the performance of the DLM with that of physicians. Receiver operating characteristic (ROC) curves, sensitivity, and specificity were used to determine the diagnostic value. Survival analysis was used to assess the prognosis of the patients without PE, stratified by DLM prediction.

**Results:** The DLM was as effective as physicians in diagnosing PE, with 70.8% sensitivity and 69.7% specificity. The area under the ROC curve of DLM was 0.778 in the testing set and up to 0.9 with D-dimer and demographic data. The non-PE patients whose ECG was misclassified as PE by DLM had higher all-cause mortality [hazard ratio (HR) 2.13 (1.51-3.02)] and risk of non-cardiovascular hospitalization [HR 1.55 (1.42-1.68)] than those correctly classified.

**Conclusions:** A DLM-enhanced ECG system may prompt PE recognition and provide prognostic outcomes in patients with false-positive predictions.

**Key Words:** Deep learning model • Electrocardiogram • Pulmonary embolism

## INTRODUCTION

Pulmonary embolism (PE) is a potentially fatal disease that presents with nonspecific symptoms and signs.<sup>1</sup> Massive PE-compromised hemodynamic stability may require urgent interventions to reverse the progressive worsening of circulatory function. An early diagnosis can be lifesaving, although it is highly challenging. Computed tomographic pulmonary angiography (CTPA) is the method of choice to establish a definite diagnosis, especially in emergency settings, owing to its high accuracy and short acquisition time.<sup>2</sup> However, it often takes hours for patients in the emergency department (ED) to un-

Received: December 11, 2022 Accepted: April 10, 2023

<sup>1</sup>Department of Internal Medicine; <sup>2</sup>Division of Cardiology, Department of Internal Medicine, Tri-Service General Hospital, National Defense Medical Center; <sup>3</sup>Graduate Institute of Life Sciences; <sup>4</sup>School of Public Health; <sup>5</sup>School of Medicine, National Defense Medical Center; <sup>6</sup>Department of Family and Community Medicine, Department of Internal Medicine; <sup>7</sup>Planning and Management Office; <sup>8</sup>Division of Colorectal Surgery, Department of Surgery, Tri-Service General Hospital, National Defense Medical Center, Taipei, Taiwan.

Corresponding author: Dr. Pang-Yen Liu, Division of Cardiology, Department of Internal Medicine, Tri-Service General Hospital, National Defense Medical Center, No. 325, Sec. 2, Chenggong Rd., Neihu, Taipei 114, Taiwan. Tel: 886-2-8792-3311 ext. 13589; Fax: 886-2-6601-2656; E-mail: liupydr@gmail.com

**Abbreviations**

AF	Atrial fibrillation
ALT	Alanine aminotransferase
AMI	Acute myocardial infarction
AST	Aspartate aminotransferase
AUC	Areas under the curve
AV	Atrioventricular
BMI	Body mass index
BNP	Brain natriuretic peptide
BUN	Blood urea nitrogen
CAD	Coronary artery disease
CI	Confidence interval
CK	Creatine kinase
CKD	Chronic kidney disease
Cl	Chloride
COPD	Chronic obstructive pulmonary disease
Cr	Creatinine
CT	Computed tomography
CTPA	Computed tomographic pulmonary angiography
DLM	Deep learning model
DM	Diabetes mellitus
ECG	Electrocardiogram
ED	Emergency department
eGFR	Estimated glomerular filtration rate
FOBT	Fecal occult blood test
GLU	Fasting glucose
Hb	Hemoglobin
HF	Heart failure
HLP	Hyperlipidemia
HR	Hazard ratio
HTN	Hypertension
ICD	International Classification of Diseases
IRB	Institutional Review Board
K	Potassium
Mg	Magnesium
Na	Sodium
OR	Odds ratio
PE	Pulmonary embolism
PLT	Platelet
QRSd	QRS duration
QTc	Corrected QT interval
RBBB	Complete right bundle branch block
ROC	Receiver operating characteristic
SMOTE	Synthetic minority over-sampling technique
TC	Total cholesterol
tCa	Total calcium
TG	Triglyceride
Tnl	Troponin I
WBC	White blood cell count
XGB	eXtreme gradient boosting

dergo computed tomography (CT) scan examinations. A previous retrospective study disclosed that the median time of this interval is approximately 3.5 hours for symptomatic PE.<sup>3</sup> The limitation of CTPA is that it is expensive and is only used in highly suspected cases. Therefore, the development of routine examinations to screen potential PE cases is necessary in clinical practice.

Electrocardiogram (ECG) is a rapidly acquired and widely used diagnostic tool that is routinely used in the ED. Several ECG changes caused by acute PE have been reported since the early 20<sup>th</sup> century.<sup>4</sup> Sinus tachycardia, atrial arrhythmia, inverted T waves in leads V1-V4, S1Q3T3 pattern, incomplete or complete right bundle branch block, and other ECG characteristics are thought to be associated with PE. However, the predictive value of PE detection varies among studies and depends on disease severity and the target population.<sup>5,6</sup> The scoring systems developed by Sreeram et al.<sup>7</sup> and Daniel et al.<sup>5</sup> have low sensitivity, even in detecting PE with severe pulmonary hypertension in the latter study. The diagnosis of PE using ECG remains challenging. In recent decades, neural network research has flourished, and deep learning models (DLMs) have been utilized in various fields. ECG interpretation mainly involves detecting and extracting morphological features, which are the strengths of neural networks. DLM-enhanced ECG interpretation has been applied to detect many disorders.<sup>8-18</sup> An adequately trained DLM should be able to detect specific ECG patterns for PE diagnosis, and may even reveal previously unrecognized features diagnosed by humans.

The current first-line clinical evaluation for suspected PE relies on the clinical presentation and laboratory data, especially D-dimer levels. D-dimer testing has been reported to have a high negative predictive value, and a normal D-dimer level makes the diagnosis of acute PE unlikely.<sup>2</sup> However, a positive D-dimer test result cannot be used to confirm the diagnosis of PE due to its low positive predictive value. Various causes including infection, inflammation, cancer, and pregnancy, can also induce elevation of plasma D-dimer levels and account for a large proportion of people visiting the ED. Considering that other conditions associated with high D-dimer levels may not exhibit ECG changes, ECG may provide additional supportive evidence for identifying potential PE cases with high D-dimer levels. Moreover, DLM-enhanced ECG interpretation has been demonstrated to be able to iden-

tify the predictors of cardiovascular diseases by false-positive predictions,<sup>19</sup> which further provides clinical prognostic value.

In this study, we aimed to develop a DLM-enhanced PE detection system and compare its accuracy to that of physicians, which may help to prompt clinical awareness and allow for the early diagnosis of PE. We also validated the proposed detection system in a subset of patients with high D-dimer levels.

## METHODS

### Population

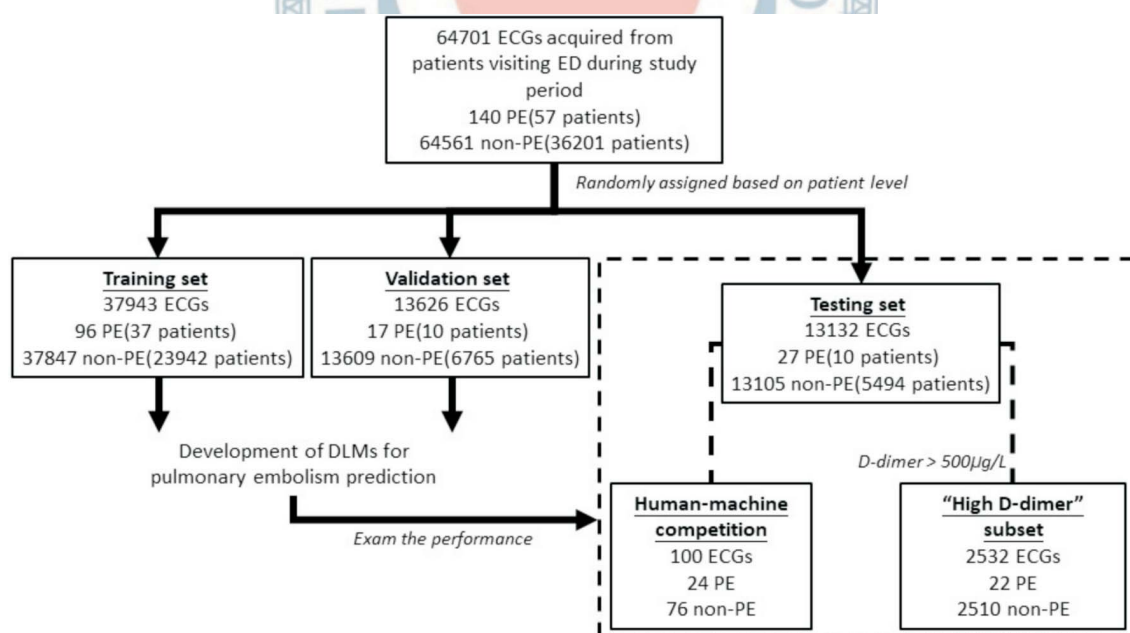
In this single-center retrospective study, all data were collected from the Tri-Service General Hospital, Taipei, Taiwan, and the Institutional Review Board of Tri-Service General Hospital approved the study (IRB NO. C202005055). Patients who visited the ED between December 2011 and December 2019 were included in this study. Cases of PE were identified according to an International Classification of Diseases, Ninth Revision (ICD-9-CM) diagnosis code of 415.x or an ICD-10 code of I26.xx. The diagnoses of the patients were confirmed by

reviewing their medical records. PE was defined on the basis of CTPA, and suspected PE without CT was excluded. Finally, 140 PE ECGs from 57 patients were included in this study. Patients without PE were recruited from those who visited the ED during the same period without the diagnoses codes mentioned above, and 64,561 non-PE ECGs from 36,201 patients were included in this study.

Figure 1 shows the process of generating the training, validation, and testing sets. The ECGs were assigned at random, based on the patient level; therefore, no ECG overlap was present in these datasets. The training set included 96 PE ECGs from 37 patients and 37,847 non-PE ECGs from 23,942 patients. The validation set included 17 PE ECGs from 10 patients and 13,609 non-PE ECGs from 6,765 patients. These were used to develop a DLM for PE detection. An independent testing set including 27 PE ECGs from 10 patients and 13,105 non-PE ECGs from 5,494 patients was used to validate the DLM performance. Within the testing set, the patients with D-dimer values higher than 500 µg/L during their ED visits were included in a “high D-dimer” subset.

### Data source

ECGs were recorded using a Philips 12-lead ECG ma-



**Figure 1.** Establishment of the training, validation, and testing sets. This figure shows the schematic presentation of the dataset creation and analysis strategy, which was devised to assure a robust and reliable dataset for training, validating, and testing of the network. Once an ECG was placed in one of the datasets, that ECG was used only in that set, avoiding ‘cross-contamination’ among the training, validation, and testing datasets. The details of the flow chart and how each of the datasets was used are described in the Methods. DLM, deep learning model; ECG, electrocardiogram; ED, emergency department; PE, pulmonary embolism.

chine (PH080A) with a 500-Hz sampling frequency and 10 s in each lead. Quantitative measurements and findings from the final ECG clinical reports were extracted to identify 31 diagnostic pattern classes and eight continuous ECG measurements. The eight ECG measurements included heart rate, PR interval, QRS duration, QT interval, correct QT interval, P wave axis, RS wave axis, and T wave axis. Data for these variables were 90-100% complete, and missing values were imputed using multiple imputations with the fully conditional specification method using the Multivariate Imputation via Chained Equations algorithm, as described by Van Buuren and Groothuis-Oudshoorn.<sup>20</sup> Each variable had its own imputation model. Built-in imputation models are provided for continuous data (predictive mean matching, normal), binary data (logistic regression), unordered categorical data (polynomial logistic regression), and ordered categorical data (proportional odds). Patterns included abnormal T wave (indicating reduced T wave amplitude, both absolute and relative to the QRS, and negative T waves), atrial fibrillation, atrial flutter, atrial premature complex, complete atrioventricular (AV) block, complete left bundle branch block, complete right bundle branch block (RBBB), first degree AV block, incomplete left bundle branch block, incomplete RBBB, ischemia/infarction, junctional rhythm, left anterior fascicular block, left atrial enlargement, left axis deviation, left posterior fascicular block, left ventricular hypertrophy, low QRS voltage, pacemaker rhythm, prolonged QT interval (QTc > 485 ms by either the Bazett formula or Fridericia formula), right atrial enlargement, right ventricular hypertrophy, second degree AV block, sinus bradycardia, sinus pause, sinus rhythm, sinus tachycardia, supraventricular tachycardia (extreme tachycardia, > 220-age, with narrow QRS complex, QRSd < 120 ms, which did not fulfill the criteria of other supraventricular rhythms, such as sinus tachycardia, atrial fibrillation, atrial flutter, and junctional tachycardia), ventricular premature complex, ventricular tachycardia, and Wolff–Parkinson–White syndrome. The 31 clinical diagnostic patterns were parsed from the structured findings statements based on key phrases that are standard within the Philips system. Detailed definitions of the diagnostic patterns mentioned above are described in the Philips DXL ECG Algorithm Physician's Guide. The corresponding electronic medical records associated with each ECG in our hospital, includ-

ing medical records, nursing records, procedure records, laboratory data, and imaging examinations, were also collected for subsequent covariate and outcome extraction.

### Implementation of the machine learning model

We previously developed an 82-layer convolutional layer and attention mechanism architecture called ECG12Net. Technology details including model architecture, data augmentation, and model visualization have been described previously.<sup>17</sup> Based on the same architecture, we trained a new DLM to estimate the ECG-based possibility of PE. Each original ECG signal length was considered as a  $12 \times 5000$  matrix. We randomly cropped 1,024 sequences as inputs in the training process. To perform the random cropping process, a starting point was randomly selected between the 1<sup>st</sup> and 3977<sup>th</sup> data points of the ECG raw data using R's built-in random-sampling function. Subsequently, a segment with a length of 1,023 data points, followed by the starting point, was cropped as a training sample. For the inference stage, nine overlapping lengths of 1,024 sequences based on interval sampling were used to generate predictions that were averaged as the final prediction as previously described.<sup>11,13</sup>

An oversampling process of directly duplicating the minor sample was used in the training step because of the low rate of PE in our dataset. An ECG signal is a continuous sequence and cannot be directly defined as an independent feature vector. Thus, it is not suitable to generate new ECG signal samples using an interpolating algorithm such as the Synthetic Minority Over-sampling Technique (SMOTE) method.

The settings for the training model were as follows: (1) Adam optimizer with standard parameters ( $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ ) and a batch size of 36 for optimization, (2) learning rate of 0.001, and (3) a weight decay of  $10^{-4}$ . The 100th epoch model was used as the final model, which presented performance in the testing set that was evaluated only once.

To compare the usage of ECG voltage-time traces and the corresponding clinically reported ECG measures, we trained an eXtreme gradient boosting (XGB) model and elastic net using 31 diagnostic pattern classes and eight ECG measurements to recognize PE in the training set. Moreover, the XGB model and elastic net also pro-

vided corresponding variable importance rankings to explore the relationships between the explainable features and PE.

### Study covariates and outcomes

The study covariates, including demographics, disease histories, and laboratory test results, were obtained from the electronic medical records mentioned previously. Patient demographics such as sex, age, body height, body weight, body mass index (BMI), systolic blood pressure, diastolic blood pressure, and presentation of chest pain were extracted from medical and nursing records at the ED visit (Supplementary Table 1). We used ICD-9 and ICD-10 codes to define diabetes mellitus, hypertension, hyperlipidemia, chronic kidney disease, coronary artery disease, stroke, heart failure (HF), atrial fibrillation, and chronic obstructive pulmonary disease (COPD). We collected laboratory values from tests conducted in the ED, including D-dimer, white blood cell count, hemoglobin, platelet, sodium, potassium, chloride, total calcium, magnesium, aspartate aminotransferase, alanine aminotransferase, glucose, creatine kinase, creatinine, blood urea nitrogen, troponin I, NT-pro-B-type natriuretic peptide, triglycerides, and total cholesterol. The 30-day outcomes of interest included all-cause mortality, cardiovascular disease-related mortality, cardiovascular-related hospitalization, and non-cardiovascular-related hospitalization. Mortality was defined based on the electronic medical records of our hospital. Data for live visits were censored at the patient's last known hospital live encounter to limit bias from incomplete records. All visits were divided into admitted or non-admitted groups according to hospitalization outcomes.

### Human-machine competition

To evaluate the performance of our DLM, we conducted a human-machine competition using a testing cohort. The database contained 100 ECGs, including 24 PE cases and 76 non-PE cases. Five doctors participated in the competition (two internal medicine residents, two emergency medicine residents, and one cardiologist), all of whom completed the tests through an online standardized data entry program without patient information except the ECGs. We calculated the sensitivity, specificity, and Youden's index of the doctors' results for comparison with those of the DLM.

### Statistical analysis and model performance assessment

The analyses of patient characteristics and outcomes were based on patients, and evaluation of the model performance was based on ECGs. We analyzed the characteristics and laboratory results of the patients with and without PE in each dataset. The results are presented as means and standard deviations for continuous variables and as numbers and percentages for categorical variables. We used the Student's t-test or chi-square test to compare the results between two groups, as appropriate, and p values < 0.05 were considered statistically significant. Statistical analysis was performed using R version 3.4.4, and the package MXNet version 1.3.0 was used to implement our DLM.

In the primary analysis, we compared the performance of our DLM to that of human experts and two traditional machine learning-based algorithms: the XGB model and elastic net. In addition, an integrated XGB model that combined demographic data (age, BMI, and sex), DLM prediction, and D-dimer values was established through 5-fold cross-validation. Receiver operating characteristic (ROC) curves, areas under the curve (AUCs), and the McNemar test were applied to evaluate the performance in PE recognition of the DLMs, integrated model, and machine-learning algorithms. The operating point was selected on the basis of the maximum Youden's index derived from the validation set. To identify the relationships between clinical characteristics and PE, and the characteristics that led to misdiagnosis by the DLMs, logistic regression was used to calculate the odds ratio (OR) of each clinical characteristic.

In the secondary analysis, the non-PE ECGs in the testing set were separated into "false-positive" cases which were identified as PE ECGs by DLM, and "true-negative" cases which were correctly identified as non-PE ECGs. Kaplan-Meier survival analysis was performed at the patient level with the available follow-up data stratified by the DLM prediction for each outcome of interest. Data were censored based on recent encounters. Hazard ratios (HRs) were calculated using a Cox proportional hazard model, and values with 95% confidence intervals (95% CIs) were reported for all data.

## RESULTS

The corresponding patient characteristics and laboratory data for each ECG set are presented in Table 1. In

**Table 1.** Corresponding patient demographics of the training, validation and testing sets

	Training set			Validation set			Testing set		
	PE (n = 37)	Non-PE (n = 23942)	p value	PE (n = 10)	Non-PE (n = 6765)	p value	PE (n = 10)	Non-PE (n = 5494)	p value*
<b>Demographic data</b>									
Sex (male)	17 (47.2%)	12058 (50.4%)	0.707	4 (44.4%)	3401 (50.3%)	0.752 <sup>#</sup>	7 (50.0%)	2929 (53.4%)	0.802 <0.001
Age (years)	65.8 ± 17.1	58.9 ± 19.9	0.038	66.1 ± 16.9	61.2 ± 19.6	0.426 <sup>#</sup>	60.4 ± 17.7	60.9 ± 19.5	0.984 <sup>#</sup> <0.001
BMI (kg/m <sup>2</sup> )	24.7 ± 3.8	24.2 ± 5.1	0.598	24.9 ± 4.8	24.2 ± 5.1	0.457 <sup>#</sup>	25.8 ± 4.2	24.2 ± 5.6	0.309 <sup>#</sup> 0.825
<b>Disease histories</b>									
AMI	0 (0.0%)	600 (2.5%)	1.000 <sup>#</sup>	0 (0.0%)	193 (2.9%)	1.000 <sup>#</sup>	1 (7.1%)	183 (3.3%)	0.379 <sup>#</sup> 0.002
Stroke	3 (8.3%)	3335 (13.9%)	0.333	2 (22.2%)	1182 (17.5%)	0.662 <sup>#</sup>	2 (14.3%)	1037 (18.9%)	1.000 <sup>#</sup> <0.001
CAD	8 (22.2%)	4156 (17.4%)	0.441	0 (0.0%)	1400 (20.7%)	0.216 <sup>#</sup>	6 (42.9%)	1234 (22.5%)	0.101 <sup>#</sup> <0.001
HF	5 (13.9%)	1367 (5.7%)	0.053 <sup>#</sup>	0 (0.0%)	596 (8.8%)	1.000 <sup>#</sup>	4 (28.6%)	534 (9.7%)	0.041 <sup>#</sup> <0.001
AF	4 (11.1%)	963 (4.0%)	0.056 <sup>#</sup>	0 (0.0%)	401 (5.9%)	1.000 <sup>#</sup>	1 (7.1%)	316 (5.8%)	0.565 <sup>#</sup> <0.001
DM	6 (16.7%)	4508 (18.8%)	0.740	2 (22.2%)	1547 (22.9%)	1.000 <sup>#</sup>	4 (28.6%)	1432 (26.1%)	0.767 <sup>#</sup> <0.001
HTN	11 (30.6%)	7315 (30.6%)	0.999	2 (22.2%)	2539 (37.5%)	0.498 <sup>#</sup>	7 (50.0%)	2157 (39.3%)	0.413 <0.001
CKD	1 (2.8%)	1789 (7.5%)	0.519 <sup>#</sup>	0 (0.0%)	750 (11.1%)	0.610 <sup>#</sup>	1 (7.1%)	688 (12.5%)	1.000 <sup>#</sup> <0.001
HLP	7 (19.4%)	5872 (24.5%)	0.479	3 (33.3%)	1948 (28.8%)	0.723 <sup>#</sup>	3 (21.4%)	1592 (29.0%)	0.769 <sup>#</sup> <0.001
COPD	7 (19.4%)	3266 (13.6%)	0.327 <sup>#</sup>	1 (11.1%)	1190 (17.6%)	1.000 <sup>#</sup>	4 (28.6%)	942 (17.2%)	0.280 <sup>#</sup> <0.001
<b>Laboratory data</b>									
D-dimer (µg/L)	9258.0 ± 8947.1	2195.5 ± 5222.2	<0.001	21000.0 ± 13659.4	2166.0 ± 4736.9	<0.001 <sup>#</sup>	8830.1 ± 6711.7	2712.5 ± 5679.1	<0.001 <sup>#</sup> 0.002
eGFR (mL/min/1.73m <sup>2</sup> )	84.5 ± 29.1	82.8 ± 38.6	0.790	62.2 ± 28.5	79.4 ± 38.8	0.125 <sup>#</sup>	61.1 ± 32.6	79.0 ± 40.8	0.015 <sup>#</sup> <0.001
Cr (mg/dL)	0.9 ± 0.3	1.3 ± 1.6	0.179	1.3 ± 0.8	1.4 ± 1.9	0.304 <sup>#</sup>	1.8 ± 2.2	1.5 ± 1.9	0.032 <sup>#</sup> <0.001
BUN (mg/dL)	20.2 ± 18.0	21.9 ± 19.0	0.607	24.6 ± 11.5	23.5 ± 20.7	0.169 <sup>#</sup>	17.2 ± 8.9	23.6 ± 20.8	0.341 <sup>#</sup> <0.001
Na (mmol/L)	137.8 ± 4.0	136.9 ± 4.6	0.263	137.6 ± 4.4	136.7 ± 4.9	0.483 <sup>#</sup>	135.6 ± 3.5	136.8 ± 5.1	0.143 <sup>#</sup> 0.018
K (mmol/L)	4.0 ± 0.5	3.9 ± 0.5	0.174	3.6 ± 0.4	3.9 ± 0.6	0.161 <sup>#</sup>	4.0 ± 0.4	3.9 ± 0.6	0.575 <sup>#</sup> <0.001
Cl (mmol/L)	105.1 ± 5.2	102.7 ± 5.5	0.024	103.9 ± 7.1	102.5 ± 5.5	0.429 <sup>#</sup>	107.2 ± 7.0	102.6 ± 5.8	0.031 <sup>#</sup> 0.024
tCa (mg/dL)	8.6 ± 0.4	8.5 ± 0.7	0.360 <sup>#</sup>	8.5 ± 0.8	8.5 ± 0.7	0.654 <sup>#</sup>	8.2 ± 0.6	8.5 ± 0.7	0.183 <sup>#</sup> 0.912
Mg (mg/dL)	2.1 ± 0.3	2.1 ± 0.4	0.955 <sup>#</sup>	2.1 ± NA	2.1 ± 0.3	0.896 <sup>#</sup>	2.1 ± 0.3	2.1 ± 0.4	0.837 <sup>#</sup> 0.479
TnI (pg/mL)	408.7 ± 1198.9	242.5 ± 3317.6	0.795	193.2 ± 198.5	139.7 ± 2103.3	<0.001 <sup>#</sup>	217.3 ± 380.6	267.5 ± 3127.5	<0.001 <sup>#</sup> 0.066
CK (U/L)	185.7 ± 367.3	195.3 ± 736.6	0.946	261.0 ± 492.6	178.9 ± 721.9	0.974 <sup>#</sup>	96.9 ± 102.8	182.7 ± 815.7	0.221 <sup>#</sup> 0.280
BNP (ng/mL)	223.1 ± 148.7	387.6 ± 762.0	0.104 <sup>#</sup>	387.8 ± 271.4	443.8 ± 868.7	0.269 <sup>#</sup>	849.1 ± 1195.3	515.8 ± 955.7	0.141 <sup>#</sup> <0.001
GLU (g/dL)	157.4 ± 70.1	144.8 ± 81.4	0.441	157.5 ± 58.1	143.2 ± 76.6	0.304 <sup>#</sup>	145.7 ± 72.5	144.9 ± 82.0	0.781 <sup>#</sup> 0.438
Hb (g/dL)	12.7 ± 1.9	12.9 ± 2.3	0.456	12.5 ± 1.5	12.8 ± 2.4	0.508 <sup>#</sup>	13.3 ± 2.0	12.6 ± 2.4	0.249 <sup>#</sup> <0.001
WBC (10 <sup>3</sup> /µL)	9.6 ± 3.5	9.5 ± 6.0	0.897	10.7 ± 4.6	9.4 ± 4.7	0.164 <sup>#</sup>	9.5 ± 2.8	9.4 ± 8.7	0.365 <sup>#</sup> 0.473
PLT (10 <sup>3</sup> /µL)	241.7 ± 97.3	241.9 ± 88.4	0.992	219.7 ± 60.0	241.4 ± 88.8	0.560 <sup>#</sup>	207.3 ± 91.1	239.1 ± 91.2	0.073 <sup>#</sup> 0.101
AST (U/L)	43.1 ± 44.6	45.5 ± 145.1	0.923	59.6 ± 40.2	36.5 ± 86.2	<0.001 <sup>#</sup>	23.3 ± 9.8	39.3 ± 99.4	0.755 <sup>#</sup> <0.001
ALT (U/L)	33.4 ± 46.7	34.5 ± 110.9	0.958	50.5 ± 51.7	31.0 ± 65.7	0.028 <sup>#</sup>	14.6 ± 6.9	33.3 ± 113.1	0.275 <sup>#</sup> 0.083
TG (g/L)	110.9 ± 48.9	128.3 ± 160.1	0.834 <sup>#</sup>	117.8 ± 49.0	125.5 ± 142.5	0.576 <sup>#</sup>	119.8 ± 45.3	124.1 ± 131.8	0.426 <sup>#</sup> 0.698
TC (g/L)	158.8 ± 24.6	155.2 ± 50.5	0.724	174.0 ± 51.0	152.7 ± 49.9	0.272 <sup>#</sup>	164.0 ± 44.0	151.1 ± 46.1	0.180 <sup>#</sup> 0.001

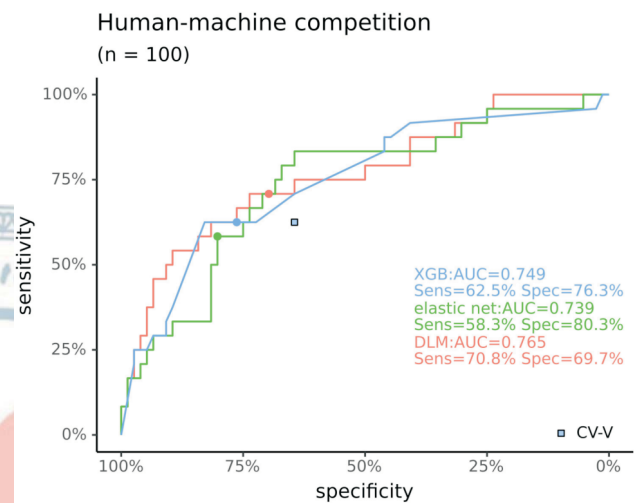
\* Hypothesis test of the difference between the training, validation, and testing sets. <sup>#</sup> p value calculated with n < 25 in continuous variables or by Fisher's exact test for categorical variables.

AF, atrial fibrillation; ALT, alanine aminotransferase; AMI, acute myocardial infarction; AST, aspartate aminotransferase; BMI, body mass index; BNP, brain natriuretic peptide; BUN, blood urea nitrogen; CAD, coronary artery disease; CK, creatine kinase; CKD, chronic kidney disease; COPD, chronic obstructive pulmonary disease; Cr, creatinine; DM, diabetes mellitus; eGFR, estimated glomerular filtration rate; GLU, fasting glucose; Hb, hemoglobin; HF, heart failure; HLP, hyperlipidemia; HTN, hypertension; K, potassium; Mg, Magnesium; Na, sodium; PLT, platelet; TC, total cholesterol; tCa, total calcium; TG, triglyceride; TnI, troponin I; WBC, white blood cell count.

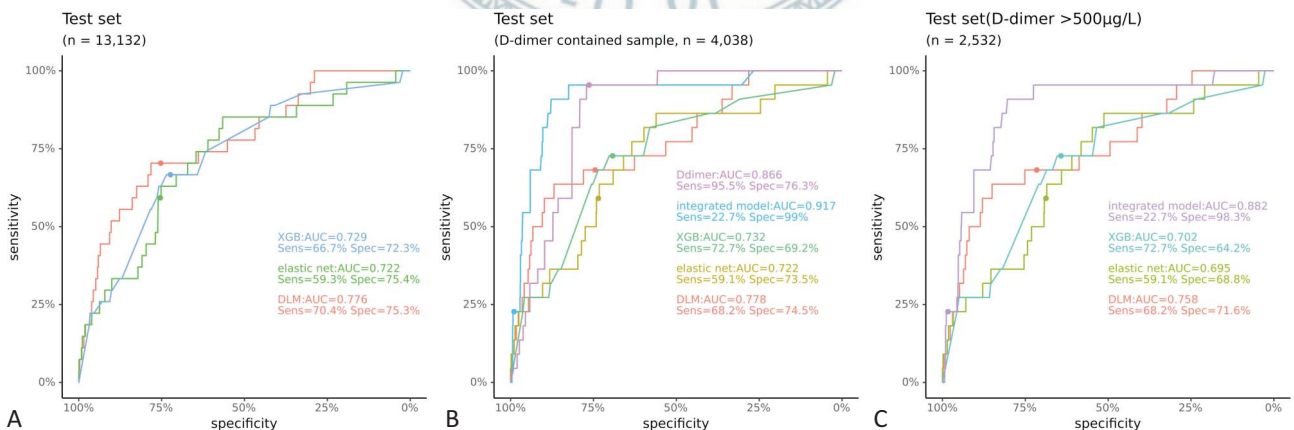
the training set, the patients with PE were older than those without PE (65.8 vs. 58.9 years,  $p = 0.038$ ). There were no statistically significant differences in sex and BMI between the three sets. Among all three sets, D-dimer values were significantly higher in the patients with PE (9258.0 vs. 2195.5  $\mu\text{g}/\text{mL}$ ,  $p < 0.001$  in the training set, 21000.0 vs. 2166.0  $\mu\text{g}/\text{mL}$ ,  $p < 0.001$  in the validation set, and 8830.1 vs. 2712.5  $\mu\text{g}/\text{mL}$ ,  $p < 0.001$  in the testing set).

The PE recognition performances of the human physicians and each model are presented in Figure 2. In the human-machine competition, the AUCs of the DLM, XGB model, and elastic net were 0.765, 0.749, and 0.739, respectively. When choosing the cutoff value with the maximum Youden's index (0.5704) from the validation set, the DLM yielded 70.8% sensitivity and 69.7% specificity. The overall performance of all participating human physicians had a sensitivity ranging from 45.8% to 79.2% and specificity ranging from 48.7% to 80.3%. The attending cardiologist had a Youden's index of 0.27 with a 62.5% sensitivity and approximately 64.4% specificity. There was no significant difference between the human and DLM performance in the competition. The detailed results of the human-machine competition are provided in Supplementary Figure 1. Among the whole testing test, the DLM, XGB model, and elastic net had AUCs of 0.776, 0.729, and 0.722, respectively (Figure 3). When the analysis focused on ECGs with corresponding D-dimer values in the testing set, the AUCs of D-dimer, integrated model, DLM, XGB model, and elastic net were

0.866, 0.917, 0.778, 0.732, and 0.722, respectively. The performance of the integrated model was statistically better than that of D-dimer alone ( $p = 0.01$ ). In the population with a D-dimer value  $> 500 \mu\text{g}/\text{L}$ , the AUC values were 0.882, 0.758, 0.702, and 0.695, for the integrated model, DLM, XGB model, and elastic net, respectively. The performance of each lead in the testing set is presented in Supplementary Figure 2, and their AUCs ranged from 0.5558 to 0.7154, which was lower than that of the combined results of all 12 leads.



**Figure 2.** The performance of DLM, human experts, and machine-learning algorithms in identifying PE ECG in a human-machine competition. The operating point was selected based on the maximum Youden's index obtained from the validation set. The sensitivity and specificity were calculated using the testing set. CV-V: attending cardiologist. AUC, area under the curve; DLM, deep learning model; ECG, electrocardiogram; PE, pulmonary stenosis; XGB, eXtreme gradient boosting.



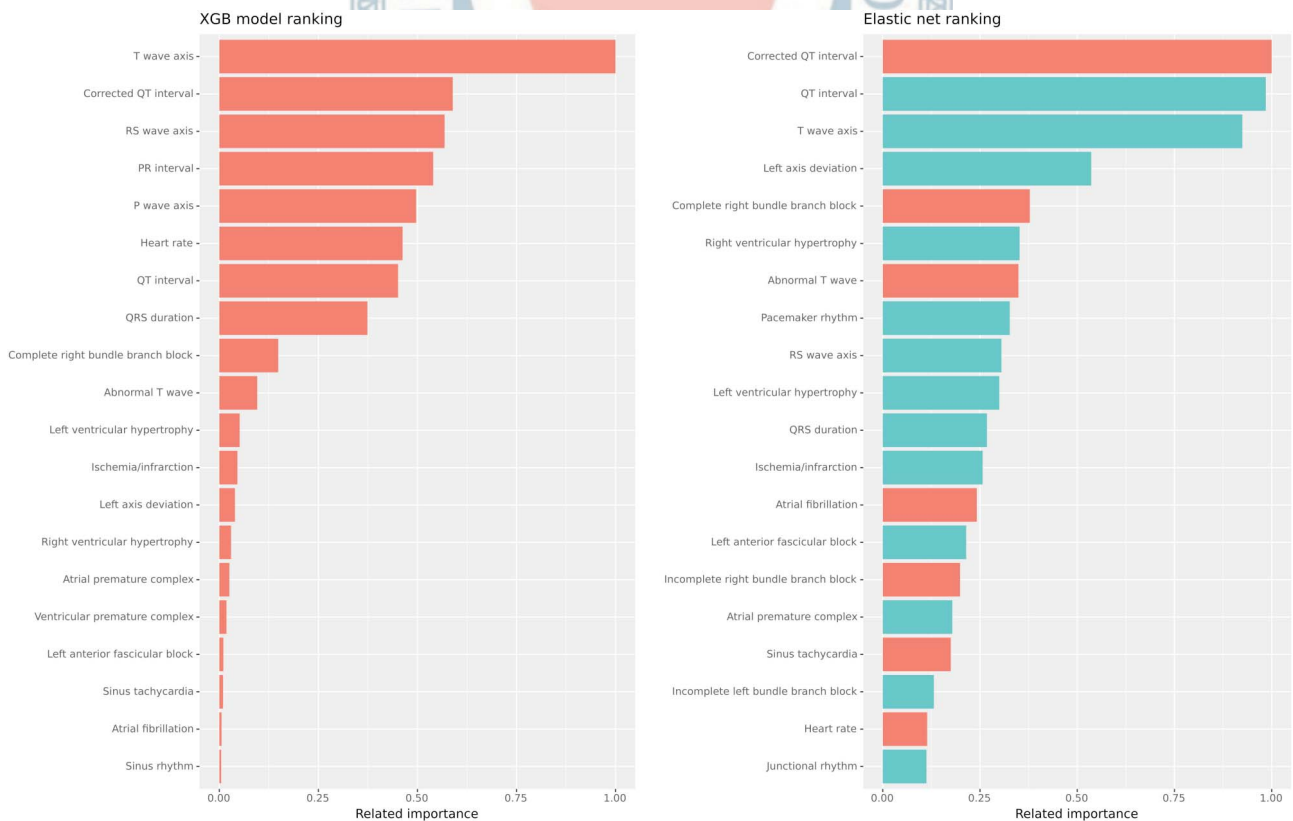
**Figure 3.** The performance of DLM, machine-learning algorithms, and integrated model in identifying PE ECG in the testing set and "High D-dimer" subset. The operating point was selected based on the maximum Youden's index obtained from the validation set. The sensitivity and specificity were calculated using the testing set. AUC, area under the curve; DLM, deep learning model; ECG, electrocardiogram; PE, pulmonary stenosis; XGB, eXtreme gradient boosting.

A stratified analysis of the clinical picture of the patients with PE in the training set is presented as the forest plot in Supplementary Figure 3. Most of the differences between the patients with and without PE were non-significant, except for age, HF, serum creatinine, sodium, and chloride levels. When focusing on misclassified non-PE ECGs in the testing set, the results of the stratified analysis (Supplementary Figure 4) indicated that ECGs of older patients (OR 1.2, 95% CI 1.13-1.28) and those with HF (OR 1.39, 95% CI 1.14-1.68), atrial fibrillation (OR 1.69, 95% CI 1.33-2.14), and COPD (OR 1.20, 95% CI 1.02-1.40) had a higher chance of being misclassified as PE. Regarding laboratory values, this group of patients had lower serum electrolyte levels, lower hemoglobin values, and higher WBC counts.

To further understand the key features in PE recognition, we calculated the importance of various ECG features in traditional machine-learning models, as shown in Figure 4. The highest ranked features in the XGB model were “T wave axis,” “corrected QT interval,” and “RS

wave axis.” In the elastic net, “corrected QT interval,” “T wave axis,” and “QT interval” had the highest relative importance.

We reviewed the ECG presentations in the human-machine competition (Figure 5). Figure 5A shows that sinus tachycardia and the S1Q3T3 pattern were correctly recognized as PE by all physicians and the DLM. ECG has the morphologies of a prolonged correct QT interval (QTc) and negative T wave axis. Interestingly, the saliency map revealed that the DLM mainly focused on the QT and PR segments (Figure 5B). The PE ECG in Figure 5C presents a normal sinus rhythm, prolonged QTc, and left axis deviation of the QRS complex (-25°) and T wave (-82°), which were misdiagnosed by physicians but correctly recognized by the DLM with a focus on the QT segment in the saliency map (Figure 5D). The ECGs in Figure 5E and F were obtained from a patient without PE. They show atrial fibrillation with RBBB morphology and a premature ventricular complex. The QRS complex had an extreme axis deviation (-93°). No specific seg-



**Figure 4.** Related feature importance ranking in the XGB model (information gain) and elastic net (standard coefficient). There are only the top 20 important variables in each model. The red color demonstrates the positive relationship between variables and PE, and the blue color, in contrast, demonstrates the negative relationship. XGB, eXtreme gradient boosting.



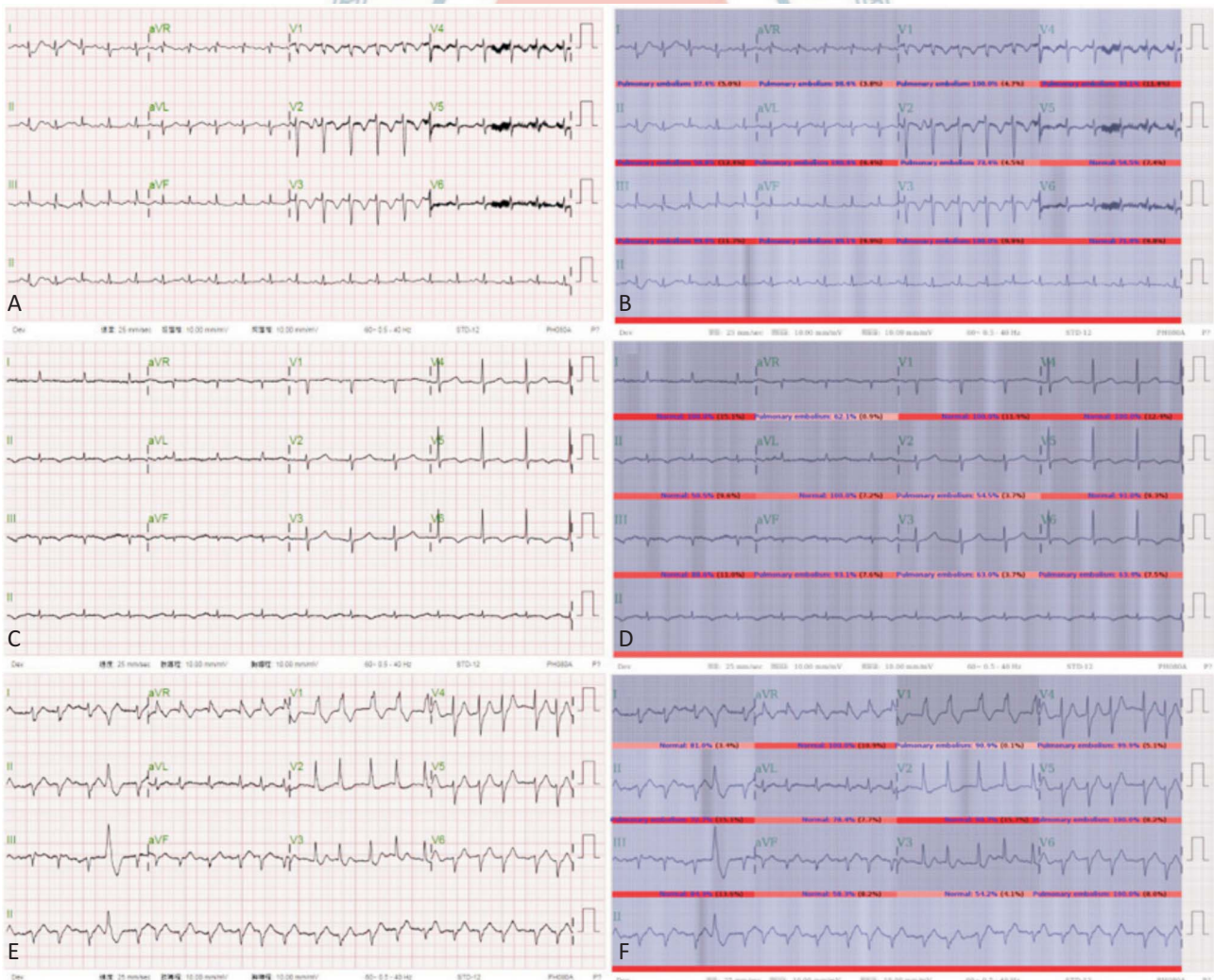
ment or pattern focused on by the DLM could be identified in Figure 5F. This ECG was considered non-PE by all physicians but misclassified as PE by the DLM.

We then performed a 30-day outcome analysis (Figure 6), including mortality and hospitalization events. Cases of non-PE that were misidentified as PE, meaning “false-positive” cases, were compared with “true-negative” cases in which non-PE cases were correctly classified. CVD mortality was not significantly different between the two groups (HR 1.31, 95% CI 0.34-5.06,  $p = 0.696$ ). However, the “false-positive” group had significantly higher all-cause mortality (HR 2.13, 95% CI 1.51-3.02,  $p < 0.0001$ ) than the “true-negative” group. There was no difference in CV hospitalization (HR 1.07, 95% CI

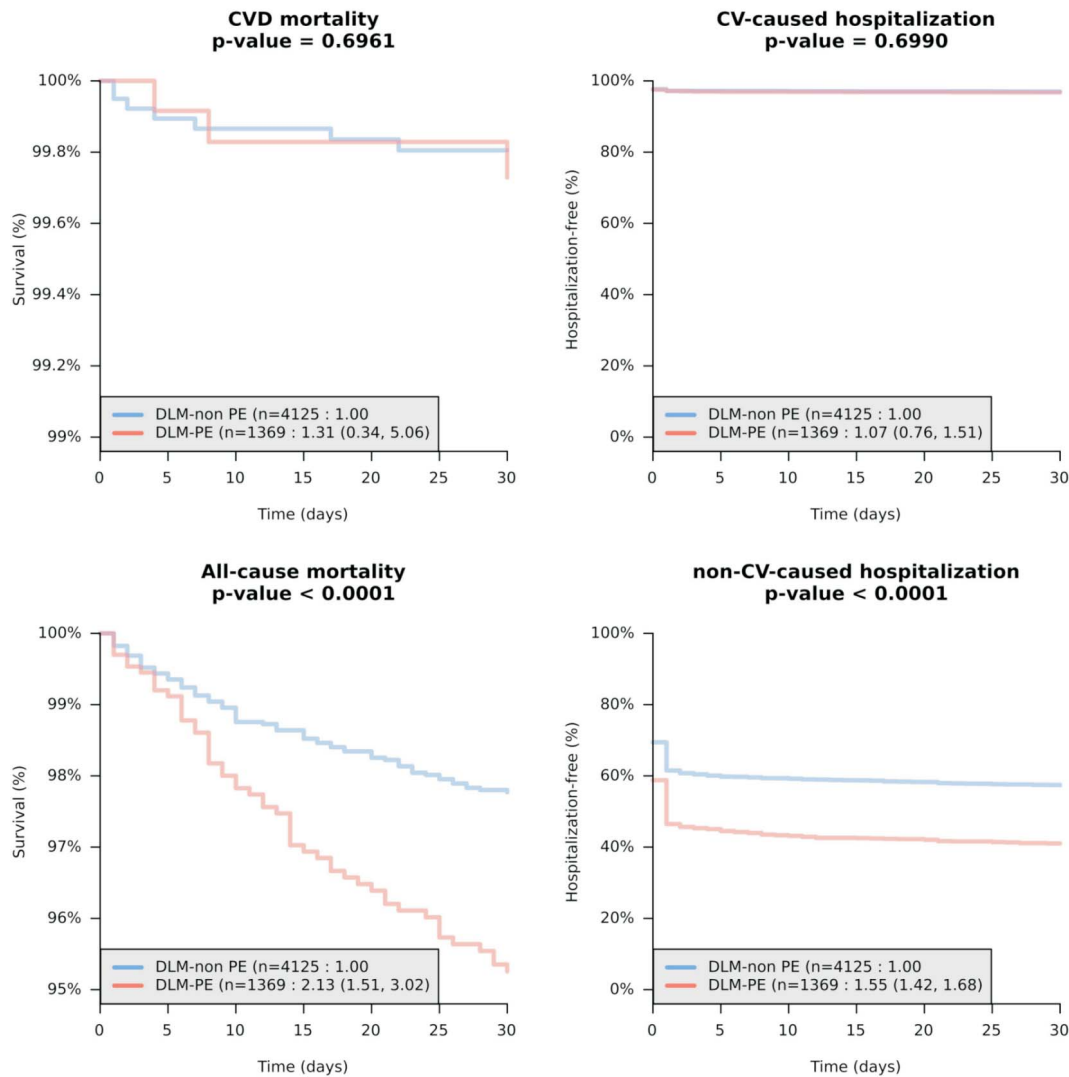
0.76-1.51), but there was a significantly higher non-CV hospitalization rate in the “false-positive” group (HR 1.55, 95% CI 1.42-1.68).

**DISCUSSION**

In this study, we developed a DLM with compatible performance in distinguishing PE ECGs from first-line physicians and traditional machine learning algorithms. In addition, we found that an integrated model including DLM, demographic data, and D-dimer values could provide better performance than each model alone. In further analysis of the clinical characteristics of the mis-



**Figure 5.** Three ECGs, which were identified as PE by DLM, and their saliency map. (A, B) From patient with PE and was recognized as PE by all physicians and DLM. (C, D) From patient with PE and was thought to be non-PE by three of five physicians but as PE by DLM. (E, F) From patients without PE classified as non-PE by all five physicians but as PE by DLM. DLM, deep learning model; ECG, electrocardiogram; PE, pulmonary stenosis.



**Figure 6.** 30-day outcomes of developing mortality and hospitalization events in patients without PE stratified by DLM prediction. The patients without PE who were misidentified as PE by DLM were labeled “DLM-PE”. In contrast, those who were not were labeled “DLM-non PE”. CV, cardiovascular; DLM, deep learning model; PE, pulmonary stenosis.

classified cases, we found that some physical factors and comorbidities influenced PE recognition in the DLM. The results of the survival analysis indicated that patients without PE whose ECGs were recognized as PE by the DLM had higher mortality and hospitalization rates related to underlying comorbidities other than cardiovascular diseases.

Regarding conventional screening tools, examinations such as mammography and fecal occult blood tests (FOBTs) have similar costs and acquisition times as ECGs. In a previous report, digital mammography yielded a sensitivity ranging from 69% to 86% and specificity ranging from 57% to 94%.<sup>21</sup> FOBTs can achieve a sensitivity

or specificity above 90%, but at the expense of the other, which is lower than 40%.<sup>22</sup> Our DLM had a sensitivity and specificity of more than 70% for PE recognition. When adding histories, clinical symptoms and signs, and laboratory examinations, it can be effective in helping with the early identification of PE in either the ED or places where advanced tools are not available.

After a detailed review of ECGs in the human-machine competition, some morphological features were found to be helpful in PE recognition. In addition to prolonged QTc and negative T wave axis, a negative QRS complex axis seemed to play a role in PE recognition in the DLM, similar to the importance ranking in the XGB

model. Features such as tachycardia and RBBB, which are less important in the XGB and elastic net models, may also have been considered in the DLM. These findings may partially explain why the XGB and elastic net models had similar performances as the DLM, and highlight the blind spots in PE recognition by humans. However, it remains unclear which features were identified by the DLM and how the relative importance of each feature led to better performance of the DLM compared to the traditional machine learning algorithms.

Similar to the role of D-dimer, the manifestations of ECG provide unspecific clues to diagnose PE, which is also present in the ECGs of patients with other morbidities. Our stratified analysis showed that patients with older age, HF, atrial fibrillation, and COPD were more likely to be misclassified as PE by the DLM. These false-positive patients had a higher risk of all-cause mortality and non-CV hospitalization compared to the true-negative patients, indicating that the DLM learned to identify the PE-associated physiologic or anatomic abnormalities, including right ventricular overloading, tachycardia, and abnormalities of ST interval, during the training process. The abnormal ECG patterns revealed morbidities which may have included but are not limited to the diseases mentioned above in those patients. Taken together, these results may partly explained the higher rates of complications, non-CV hospitalizations, and mortality in the false-positive patients.

The concept of “point-of-care testing” has developed over the years. The use of DLMs in critical illness recognition, such as PE, can provide timely warnings to healthcare providers and allow rapid-response teams to become involved in patient care much earlier. The high availability of ECG makes our DLM-enhanced PE detection system especially helpful in places where healthcare resources, such as blood examinations and advanced imaging examinations are lacking. Additionally, a combination of ECG and laboratory data can be used to develop a rapid rule/out protocol of PE in the future. This can help to reduce unexpected complications and mortality during clinical practice. When not used for PE identification, DLMs can still provide a predictive value for clinical prognosis. Clinicians can adjust their therapeutic strategy and strength to achieve better control of patient comorbidities and reduce feature deaths.

Our study has some limitations. First, this was a sin-

gle-hospital retrospective study. A prospective study of multiple independent emergency services will be helpful in validating our DLM-enhanced PE detection system. Second, the number of patients with PE was small due to the rarity of the disease. Consequently, the small number of ECGs for DLM development and validation may have influenced the final performance of the DLM. Third, the number of physicians, especially attending physicians, who joined the competition was limited and may not represent actual human performance. Finally, traditional machine learning models revealed some relationships between the explainable features and ECG morphologies. The “black box” effect of exact ECG morphology identified by the DLM still remains.

## CONCLUSIONS

In summary, our proposed DLM-enhanced PE detection system was shown to be an effective and automatic tool to rapidly screen patients with potential PE in either in-hospital or out-of-hospital settings, and could be used to promptly alert first-line physicians. False-positive recognition in patients without PE could help healthcare professionals to predict the prognosis and help guide treatment strategies.

## ACKNOWLEDGEMENT

The authors would like to thank all colleagues who contributed to this study. The work was supported by the Tri-Service General Hospital Medical Research Foundation Grant TSGH-PH-E-111017.

## DECLARATION OF CONFLICT OF INTEREST

All the authors declare no conflict of interest.

## REFERENCES

1. Goldhaber SZ. Deep Venous Thrombosis and Pulmonary Thromboembolism, In: Jameson J, Fauci AS, Kasper DL, Harrison's Principles of Internal Medicine, 21th ed, McGraw-Hill Education,

- 2022:2091.
- Konstantinides SV, Meyer G, Becattini C, et al. 2019 ESC Guidelines for the diagnosis and management of acute pulmonary embolism developed in collaboration with the European Respiratory Society (ERS). *Eur Heart J* 2020;41:543-603.
  - Bach AG, Bandzauner R, Nansalmaa B, et al. Timing of pulmonary embolism diagnosis in the emergency department. *Thromb Res* 2016;137:53-7.
  - McGinn S, White PD. Acute cor pulmonale resulting from pulmonary embolism. *JAMA* 1935;104:1473.
  - Daniel KR, Courtney DM, Kline JA. Assessment of cardiac stress from massive pulmonary embolism with 12-lead ECG. *Chest* 2001;120:474-81.
  - Rodger M, Makropoulos D, Turek M, et al. Diagnostic value of the electrocardiogram in suspected pulmonary embolism. *Am J Cardiol* 2000;86:807-9.
  - Sreeram N, Cheriex EC, Smeets JLRM, et al. Value of the 12-lead electrocardiogram at hospital admission in the diagnosis of pulmonary embolism. *Am J Cardiol* 1994;73:298-303.
  - Liu WT, Lin CS, Tsao TP, et al. A deep-learning algorithm-enhanced system integrating electrocardiograms and chest X-rays for diagnosing aortic dissection. *Can J Cardiol* 2022;38:160-8.
  - Lee CC, Lin CS, Tsai CS, et al. A deep learning-based system capable of detecting pneumothorax via electrocardiogram. *Eur J Trauma Emerg Surg* 2022;48:3317-26.
  - Chang CH, Lin CS, Luo YS, et al. Electrocardiogram-based heart age estimation by a deep learning model provides more information on the incidence of cardiovascular disorders. *Front Cardiovasc Med* 2022;9:754909.
  - Liu WC, Lin CS, Tsai CS, et al. A deep-learning algorithm for detecting acute myocardial infarction. *EuroIntervention* 2021;17:765-73.
  - Lin C, Lin CS, Lee DJ, et al. Artificial intelligence assisted electrocardiography for early diagnosis of thyrotoxic periodic paralysis. *J Endocr Soc* 2021;5:bvab120.
  - Chang DW, Lin CS, Tsao TP, et al. Detecting digoxin toxicity by artificial intelligence-assisted electrocardiography. *Int J Environ Res Public Health* 2021;18:3839.
  - Liu WC, Lin C, Lin CS, et al. An artificial intelligence-based alarm strategy facilitates management of acute myocardial infarction. *J Pers Med* 2021;11:1149.
  - Lou YS, Lin CS, Fang WH, et al. Artificial intelligence-enabled electrocardiogram estimates left atrium enlargement as a predictor of future cardiovascular disease. *J Pers Med* 2022;12:315.
  - Lin CS, Lee YT, Fang WH, et al. Deep learning algorithm for management of diabetes mellitus via electrocardiogram-based glycosylated hemoglobin (ECG-HbA1c): a retrospective cohort study. *J Pers Med* 2021;11:725.
  - Lin CS, Lin C, Fang WH, et al. A deep-learning algorithm (ECG12Net) for detecting hypokalemia and hyperkalemia by electrocardiography: algorithm development. *JMIR Med Inform* 2020;8:e15931.
  - Lin C, Chau T, Lin CS, et al. Point-of-care artificial intelligence-enabled ECG for dyskalemia: a retrospective cohort analysis for accuracy and outcome prediction. *NPJ Digit Med* 2022;5:8.
  - Attia ZI, Harmon DM, Behr ER, et al. Application of artificial intelligence to the electrocardiogram. *Eur Heart J* 2021;42:4717-30.
  - Buuren SV, Groothuis-Oudshoorn K. Mice: multivariate imputation by chained equations in R. *Share spsurvey: Spatial Sampling Design and Analysis in R*. Dumelle M, Kincaid T, Olsen AR, Weber M. *J Stat Softw* 2011;45:1-67.
  - Phi XA, Tagliafico A, Houssami N, et al. Digital breast tomosynthesis for breast cancer screening and diagnosis in women with dense breasts - a systematic review and meta-analysis. *BMC Cancer* 2018;18:380.
  - Burch JA, Soares-Weiser K, St John DJ, et al. Diagnostic accuracy of faecal occult blood tests used in screening for colorectal cancer: a systematic review. *J Med Screen* 2007;14:132-7.

## SUPPLEMENTARY MATERIALS

Supplementary Table 1. Corresponding patient demographics of the training, validation, and testing sets

	Human-machine competition			“High D-dimer” subset		
	PE (n = 10)	Non-PE (n = 66)	p value	PE (n = 9)	Non-PE (n = 715)	p value
Demographic data						
Sex (male)	7 (53.8%)	44 (65.7%)	0.531*	6 (50.0%)	396 (53.0%)	0.836
Age (years)	63.0 ± 15.5	61.6 ± 24.1	0.901*	59.2 ± 18.2	71.0 ± 17.0	0.023*
BMI (kg/m <sup>2</sup> )	26.7 ± 4.0	22.5 ± 3.5	0.062*	26.0 ± 4.6	24.5 ± 8.3	0.355*
Disease histories						
AMI	1 (7.7%)	6 (9.0%)	1.000*	0 (0.0%)	53 (7.1%)	1.000*
Stroke	2 (15.4%)	12 (17.9%)	1.000*	2 (16.7%)	213 (28.5%)	0.525*
CAD	6 (46.2%)	28 (41.8%)	0.771	4 (33.3%)	258 (34.5%)	1.000*
HF	4 (30.8%)	16 (23.9%)	0.727*	3 (25.0%)	163 (21.8%)	0.731*
AF	1 (7.7%)	10 (14.9%)	0.682*	0 (0.0%)	91 (12.2%)	0.378*
DM	4 (30.8%)	21 (31.3%)	1.000*	2 (16.7%)	284 (38.0%)	0.228*
HTN	7 (53.8%)	33 (49.3%)	0.762	5 (41.7%)	416 (55.7%)	0.332
CKD	1 (7.7%)	8 (11.9%)	1.000*	0 (0.0%)	168 (22.5%)	0.079*
HLP	3 (23.1%)	19 (28.4%)	1.000*	2 (16.7%)	272 (36.4%)	0.228*
COPD	4 (30.8%)	18 (26.9%)	0.745*	4 (33.3%)	193 (25.8%)	0.520*
Laboratory data						
D-dimer (μg/L)	8035.5 ± 6420.2	2734.6 ± 5580.5	0.002*	8830.1 ± 6711.7	4529.2 ± 6966.9	0.001*
eGFR (mL/min/1.73 m <sup>2</sup> )	53.7 ± 17.9	71.8 ± 33.1	0.032*	67.3 ± 30.1	64.5 ± 41.7	0.737*
Cr (mg/dL)	1.9 ± 2.3	1.5 ± 1.5	0.114*	1.1 ± 0.3	1.9 ± 2.2	0.834*
BUN (mg/dL)	18.0 ± 8.9	30.3 ± 26.1	0.174*	16.1 ± 8.6	32.4 ± 28.6	0.009*
Na (mmol/L)	135.8 ± 3.6	135.5 ± 4.3	0.920*	136.3 ± 3.1	135.6 ± 6.3	0.906*
K (mmol/L)	3.9 ± 0.5	4.1 ± 0.7	0.844*	3.8 ± 0.3	4.1 ± 0.8	0.402*
Cl (mmol/L)	106.8 ± 8.0	99.7 ± 5.4	0.128*	110.3 ± 2.3	101.6 ± 7.1	0.012*
tCa (mg/dL)	8.2 ± 0.6	8.7 ± 0.5	0.014*	8.2 ± 0.6	8.5 ± 0.6	0.324*
Mg (mg/dL)	2.1 ± 0.3	2.1 ± 0.3	0.731*	2.0 ± 0.3	2.1 ± 0.4	0.512*
TnI (pg/mL)	113.7 ± 172.7	147.9 ± 672.0	0.016*	258.0 ± 413.3	874.3 ± 6456.4	0.045*
CK (U/L)	98.5 ± 108.3	122.4 ± 133.3	0.327*	103.3 ± 113.4	225.7 ± 1011.0	0.493*
BNP (ng/mL)	941.4 ± 1229.4	816.1 ± 1012.4	0.503*	563.5 ± 633.9	637.4 ± 1056.8	0.690*
GLU (g/dL)	150.9 ± 75.7	147.0 ± 62.3	0.735*	135.5 ± 70.3	166.9 ± 99.8	0.105*
Hb (g/dL)	13.3 ± 2.1	12.8 ± 2.4	0.458*	14.1 ± 0.9	11.6 ± 2.6	0.001*
WBC (10 <sup>3</sup> /μL)	9.4 ± 2.9	8.8 ± 3.2	0.390*	9.4 ± 2.6	11.2 ± 20.5	0.912*
PLT (10 <sup>3</sup> /μL)	209.7 ± 95.2	200.8 ± 66.5	0.833*	182.7 ± 49.3	231.4 ± 108.3	0.119*
AST (U/L)	23.5 ± 10.1	33.6 ± 57.6	0.704*	23.2 ± 10.3	50.5 ± 125.7	0.154*
ALT (U/L)	14.2 ± 7.9	39.9 ± 50.0	0.175*	14.6 ± 6.9	33.7 ± 79.8	0.263*
TG (g/L)	127.4 ± 46.2	174.7 ± 135.1	0.947*	119.8 ± 45.3	111.2 ± 70.1	0.323*
TC (g/L)	164.9 ± 45.8	150.0 ± 42.0	0.266*	170.2 ± 40.3	141.3 ± 44.8	0.020*

\* p value calculated with  $n < 25$  in continuous variables or by Fisher's exact test for categorical variables.

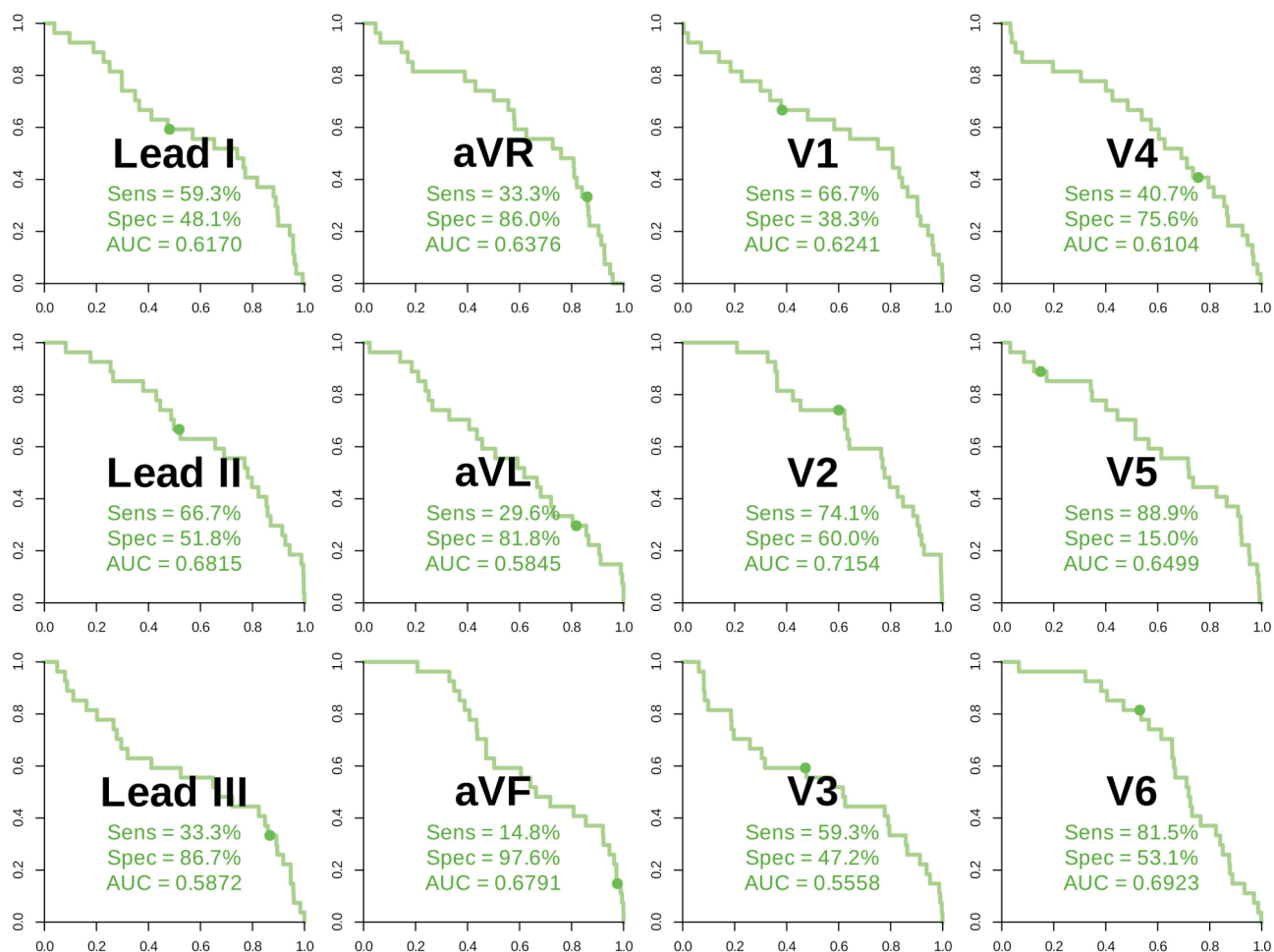
AF, atrial fibrillation; ALT, alanine aminotransferase; AMI, acute myocardial infarction; AST, aspartate aminotransferase; BMI, body mass index; BNP, brain natriuretic peptide; BUN, blood urea nitrogen; CAD, coronary artery disease; Cl, chloride; CK, creatine kinase; CKD, chronic kidney disease; COPD, chronic obstructive pulmonary disease; Cr, creatinine; DM, diabetes mellitus; eGFR, estimated glomerular filtration rate; GLU, fasting glucose; Hb, hemoglobin; HF, heart failure; HLP, hyperlipidemia; HTN, hypertension; K, potassium; Mg, Magnesium; Na, sodium; PE, pulmonary embolism; PLT, platelet; TC, total cholesterol; tCa, total calcium; TG, triglyceride; TnI, troponin I; WBC, white blood cell count.

**A: Pulmonary embolism ECGs (n = 24)**

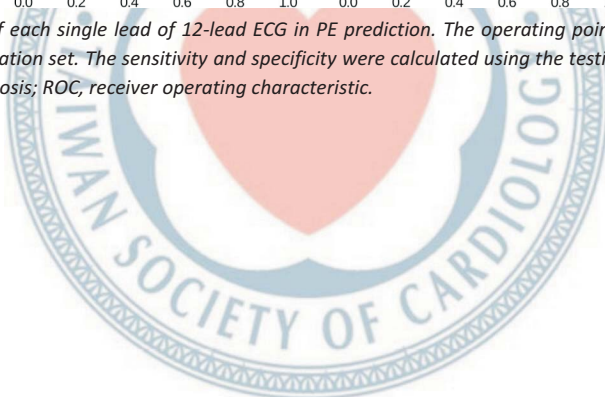
risk	risk	risk	norm	risk	0.000
norm	norm	norm	norm	norm	0.000
norm	norm	norm	norm	risk	0.001
norm	norm	norm	norm	norm	0.003
norm	risk	norm	risk	risk	0.004
norm	risk	risk	norm	norm	0.007
risk	risk	norm	norm	risk	0.013
norm	norm	norm	risk	norm	0.028
norm	risk	risk	risk	risk	0.029
norm	risk	risk	norm	norm	0.038
risk	risk	risk	risk	risk	0.049
risk	risk	risk	risk	risk	0.057
risk	risk	risk	norm	risk	0.059
risk	risk	risk	risk	risk	0.075
risk	risk	norm	risk	norm	0.080
risk	risk	risk	norm	risk	0.083
risk	risk	risk	norm	risk	0.087
risk	norm	norm	norm	risk	0.096
risk	risk	risk	norm	risk	0.103
risk	risk	risk	risk	risk	0.111
norm	risk	norm	norm	risk	0.137
risk	risk	risk	risk	norm	0.156
risk	risk	risk	risk	risk	0.187
risk	risk	risk	risk	risk	0.280
CV-R2	ER-R3	CV-V9	ER-R4	CV-R3	DLM

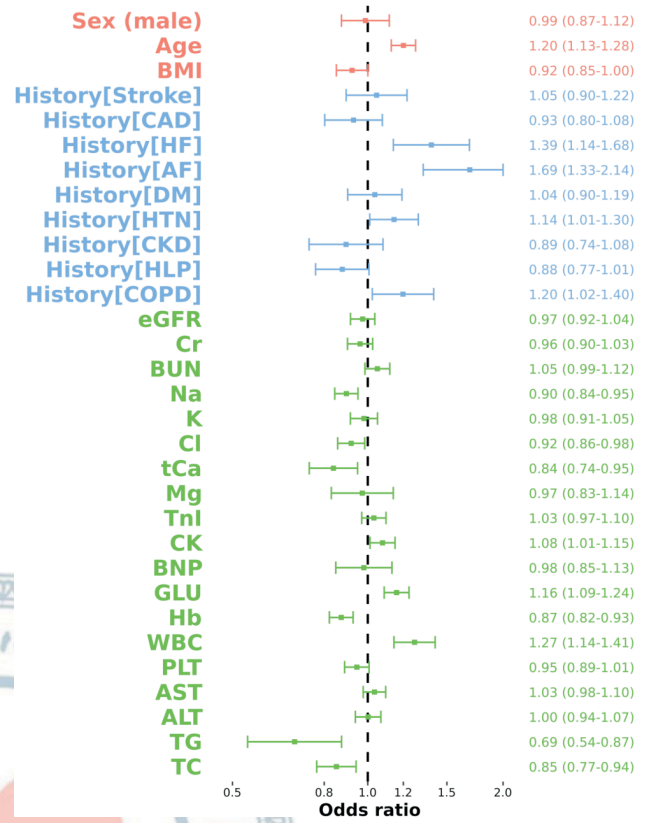
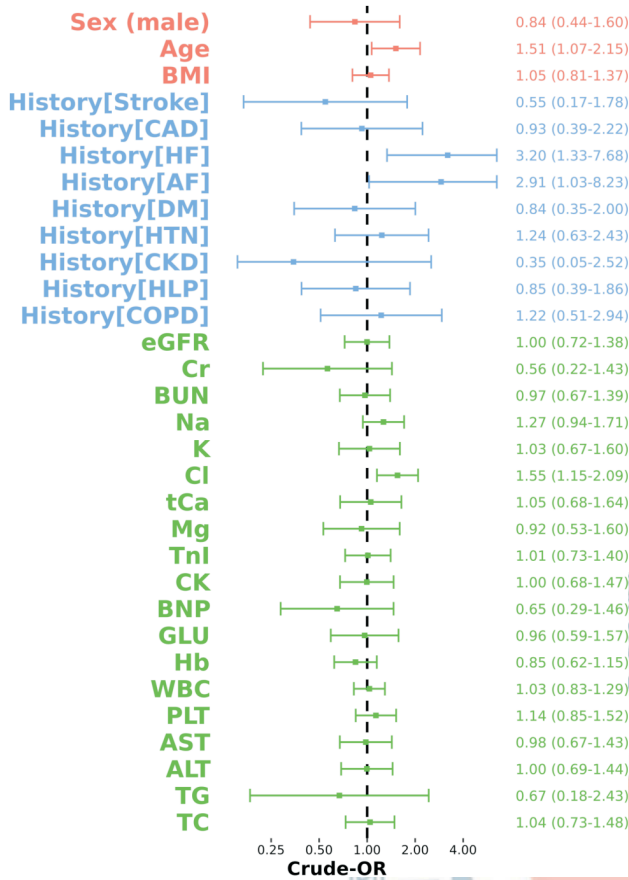
**B: non-pulmonary embolism ECGs (n = 76)**

**Supplementary Figure 1.** The detailed answers of the human-machine competition from each human expert and the DLM prediction. ECGs that were thought to be PE were labeled “risk” with a red background, and those that were not were labeled “normal” with a green background. The DLM prediction value of each ECG is labeled in the “DLM” column. DLM, deep learning model; ECG, electrocardiogram; PE, pulmonary stenosis



**Supplementary Figure 2.** The ROC of each single lead of 12-lead ECG in PE prediction. The operating point was selected based on the maximum Youden's index obtained from the validation set. The sensitivity and specificity were calculated using the testing set. AUC, area under the curve; ECG, electrocardiogram; PE, pulmonary stenosis; ROC, receiver operating characteristic.





**Supplementary Figure 3.** Forest plot presents the stratified analysis of the characteristics of patients with and without PE in the training set. AF, atrial fibrillation; ALT, alanine aminotransferase; AST, aspartate aminotransferase; BMI, body mass index; BNP, brain natriuretic peptide; BUN, blood urea nitrogen; CAD, coronary artery disease; COPD, chronic obstructive pulmonary disease; CK, creatine kinase; DM, diabetes mellitus; eGFR, estimated glomerular filtration rate; GLU, fasting glucose; Hb, hemoglobin; HF, heart failure; HLP, hyperlipidemia; HTN, hypertension; K, potassium; Mg, Magnesium; Na, sodium; PLT, platelet; TC, total cholesterol; TG, triglyceride; Tnl, troponin I; WBC, white blood cell count.

**Supplementary Figure 4.** Forest plot presents the stratified analysis of patient characteristics in non-PE cases within the testing set according to DLM prediction. AF, atrial fibrillation; ALT, alanine aminotransferase; AST, aspartate aminotransferase; BMI, body mass index; BNP, brain natriuretic peptide; BUN, blood urea nitrogen; CAD, coronary artery disease; COPD, chronic obstructive pulmonary disease; CK, creatine kinase; CKD, chronic kidney disease; Cl, chloride; Cr, creatinine; DM, diabetes mellitus; eGFR, estimated glomerular filtration rate; GLU, fasting glucose; Hb, hemoglobin; HF, heart failure; HLP, hyperlipidemia; HTN, hypertension; K, potassium; Mg, Magnesium; Na, sodium; PLT, platelet; TC, total cholesterol; tCa, total calcium; TG, triglyceride; Tnl, troponin I; WBC, white blood cell count.