

# Large Language Models versus Classical Machine Learning: Performance in COVID-19 Mortality Prediction Using High-Dimensional Tabular Data

**Running Title:** Large Language Models vs Classical Machine Learning on Structured Data

Mohammadreza Ghaffarzadeh-Esfahani<sup>1,2</sup>, Mahdi Ghaffarzadeh-Esfahani<sup>2</sup>, Arian Salahi-Niri<sup>1</sup>, Hossein Toreyhi<sup>1</sup>, Zahra Atf<sup>3</sup>, Amirali Mohsenzadeh-Kermani<sup>2</sup>, Mahshad Sarikhani<sup>4</sup>, Zohreh Tajabadi<sup>5</sup>, Fatemeh Shojaeian<sup>6</sup>, Mohammad Hassan Bagheri<sup>2</sup>, Aydin Feyzi<sup>7</sup>, Mohammadamin Tarighatpayma<sup>4</sup>, Narges Gazmeh<sup>7</sup>, Fateme Heydari<sup>7</sup>, Hossein Afshar<sup>7</sup>, Amirreza Allahgholipour<sup>7</sup>, Farid Alimardani<sup>7</sup>, Ameneh Salehi<sup>4</sup>, Naghmeh Asadimanesh<sup>4</sup>, Mohammad Amin Khalafi<sup>4</sup>, Hadis Shabanipour<sup>7</sup>, Ali Moradi<sup>7</sup>, Sajjad Hossein Zadeh<sup>7</sup>, Omid Yazdani<sup>4</sup>, Romina Esbati<sup>4</sup>, Moozhan Maleki<sup>7</sup>, Danial Samiei Nasr<sup>4</sup>, Amirali Soheili<sup>4</sup>, Hossein Majlesi<sup>4</sup>, Saba Shahsavan<sup>4</sup>, Alireza Soheilipour<sup>4</sup>, Nooshin Goudarzi<sup>1</sup>, Erfan Taherifard<sup>8</sup>, Hamidreza Hatamabadi<sup>9</sup>, Jamil S. Samaan<sup>10</sup>, Thomas Savage<sup>11</sup>, Ankit Sakhuja<sup>12</sup>, Ali Soroush<sup>12</sup>, Girish Nadkarni<sup>12</sup>, Ilad Alavi Darazam<sup>13,14\*</sup>, Mohamad Amin Pourhoseingholi<sup>\*15</sup>, Seyed Amir Ahmad Safavi-Naini<sup>\*1,12</sup>

1. Research Institute for Gastroenterology and Liver Diseases, Shahid Beheshti University of Medical Sciences, Tehran, Iran
2. Faculty of Medicine, Isfahan University of Medical Sciences, Isfahan, Iran
3. Faculty of Business and Information Technology, Ontario Tech University, Oshawa, Canada
4. School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran.
5. Digestive Disease Research Institute, Tehran University of Medical Sciences, Tehran, Iran
6. Department of Surgery, The Johns Hopkins University, Baltimore, MD, USA
7. Student Research Committee, School of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran
8. MPH department, Shiraz University of Medical Sciences, Shiraz, Iran
9. Department of Emergency Medicine, School of Medicine, Safety Promotion and Injury Prevention Research Center, Imam Hossein Hospital, Shahid Beheshti University of Medical Sciences, Tehran, Iran
10. Karsh Division of Gastroenterology and Hepatology, Cedars-Sinai Medical Center, 8700 Beverly Blvd, Los Angeles, CA 90048.
11. Department of Medicine, Stanford University, Stanford, California, USA
12. Division of Data Driven and Digital Health (D3M), The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA
13. Infectious Diseases and Tropical Medicine Research Center, Shahid Beheshti University of Medical Sciences, Tehran, Iran
14. Department of Infectious Diseases, Loghman Hakim Hospital, Shahid Beheshti University of Medical Sciences, Tehran, Iran
15. National Institute for Health and Care Research (NIHR), Nottingham Biomedical Research Centre, Hearing Sciences, Mental Health and Clinical Neurosciences, School of Medicine, University of Nottingham, Nottingham, UK

\* **Correspondence to:** Seyed Amir Ahmad Safavi-Naini ([sdamirsa@ymail.com](mailto:sdamirsa@ymail.com)), Mohamad Amin Pourhoseingholi ([aminphg@gmail.com](mailto:aminphg@gmail.com)), and Ilad Alavi Darazam ([ilad13@yahoo.com](mailto:ilad13@yahoo.com))

**Funding:** None

**COI Disclosures:** All authors declare no conflict of interest related to this work.

**Data Transparency Statement:** Code is available at: <https://github.com/mohammad-gh009/Large-Language-Models-vs-Classical-Machine-learning> and [https://github.com/Sdamirsa/Tehran\\_COVID\\_Cohort](https://github.com/Sdamirsa/Tehran_COVID_Cohort). The datasets generated and/or analyzed during the current study are not publicly available due to privacy concerns and ethical restrictions, but are available from the corresponding author on reasonable request ([sdamirsa@ymail.com](mailto:sdamirsa@ymail.com)).

Manuscript Pages: 42; Figures: 5; Tables: 2; Supplementary File Pages: 14

## Authors' Detail

Full Name	Position	Email; ORCID	Affiliation	Contribution
Mohammadreza Ghaffarzadeh-Esfahani	Medical Student	<a href="mailto:mregahafarzadeh@gmail.com">mregahafarzadeh@gmail.com</a> ; 0009-0009-9322-5471	(a) Research Institute for Gastroenterology and Liver Diseases, Shahid Beheshti University of Medical Sciences, Tehran, Iran (b) Faculty of Medicine, Isfahan University of Medical Sciences, Isfahan, Iran	Conceptualization, Methodology, Programming, Investigation, Writing Original Draft
Mahdi Ghaffarzadeh-Esfahani	Medical Student	<a href="mailto:Mahdighafarzadeh313@gmail.com">Mahdighafarzadeh313@gmail.com</a> 0009-0006-6395-3273	Faculty of Medicine, Isfahan University of Medical Sciences, Isfahan, Iran	Investigation, Methodology
Aryan Salahi-Niri	BS	<a href="mailto:Arian.slh12@gmail.com">Arian.slh12@gmail.com</a> 0000-0002-6772-0972	Research Institute for Gastroenterology and Liver Diseases, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation, Methodology
Hossein Toreyhi	MD, Research Fellow	<a href="mailto:hoseinto@gmail.com">hoseinto@gmail.com</a>	Research Institute for Gastroenterology and Liver Diseases, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Writing Original Draft, Methodology
Zahra Atf	MS, visiting scholar	0000-0003-0642-4341; <a href="mailto:Zahra.Atf@ontariotechu.ca">Zahra.Atf@ontariotechu.ca</a>	Faculty of Business and Information Technology, Ontario Tech University, Oshawa, Canada	Investigation, Programming
Amirali Mohsenzadeh-Kermani	Medical Student	<a href="mailto:amkermani80@gmail.com">amkermani80@gmail.com</a> 0000-0002-7220-3355	Faculty of Medicine, Isfahan University of Medical Sciences, Isfahan, Iran	Investigation
Mahshad Sarikhani	MD	<a href="mailto:mahshadsarikhani9696@gmail.com">mahshadsarikhani9696@gmail.com</a> 0000-0002-0475-3064	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran.	Investigation
Zohreh Tajabadi	MD	<a href="mailto:Zohreh.tajabadi@gmail.com">Zohreh.tajabadi@gmail.com</a> 0000-0002-9854-7260	Digestive Disease Research Institute, Tehran University of Medical Sciences, Tehran, Iran	Investigation
Fatemeh Shojaeian	MD, research fellow	0000-0001-5972-9953; <a href="mailto:fshojae1@jhmi.edu">fshojae1@jhmi.edu</a>	Department of Surgery, The Johns Hopkins University, Baltimore, MD, USA	Investigation
Mohammad Hassan Bagheri	Medical Student	<a href="mailto:Hassanbagheri2002@gmail.com">Hassanbagheri2002@gmail.com</a> 0009-0000-5578-3215	Student Research Committee, School of Medicine, Isfahan University of Medical Sciences, Isfahan, Iran	Investigation
Aydin Feyzi	BS	<a href="mailto:Aydin.Feyzi.Gh@gmail.com">Aydin.Feyzi.Gh@gmail.com</a>	Student Research Committee, school of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Mohammadamin Tarighatpayma	Medical Student	<a href="mailto:tarighatamin@gmail.com">tarighatamin@gmail.com</a> ; 0009-0002-0900-7770	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Narges Gazmeh	BS	<a href="mailto:narges.gazmeh@gmail.com">narges.gazmeh@gmail.com</a>	Student Research Committee, school of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Fateme Heydari	Medical Student	<a href="mailto:fatemeheydari96@sbmu.ac.ir">fatemeheydari96@sbmu.ac.ir</a>	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Hossein Afshar	BS	<a href="mailto:hosein.afshar1119@gmail.com">hosein.afshar1119@gmail.com</a>	Student Research Committee, school of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Amirreza Allahgholipour	BS	<a href="mailto:aallahgholipourk@sbmu.ac.ir">aallahgholipourk@sbmu.ac.ir</a>	Student Research Committee, school of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Farid Alimardani	BS	<a href="mailto:alimardani80735344@gmail.com">alimardani80735344@gmail.com</a>	Student Research Committee, school of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation

Ameneh Salehi	MD	<a href="mailto:salehiamene10@gmail.com">salehiamene10@gmail.com</a>	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Naghmeh Asadimanesh	Medical Student	<a href="mailto:naghmeh_manesh@yahoo.com">naghmeh_manesh@yahoo.com</a>	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Mohammad Amin Khalafi	MD, Research Fellow	<a href="mailto:aminkhalafi1996@gmail.com">aminkhalafi1996@gmail.com</a>	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation, Project Administration
Hadis Shabanipour	BS	<a href="mailto:hadisshabanipour@gmail.com">hadisshabanipour@gmail.com</a>	Student Research Committee, school of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Ali Moradi	BS	<a href="mailto:aliaaaa.m2001@gmail.com">aliaaaa.m2001@gmail.com</a>	Student Research Committee, school of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Sajjad Hossein Zadeh	BS	<a href="mailto:sajjad.asnz@gmail.com">sajjad.asnz@gmail.com</a>	Student Research Committee, school of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Omid Yazdani	MD	<a href="mailto:Omidyazdn@yahoo.com">Omidyazdn@yahoo.com</a> ; 0000-0002-0798-3308	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran.	Investigation
Romina Esbati	MD	<a href="mailto:Romines9573@gmail.com">Romines9573@gmail.com</a> ; 0000-0002-7161-5948	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran.	Investigation
Moozhan Maleki	BS	<a href="mailto:moozhan.malaki73@gmail.com">moozhan.malaki73@gmail.com</a>	Student Research Committee, School of Nursing and Midwifery, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Danial Samiei Nasr	MD	<a href="mailto:Danial.samiei9@gmail.com">Danial.samiei9@gmail.com</a>	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Amirali Soheili	MD	<a href="mailto:amiralisoheili1375@gmail.com">amiralisoheili1375@gmail.com</a>	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Hossein Majlesi	Medical Student	<a href="mailto:majlesihossein@gmail.com">majlesihossein@gmail.com</a>	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran.	Investigation
Saba Shahsavan	Medical Student	<a href="mailto:shahsavansaba@gmail.com">shahsavansaba@gmail.com</a>	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Alireza Soheilipour	Medical Student	<a href="mailto:alirzsp77@gmail.com">alirzsp77@gmail.com</a>	School of Medicine, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Nooshin Goudarzi	Medical Student	<a href="mailto:nooshingz79@gmail.com">nooshingz79@gmail.com</a> ; 0000-0001-7252-467X	Research Institute for Gastroenterology and Liver Diseases, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Erfan Taherifard	MD	<a href="mailto:Erfantaherifard@gmail.com">Erfantaherifard@gmail.com</a> ; 0000-0002-9101-0321	MPH department, Shiraz University of Medical Sciences, Shiraz, Iran	Investigation, Validation
Hamidreza Hatamabadi	MD, Professor of Emergency Medicine	0000-0002-9085-8806; <a href="mailto:hhatamabadi@yahoo.com">hhatamabadi@yahoo.com</a>	Safety Promotion and Injury Prevention Research Center, Imam Hossein Hospital, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Investigation
Jamil S. Samaan, MD	MD, Gastroenterology clinical fellow	<a href="mailto:jamil.samaan@cshs.org">jamil.samaan@cshs.org</a> ; 0000-0002-6191-2631	David Geffen School Medicine University of California, Los Angeles, 10833 Le Conte Ave, Los Angeles, CA 90095	Reviewing and Editing the Manuscript
Thomas Savage	MD, Assistant Professor	<a href="mailto:tsavage@stanford.edu">tsavage@stanford.edu</a> ; 0000-0003-4828-5802	Department of Medicine, Stanford University, Stanford, California	Reviewing and Editing the Manuscript

Ankit Sakhuja	MD MPH, Associate Professor of Nephrology	Dr.a.sakhuja@gmail.com	Division of Data Driven and Digital Health (D3M), The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA	Reviewing and Editing the Manuscript
Ali Soroush	MD, Assistant Professor of Gastroenterology	<a href="mailto:ali.soroush.2012@gmail.com">ali.soroush.2012@gmail.com</a> ; 0000-0001-6900-5596	Division of Data Driven and Digital Health (D3M), The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA	Methodology, Validation, Reviewing and Editing the Manuscript
Girish Nadkarni	MD MPH, Associate Professor of Nephrology	girish.nadkarni@mountsinai.org	Division of Data Driven and Digital Health (D3M), The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA	Reviewing and Editing the Manuscript
Ilad Alavi Darazam	MD, Associate Professor of Infectious Diseases	<a href="mailto:ilad13@yahoo.com">ilad13@yahoo.com</a> ; 0000-0002- 4440-335X	(a) Infectious Diseases and Tropical Medicine Research Center, Shahid Beheshti University of Medical Sciences, Tehran, Iran (b) Department of Infectious Diseases, Loghman Hakim Hospital, Shahid Beheshti University of Medical Sciences, Tehran, Iran	Data Acquisition, Reviewing and Editing the Manuscript, Administration, Supervision
Mohamad Amin Pourhoseingholi	PhD, Associate Professor of Biostatistics	aminphg@gmail.com	National Institute for Health and Care Research (NIHR), Nottingham Biomedical Research Centre, Hearing Sciences, Mental Health and Clinical Neurosciences, School of Medicine, University of Nottingham, Nottingham, UK	Data Acquisition, Reviewing and Editing the Manuscript, Administration, Supervision, Validation
Seyed Amir Ahmad Safavi- Naini	MD, Research Fellow	<a href="mailto:sdamirsa@ymail.com">sdamirsa@ymail.com</a> ; 0000-0001- 9295-9283	(a) Research Institute for Gastroenterology and Liver Diseases, Shahid Beheshti University of Medical Sciences, Tehran, Iran (b) Division of Data Driven and Digital Health (D3M), The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA	Conceptualization, Methodology, Programming, Data Curation, Writing and Editing the Original Draft, Project Administration, Supervision

## Abstract

**Background:** This study aimed to evaluate and compare the performance of classical machine learning models (CMLs) and large language models (LLMs) in predicting mortality associated with COVID-19 by utilizing a high-dimensional tabular dataset.

**Materials and Methods:** We analyzed data from 9,134 COVID-19 patients collected across four hospitals. Seven CML models, including XGBoost and random forest (RF), were trained and evaluated. The structured data was converted into text for zero-shot classification by eight LLMs, including GPT-4 and Mistral-7b. Additionally, Mistral-7b was fine-tuned using the QLoRA approach to enhance its predictive capabilities.

**Results:** Among the CML models, XGBoost and RF achieved the highest accuracy, with F1 scores of 0.87 for internal validation and 0.83 for external validation. In the LLM category, GPT-4 was the top performer with an F1 score of 0.43. Fine-tuning Mistral-7b significantly improved its recall from 1% to 79%, resulting in an F1 score of 0.74, which was stable during external validation.

**Conclusion:** While LLMs show moderate performance in zero-shot classification, fine-tuning can significantly enhance their effectiveness, potentially aligning them closer to CML models. However, CMLs still outperform LLMs in high-dimensional tabular data tasks.

**Keywords:** COVID-19 mortality, Large language models, Classical machine learning, Structured data, Zero-shot classification, Fine-tuning

## Highlights

- We tested the predictive power of large language models (LLMs) on tabular data compared with classical machine learning (CML).
- XGBoost and GPT-4 were the top performers among the CML and LLMs, achieving accuracies of 87% and 62%, respectively.
- Our resource-efficient fine-tuning increased the recall of Mistral-7b from 1% to 79%, reaching an accuracy of 72% and outperforming GPT-4.
- The fine-tuned Mistral-7b logic (global feature impact) was similar to that of the CMLs and XGBoost.

## Introduction

The rapid advancement of large language models (LLMs) has revolutionized their practical applications across various domains, including medicine. These sophisticated models, trained on vast datasets, excel in a wide array of natural language processing tasks, demonstrating remarkable adaptability in assimilating specialized information from diverse medical fields (1). While primarily designed for next-word prediction, LLMs have emerged as powerful, evidence-based knowledge assistants for healthcare providers, offering valuable insights and support in clinical decision-making processes (2–4). While their main training centers on predicting the next word, LLMs can act as evidence-based knowledge helpers for healthcare providers, offering valuable insights and assistance (5).

In medical and clinical practice, machine learning models, particularly **classical machine learning (CML) models**, have gained significant traction in predicting patient outcomes, prognoses, and mortality rates. These models typically employ supervised and unsupervised learning methods, which **primarily utilize structured data** (6). However, clinical datasets often present a complex interplay of structured and unstructured information, with clinical notes serving as prime examples of the latter. Traditionally, **patient information management via machine learning has followed a two-step approach: transforming unstructured textual data into a structured format**, followed by **training CML models on these structured datasets**. This process, however, often leads to potential information loss and introduces complexities in model deployment, hindering practical application in clinical settings (7).

While the efficacy of LLMs in handling unstructured text is well documented (8), their **performance in handling structured data and their comparative effectiveness against CML models**

remain a critical area of investigation. This is particularly relevant given that much of the historical medical data are stored in structured formats that are often difficult to integrate (9). **Table 1** summarizes previous studies comparing the performance of LLMs and CML approaches in medicine (10–14). Studies reported varied results, primarily due to differences in evaluated tasks (number of input features, sample size, and prediction complexity) and transformation techniques (e.g., transforming tables into textual prompts). However, they focus on tasks with a limited number of features ( $<12$ ), fail to represent real-world medical decisions, and train instances for the models ( $<1000$ ), limiting the CMLs to reach their maximum performance.

Our study aims to address this knowledge gap by evaluating LLMs' predictive capabilities in the context of COVID-19 mortality prediction via a high-dimensional dataset and simple table-to-text transformation. By utilizing a sufficient number of training instances, we provide the opportunity for CMLs to reach their maximum performance, enabling a more robust comparison with LLMs. This investigation is designed to provide insights into CML versus LLM comparisons in real-world, time-sensitive, and complex clinical tasks.



**Table 1. Summary of studies comparing the performance of large language models and classical machine learning methods in medicine using structured data**

First Author; Year	Aim and Task	Dataset (Sample Size; #feature)	Transformation Techniques	Model, Experiment, and Metric	Zero-shot performance	Training Size: Performance
Hegselmann; 2023 (TabLLM) (10)	To transform table-to-text for binary classification of coronary artery disease and diabetes	Diabetes (768; #7); Heart (918; #11)	Template; Billion parameter LLM; Million parameter LLM			
				TabLLM-Diabetes (AUC)	0.82	32: 0.68 512: 0.78
				XGBoost-Diabetes (AUC)	-	32: 0.69 512: 0.80
				TabLLM-Heart (AUC)	0.54	32: 0.87 512: 0.92
				XGBoost-Heart (AUC)	-	32: 0.88 512: 0.92
Wang; 2024 (MediTab) (11)	To evaluate MediTab (GPT3.5) on seven medical classification tasks and compare it with TabLLM and CML	Seven datasets of breast, lung, and colorectal cancer from clinical trials (average 1451 ranging from 53 to 2968; average #3 categorical, #15 binaries; #7 numerical)	BioBERT-based model fine-tuned on transformation and GPT3.5 for sanity check			
				MediTab (Average AUC)	0: 0.82	200: 0.84
				XGBoost (Average AUC)	10: 0.64	200: 0.79
Cui; 2024 (EHR-CoAgent) (12)	To investigate the efficacy of LLMs-based disease prediction using structured EHR data generated from clinical encounters (MIMIC: acute care condition in the next hospital visit; CRADLE: CVD in diabetic patients)	MIMIC-III (11,353; #?); CRADLE (34,404; #?)	Disease, medication, and procedure codes by mapping the code value to code name (+ prompt engineering techniques)			
				EHR-CoAgent-GPT4 – MIMIC (Accuracy; F1)	0.79; 0.73	-
				GPT-4 – MIMIC (Accuracy; F1)	ZSC: 0.51; 52% Prompt engineered: 0.62; 0.58	Few-shot (N=6): 0.65; 0.64

				GPT-3.5 – MIMIC (Accuracy; F1)	ZSC: 0.78; 0.68 Prompt engineered: 0.72; 0.42	Few-shot (N=6): 0.76; 0.63
				RF – MIMIC (Accuracy; F1)	N=6: 0.69; 0.63	11,353: 0.78; 0.70
				LR – MIMIC (Accuracy; F1)	N=6: 0.48; 0.56	11,353: 0.79; 0.73
				DT – MIMIC (Accuracy; F1)	N=6: 0.71; 0.51	11,353: 0.81; 0.76
				EHR-CoAgent-GPT4 – CRADLE (Accuracy; F1)	0.70; 0.60	-
				GPT-4 – CRADLE (Accuracy; F1)	ZSC: 0.21; 0.22; Prompt engineered: 0.30; 0.29	Few-shot (N=6): 0.41; 0.40
				GPT-3.5 – CRADLE (Accuracy; F1)	ZSC: 0.56; 0.52 Prompt engineered: 0.62; 0.54	Few-shot (N=6): 0.40; 0.40
				RF – CRADLE (Accuracy; F1)	N=6: 0.66; 0.51	34,404: 0.80; 0.57
				LR – CRADLE (Accuracy; F1)	N=6: 0.54; 0.48	34,404: 0.80; 0.59
				DT – CRADLE (Accuracy; F1)	N=6: 0.31; 0.31	34,404: 0.80; 0.52
Nazary; 2024 (XAI4LLM) (13)	To evaluate the diagnostic accuracy and risk factors, including gender bias and false negative rates using LLM. In addition, a comparison with CML approaches	Heart Disease Dataset (920; #11)	Using feature name-values or transforming to simple manual textual template			
				Best XAILLM: LLM+RF (F1)	ZSC: 0.741	
				XGBoost (F1)	-	920: 0.91
				RF (F1)	-	920: 0.74

**Footnote:** Cui et al.’s and Nazary et al.’s studies were preprint publications.

## Methods

### 2.1 Ethical Consideration

The study was approved by the Institutional Review Board (IRB) of Shahid Beheshti University of Medical Sciences (IR.SBMU.RIGLD.REC.004 and IR.SBMU.RIGLD.REC.1399.058). The IRB exempted this study from informed consent. Data were pseudonymized before analysis; patients' confidentiality and data security were prioritized at all levels. The study was completed under the Helsinki Declaration (2013) guidelines. During the generation of LLM responses, using the OpenAI API and Poe Web interface, we opted out of training on OpenAI and used no training-use models in Poe to maintain the data safety of patient information.

### 2.2 Study Aim and Experimental Summary

The objective of this research is to evaluate the efficacy of CMLs in comparison to LLMs, utilizing a dataset characterized by high-dimensional tabular data. We employed a previously compiled dataset and focused our experimental efforts on the task of classifying COVID-19 mortality. As illustrated in Figure 1, the primary experiment encompasses the following:

- Assessment of the performance of seven CML models on both internal and external test sets
- The assessment of eight LLMs and two pretrained language models on the test set.
- Assessment of a trained LLM's performance on both internal and external tests.

Additionally, we investigate the performance of models necessitating training (CML and trained LLM) across varying sample sizes, coupled with an elucidation of model prediction mechanisms through SHAP analysis.

## 2.3 Study Context, Data Collection, and Dataset

This study was conducted as part of the Tehran COVID-19 cohort, which included four tertiary centers with dedicated COVID-19 wards and ICUs in Tehran, Iran. The study period was from March 2020 to May 2023 and included two phases of data collection. The protocol and results of the first phase have been published previously. The four COVID-19 peaks during this period covered the alpha, beta, delta, and Omicron variants.

All admitted patients with a positive swab test during the first two days of admission or those with CT scans and clinical symptoms were included in the study. A medical team collected the patients' symptoms, comorbidities, habitual history, vital signs at admission, and treatment protocol through the hospital information system and reviewed the medical records. Laboratory values during the first and second days of admission were collected and sorted from the hospitals' electronic laboratory records via Python (Python Software Foundation, Python Language Reference, version 4. Available at: <http://www.python.org>). Patients with a negative PCR result in the first two days of admission or with one missing clinical record in the HIS were excluded.

The dataset included the records of 9,134 patients with COVID-19. The data were filtered to include demographic information, comorbidities, vital signs, and laboratory results collected at the time of admission (first two days).

## 2.4. Data Preprocessing

### 2.4.1 Imputing and normalization

The features in the dataset were divided into categorical and numerical categories. To address the missing values in the numerical features, we used an iterative imputer from the scikit-learn library. This method employs iterative prediction for each feature, considering the multiple imputation by

chained equations (MICE) method (16,17). Missing values in the categorical features were imputed via KNN from the scikit-learn library. For optimal model performance, the dataset was normalized via a standard scaler (18). These preprocessing steps were executed independently for the input features of the training, test, and external validation sets, ensuring a consistent approach for handling missing values across the experimental sets without information leakage.

### 2.4.2 Feature Selection

The dataset comprised 81 on-admission features. The dataset was separated into external and internal validations using patient hospitals. Patients from Hospital-4 were used for external validation, whereas patients from the remaining hospitals were used for internal validation. For internal validation, we split the data with a test size of 20% and allocated 80% for training.

The output features in this study include "in-hospital mortality," "ICU admission," and "intubation," with a focus solely on "hospital mortality" as the targeted feature, excluding other output features. Of the 81 features initially available, 76 were employed for training, comprising 53 categorical features and the remaining numerical values. During data wrangment, two duplicate features were dropped.

We strategically employed the lasso method for feature selection because of its effectiveness in handling high-dimensional data. The Lasso method introduces regularization by adding a penalty term to the linear regression objective function, which encourages sparsity in the feature coefficients(15,16). This approach proved to be superior to alternative methods, facilitating notable enhancements in our results. Through the application of Lasso, we derived a refined dataset that highlighted the most impactful features on the basis of their importance, aiding dimensionality reduction. We subsequently ranked and selected the top 40 features for further analyses.

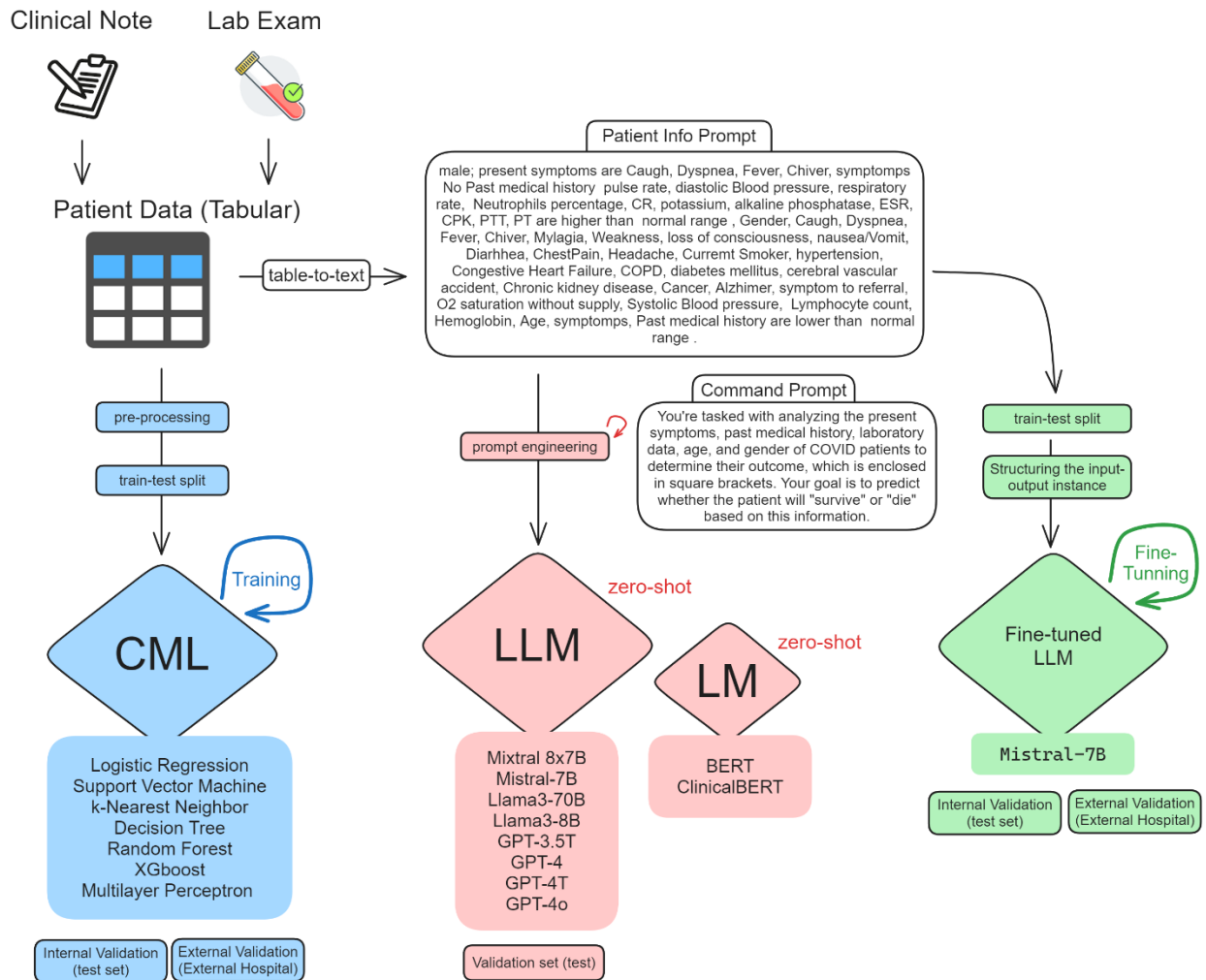
### 2.4.3 Oversampling

To address the issue of class imbalance in our dataset, we employed the synthetic minority oversampling technique (SMOTE), a widely used method in machine learning, particularly for medical diagnosis and prediction tasks (17). By applying SMOTE, we mitigated dataset imbalances, resulting in a more robust and reliable analysis for predicting mortality. SMOTE works by creating synthetic samples for the minority class instead of simply duplicating existing samples. It selects samples from the minority class and their nearest neighbors and then generates new synthetic samples by interpolating between these samples and their neighbors. This approach not only increases the number of samples in the minority class but also introduces new data points, improving dataset diversity. In our experiments, the SMOTE technique was applied to the training set ( $X_{train}$ ), increasing the number of samples from 6118 to 9760. Similarly, for the test set ( $X_{test}$ ), the sample size increased from 1530--2470, and for the external set ( $X_{ex}$ ), it increased from 1409--248.

### 2.4.4 Preparing Data for the LLM

To prepare the data for input into the LLM, we completed all the previous steps for feature selection and sampling, but normalization was not performed. As shown in Figure 1, we converted the dataset into text. We categorized the dataset features into symptoms, past medical history, age, sex, and laboratory data. For symptoms and medical history, we considered only positive data. For age, we added 'the patient's age is' before the age number. For sex, we used 'male' and 'female.' We used the normal range of laboratory data to classify the data into the normal range, higher than the normal range, and lower than the normal range. For example, if blood pressure and oxygen saturation were higher than the normal range, we used the sentence 'blood pressure and oxygen saturation are higher than the normal range.' We considered only laboratory data that were higher

or lower than the normal range. The exclusion of negative features in symptoms and past medical history, or the normal range in laboratory data, is due to limitations in LLM context windows. We then concatenated the dataset into a single paragraph for each patient, indicating their medical history.



**Figure 1. Study Design and Experimental Summary**

## 2.5 CML Predictive Performance

We employed five CML algorithms: logistic regression (LR), support vector machine (SVM), decision tree (DT), k-nearest neighbor (KNN), random forest (RF), multilayer perceptron neural network (MLP), and XGBoost. The hyperparameters were optimized via a grid search and cross-validation. The full details of training and hyperparameters are provided in **Supplementary Section 1**.

## 2.6 LLM Predictive Performance

We utilized open-source and proprietary LLMs to test their predictive power on clinical texts transformed from tabular data. First, we tested different prompts to determine the most efficient prompt to use, as well as the temperature (between 0.1 and 1). These prompts are listed in Supplementary Table S1. We then sent clinical text and commands, received the unstructured output, and extracted the selected outcome, which could be either "survive" or "die." **We used different sessions for each prediction, limiting the memory of the LLM to remembering previous generations.**

We tested open-source, open-weight models of Mistral-7b, Mixtral  $8 \times 7$  B, Llama3-8b, and Llama3-70b via the **Poe Chat Interface**. OpenAI models, including GPT-3.5T, GPT-4, GPT-4T, and GPT-4o, were utilized via the **OpenAI API**. We also tested the performance of two pretrained language models, BERT (18) and ClinicalBERT (19), which are fine-tuned versions of BERT on medical text. A list of all LLMs and times of use, as well as model parameters, is available in Supplementary Table S2.

### 2.6.1 Zero-Shot Classification



Zero-shot classification is an approach in prompt engineering in which the prompt is given to the model without any training. This approach is used in transfer learning, where a model used for different purposes is employed instead of fine-tuning a new model, thereby reducing the cost of training the new model. To perform zero-shot classification, we used eight different LLMs and two LMs. We provided each patient's history as input to predict whether the patient would die or survive and then stored the results.

### 2.6.2 Fine-tuning LLM

We fine-tuned one of the open-source LLMs, Mistral-7b-Instruct-v0.2, which is a GPT-like large language model with 7 billion parameters. It is trained on a mixture of publicly available and synthetic data and can be used for natural language processing (NLP) tasks. It is also a decoder-only model that is used for text-generation tasks. Fine-tuning an LLM is usually considered time-consuming and expensive; recently, several methods have been introduced to reduce costs. We implemented the QLoRA fine-tuning approach to optimize the LLM while minimizing computational resources (20).

The model was configured for 4-bit loading with double quantization, utilizing an "nf4" quantization type and torch.bfloat16 compute data type. A 16-layer model architecture with Lora attention and targeted projection modules was employed. We used the PEFT library to create a LoraConfig object with a dropout rate of 0.1 and task type 'CAUSAL\_LM'. The training pipeline, established via the transformer library, consisted of 4 epochs with a per-device batch size of 1 and gradient accumulation steps of 4. We utilized the "paged\_adamw\_32bit" optimizer with a learning rate of  $2e-4$  and a weight decay of 0.001. Mixed-precision training was conducted via fp16, with a maximum gradient norm of 0.3 and a warm-up ratio of 0.03. A cosine learning rate scheduler was employed, and training progress was logged every 25 steps and reported to TensorBoard. This

methodology, which combines QLoRA with the bits and bytes library, enables efficient enhancement of our language model while significantly reducing resource requirements, demonstrating superior performance across various instruction datasets and model scales. A more detailed description is provided in **Supplementary Section S2**.

## **2.6 CML and LLM Performance on Different Sample Sizes**

To investigate the influence of training sample sizes on model performance, we conducted a series of experiments using varying sample sizes: 20, 100, 200, 400, 1000, and 2476. Multiple models were trained using these sample sizes, and their performance was evaluated on the basis of the F1 score and accuracy metrics via an internal test set. The objective of this exploration was to gain valuable insights into the correlation between the volume of training data and the accuracy of predictive models.

## **2.7 Evaluation and analysis**

The accuracy of the outputs was assessed by comparing them against a ground truth that categorized outcomes as either mortality or survival. Outputs from the LLM were similarly classified. If an LLM initially produced an undefined result, the prompt was repeatedly presented up to five times to elicit a defined prediction; these instances are documented in Supplementary Table S1. We evaluated the models' performance via five critical metrics: specificity, recall, accuracy, precision, and F1 score. To optimize our models, we employed a grid search strategy with accuracy as the primary criterion. We further employed cross-validation to ensure the robustness and reliability of our model selection process. Additionally, the area under the receiver operating characteristic curve (AUC) was used to illustrate the predictive capacity of each model.

## **2.9 Explainability**

In our study, we employed SHAP (SHapley Additive exPlanations) values to examine both the total (global) and individual (granular) impacts of features on model predictions. We normalized the numerical data via a standard scaler and adopted a model-agnostic methodology. This approach involved employing XGBoost as the explainer model, which was chosen for its robust performance, as demonstrated in prior research and our own findings. SHAP values provide a clear, quantitative assessment of how each feature influences individual predictions, enhancing transparency in the model's decision-making process.

For our analysis, we used the test set for each model, generated SHAP values for every prediction, and computed the mean and standard deviation of the absolute SHAP scores. We then converted SHAP scores from a range of 0 to 1 into "global impact percentages" by dividing each feature's score by the total score of all features and multiplying by 100. We calculated the average impact percentages for both CMLs and LLMs by first averaging the SHAP scores and then determining the impact percentages. To compute the standard deviation of the impact percentages, we adjusted the average standard deviation of CML/LLM via a multiplication factor derived from the ratio of the impact score to the SHAP mean. The global impact percentage represents the proportion of each feature's impact on the predicted class across the entire dataset. A violin plot visually represents the variability of each input feature's effect on the output.

## Results

Our study initially included a dataset of 9,057 patients, with a mean age of  $58.40 \pm 19.81$  years and a male–female ratio of 1.19. The overall mortality rate in this group was 25.11% (N=1818). We utilized an internal validation test set and an external validation set comprising 2,470 and 2,248 participants, respectively, each with a mortality rate of 50%. Additionally, the validation set for

zero-shot classification included 590 patients randomly selected from the internal validation test set, with a mean age of  $63.85 \pm 18.37$  years, a male-to-female ratio of 355:255, and a mortality rate of 50% (mortality count = 295). The performance metrics of the CMLs and LLMs are detailed in Table 2.

### 3.1 Classic Machine Learning Predictive Performance

XGBoost and RF were the top-performing models in terms of accuracy, achieving scores of 86.28% and 86.52%, respectively, as detailed in Table 2. These models also excelled in precision, recall, specificity, and F1 scores, all surpassing 85%. The MLP also delivered an acceptable performance, with an accuracy of 75.87%. When the models were applied to the external validation set, a slight decline in the AUC of 2–5% was observed. Supplementary Figures S1 and S2 depict the performance of the CMLs on the internal validation test set and the external validation set, respectively. SVM, KNN, and DT showed consistent performance across both validation sets, confirming their reliability in generalizing to unseen data.

**Table 2. Model results on the internal validation test set and external validation dataset**

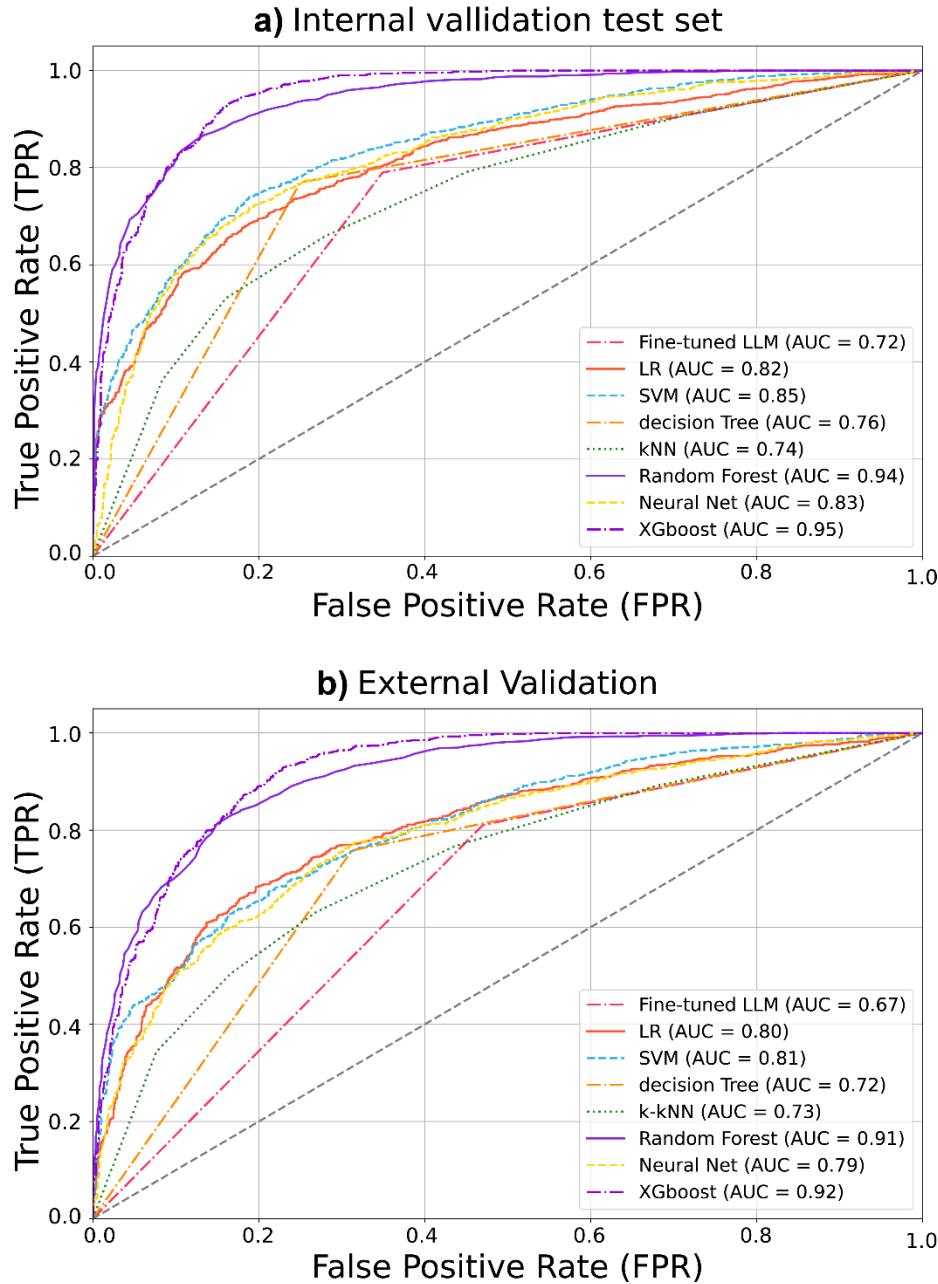
	Approach	Model	Accuracy	Precision	Recall	Specificity	F1	AUC
<b>Internal Validation</b>								
	CML	LR	0.74	0.74	0.75	0.74	0.74	0.82
	CML	SVM	0.76	0.73	0.81	0.71	0.77	0.85
	CML	DT	0.76	0.76	0.76	0.76	0.76	0.76
	CML	KNN	0.69	0.71	0.66	0.72	0.68	0.74
	CML	RF	0.86	0.84	<b>0.89</b>	0.83	<b>0.87</b>	0.94
	CML	Xgboost	<b>0.87</b>	<b>0.88</b>	0.84	<b>0.89</b>	0.86	<b>0.95</b>
	CML	MLP	0.76	0.73	0.81	0.70	0.77	0.83
	Fine-tuned LLM	Fine-tuned Mistral-7b	0.72	0.69	0.79	0.65	0.74	0.72
<b>External Validation</b>								
	CML	LR	0.73	0.73	0.74	0.72	0.73	0.80
	CML	SVM	0.74	0.73	0.75	0.72	0.74	0.81
	CML	DT	0.73	0.72	0.73	0.72	0.73	0.72
	CML	KNN	0.68	0.71	0.62	0.75	0.66	0.73
	CML	RF	<b>0.82</b>	0.79	<b>0.86</b>	0.77	<b>0.83</b>	0.91
	CML	Xgboost	<b>0.82</b>	<b>0.85</b>	0.78	<b>0.86</b>	0.82	<b>0.92</b>
	CML	MLP	0.71	0.69	0.76	0.67	0.72	0.79
	Fine-tuned LLM	Fine-tuned Mistral-7b	0.69	0.69	0.68	0.69	0.69	0.67
<b>ZSC Validation*</b>								
	LLM	Mistral-7b	0.51	0.80	0.01	<b>1.00</b>	0.03	0.51
	LLM	Mixtral-8x7b	0.52	0.94	0.05	<b>1.00</b>	0.10	0.52
	LLM	Llama-3-8b	0.54	<b>1.00</b>	0.08	<b>1.00</b>	0.15	0.54
	LLM	Llama-3-70b	0.54	0.89	0.08	0.99	0.15	0.54
	LLM	gpt-3.5-turbo	0.50	0.49	0.17	0.83	0.25	0.50
	LLM	gpt-4	<b>0.62</b>	0.84	<b>0.28</b>	0.95	<b>0.43</b>	<b>0.62</b>
	LLM	gpt-4-turbo	0.57	0.82	0.19	0.96	0.31	0.57
	LLM	gpt-4o	0.50	1.00	0.00	1.00	0.01	0.50
	LM	BERT	0.50	1.00	0.00	1.00	0.01	0.50
	LM	ClinicalBERT	0.50	1.00	0.00	1.00	0.01	0.50

Footnote: \* ZSC validation dataset was created using a random sample of the internal validation dataset. Abbreviations: AU-ROC, area under the receiver operating characteristic curve; ZSC validation: zero-shot classification; CML, classical machine learning; LLM, large language model; LR, logistic regression; SVM, support vector machine; DT, decision tree; KNN, K-nearest neighbors; RF, random forest; XGBoost, eXtreme gradient boosting; MLP, multilayer perceptron

### 3.2 LLM: Zero-shot classification and fine-tuned Mistral-7b

Table 2. The zero-shot classification results showed variability among the models, with GPT-4 outperforming the other models by achieving an accuracy of 0.62 and an F1 score of 0.43 and recording the highest recall at 0.28 among the LLMs. Generally, LLMs exhibited low recall rates, predominantly classifying predictions as "mortality." The open-source models, including Llama-3-70B, Llama-3-8B, Mistral-7b-Instruct, and Mistral-8x7b-Instruct-v0.1, had F1 scores ranging from 0.03 (Mistral 7b) to 0.15 (L Llama3-8B and Llama3--70b). Notably, the gpt-°model showed limited effectiveness, with an F1 score of 0.01, indicating a challenge in distinguishing between true positives and true negatives. The pretrained language models – BERT and ClinicalBERT – also labeled all outcomes as dies, failing to provide predictive power.

Fine-tuning Mistral-7b significantly improved its performance, increasing the F1 score from 0.03 to 0.74 in the internal test set and to 0.69 in the external test set. This fine-tuned version also demonstrated a high recall rate of 78.98%, a substantial increase from 1% in zero-shot classification, showing its ability to accurately identify a greater proportion of actual survival instances. This consistency between internal and external validations highlights the generalizability of the fine-tuned Mistral-7b in mortality prediction. The detailed results are presented in Supplementary Figure S4. Also, Figure 2 sets out the ROC curves and AUC scores for internal and external validation of COVID-19 mortality prediction models.



**Figure 2. ROC curves and AUC scores for internal and external validation of COVID-19 mortality prediction models, including the fine-tuned LLM (Mistral 7b) and 7 classical machine learning algorithms.**

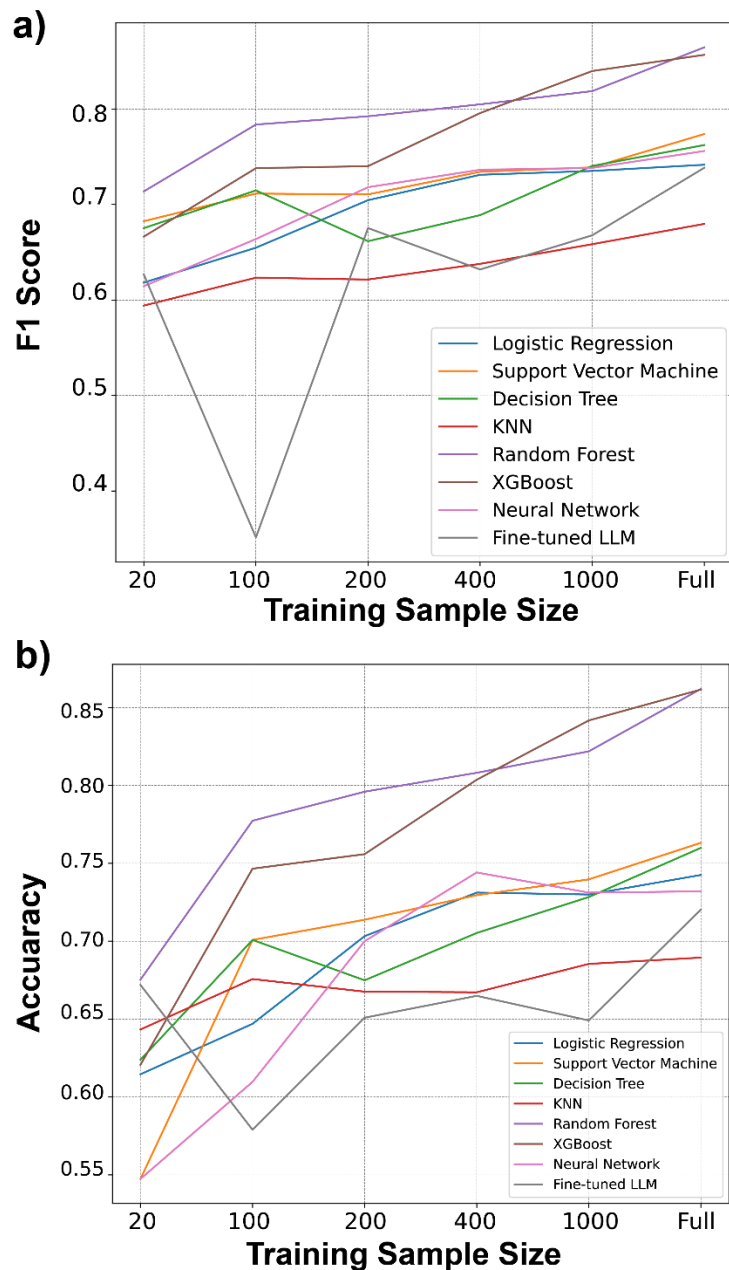
*Abbreviations: LLM, Large Language Model; LR, Logistic Regression; SVM, Support Vector Machine; kNN, k-Nearest Neighbors; TPR, True Positive Rate; FPR, False Positive Rate; AUC, Area Under the Curve.*

### 3.3 Comparing models on different training sample sizes

To evaluate the impact of training sample size on model efficacy, experiments were conducted across various sample sizes. Figure 3 shows that the performance of all CMLs increased as the size of the training set increased. XGBoost demonstrated the strongest performance across all categories: small (100 samples), medium (400--1000 samples), and full training set sizes (2476 samples). Notably, the MLP neural network and SVM exhibited the most significant performance improvements, with accuracies increasing from 55% with 20 training samples to 73% and 77%, respectively.

In contrast, while the zero-shot performance of GPT-4 reached an F1 score of 0.43, CMLs still surpassed both zero-shot classification and fine-tuned LLMs in predicting COVID-19 mortality. During the fine-tuning of Mistral-7b, notable performance degradation occurred in scenarios with small training sizes, leading to a loss of broader model understanding, an effect termed "negative transfer."





**Figure 3. Performance comparison of classical machine learning models and fine-tuned large language models in COVID-19 mortality prediction using different sample sizes.**

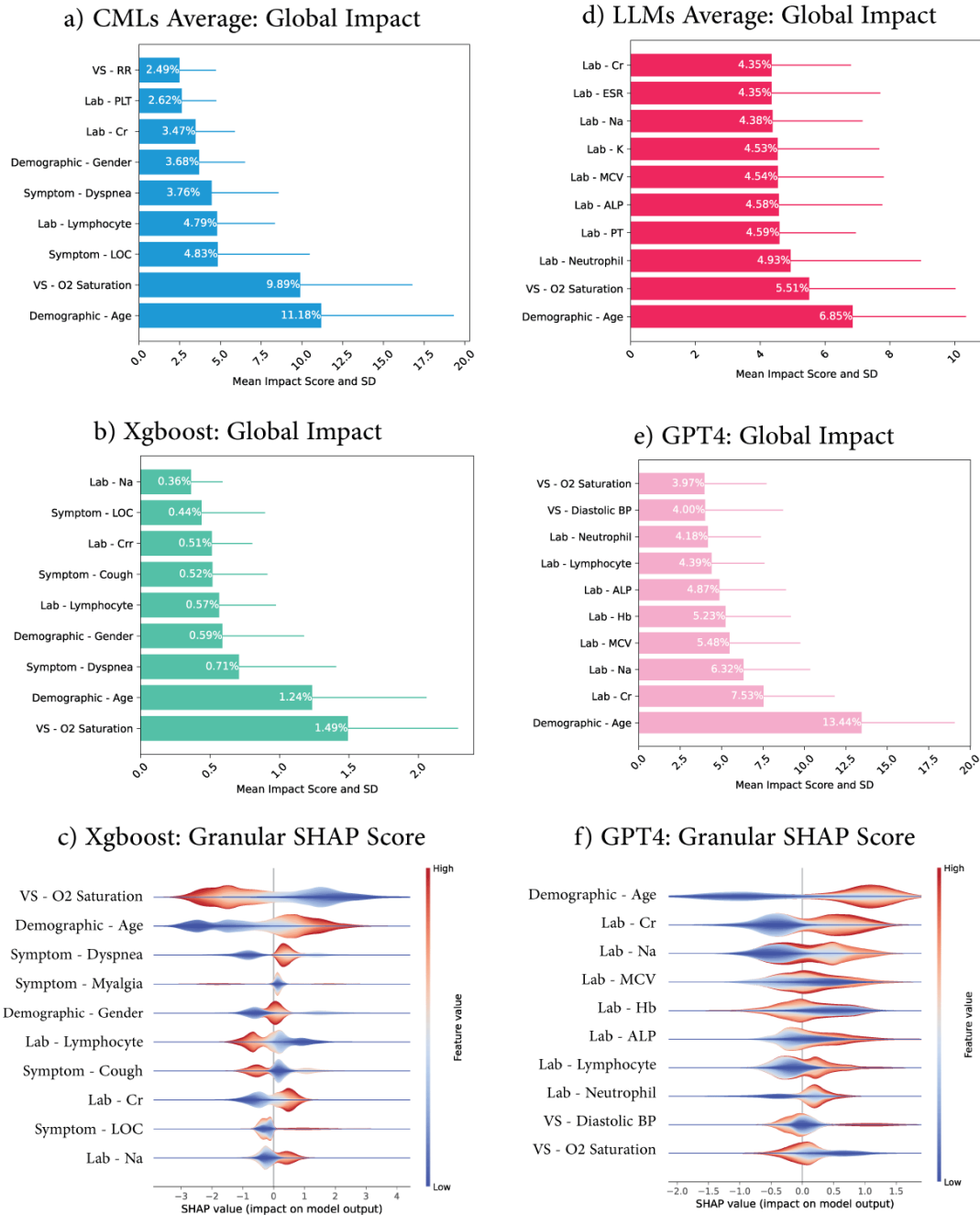
*Caption: This figure presents the performance of seven conventional machine learning (CML) models and a fine-tuned large language model (LLM) in predicting COVID-19 mortality. The performance metrics evaluated are the F1 score and accuracy across different training sample sizes (20, 100, 200, 400, 1000, and the full 2047 training samples). The CML models include logistic regression, support vector machines, decision trees, k-nearest neighbors (KNNs), random forests, XGBoost, and neural networks. Compared with the conventional models, XGBoost consistently achieves superior performance in terms of both the F1 score and accuracy, especially as the training sample size increases.*

### 3.4 Explainability: Impact of Features on Prediction

As shown in Supplementary Figure S5, while the global impact of features among CMLs exhibits similar patterns, with many of the top 10 impactful features being consistent, the granular impact differs significantly. For example, in the context of O2 saturation levels in patients, XGBoost, RF, DT, and MCP consider both high (increasing mortality risk) and low (increasing survival chance) levels to be significant, whereas KNN and LR focus only on low saturation levels. According to Figure 4.a, the most influential features are age (11.18%) and O2 saturation (9.89%), followed by LOC (4.83%), lymphocyte count (4.79%), dyspnea (3.76%), and sex (3.68%).

Conversely, the influence of features in LLMs, particularly in lower-performing models such as Mistralb-7b and GPT4o, appears less coherent, as illustrated in Supplementary Figure S6. This inconsistency contributes to noise in the average feature impact among LLMs (Figure 4.d). Nonetheless, age (6.58%) and O2 saturation (5.51%) remained the most significant features, with a series of laboratory tests, including neutrophil count, PT, ALP, MCV, K, Na, ESR, and Cr, revealing impacts in the 4%--5% range.

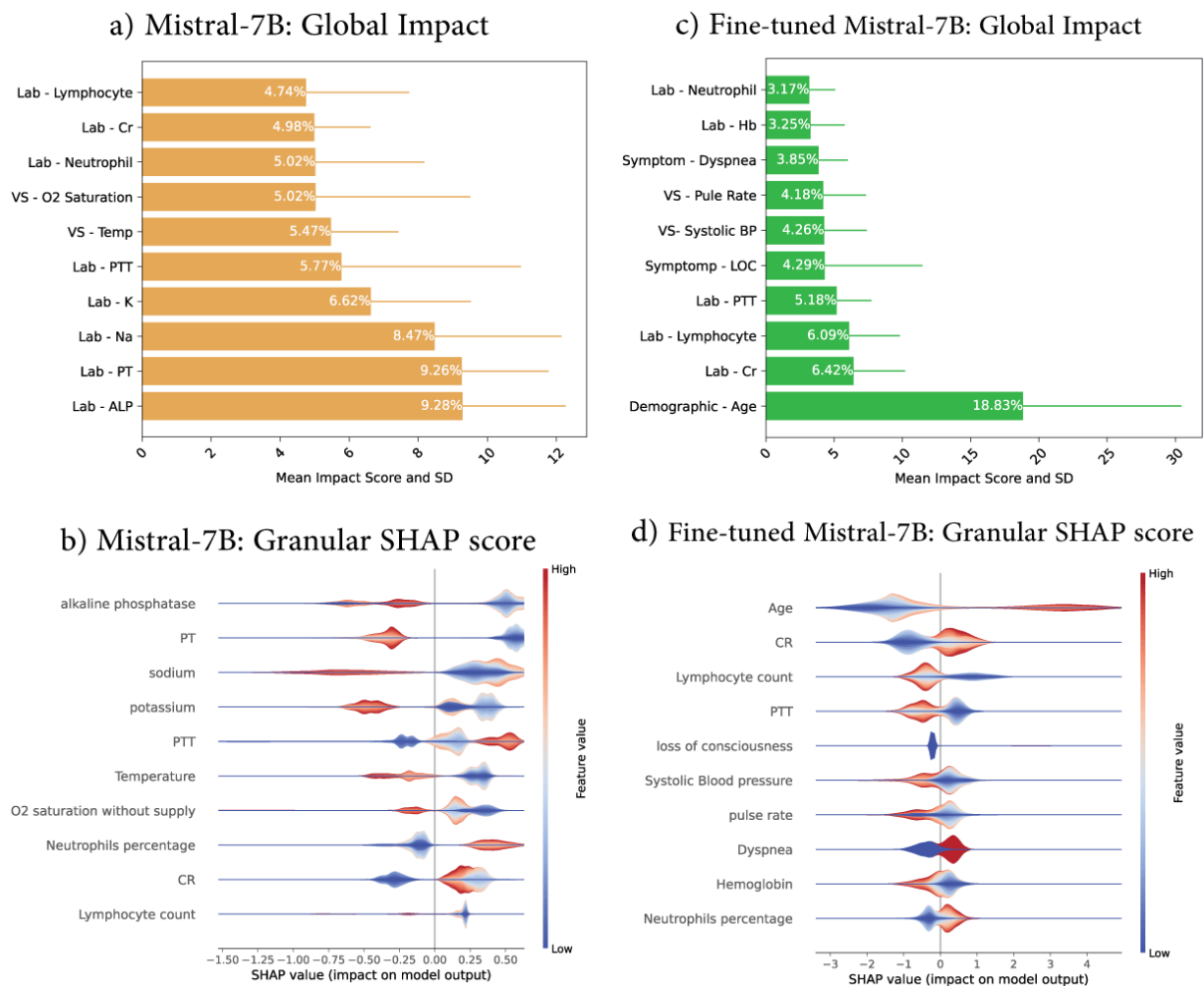
When comparing the top performers among CMLs and LLMs—XGBoost and GPT4—the patterns of global (Fig. 4.b and Fig. 4.e) and granular (Fig. 4.c and Fig. 4.f) impacts diverge, with XGBoost displaying more specific impacts and GPT4 showing broader ranges of impact.



**Figure 4. SHAP analysis of feature importance in COVID-19 mortality prediction models, including global feature impact in classical machine learning (CML) average (a), CML best performing XGBoost (b), large language models (LLM) average (e), LLM best performing GPT-4 (d), and the granular impact of XGBoost and GPT-4**

*Abbreviations: CML, classical machine learning; LLM, large language model; VS, vital sign; RR, respiratory rate; LOC, level of consciousness; ALP, alkaline phosphatase; PT, prothrombin time; MCV, mean corpuscular volume; Hb, hemoglobin; Cr, creatinine; Na, sodium; BP, blood pressure.*

Figure 5 illustrates how fine-tuning Mistral-7b altered the impact of features at both the global and granular levels. This refinement in prediction logic aligned the top 10 most important features more closely with those of CMLs, resulting in more equitable impact percentages among features and enhanced granularity.



**Figure 5. SHAP Analysis global and granular feature impact for Mistral-7b (a and b) and Fine-tuned Mistral-7b Model (c and d) in COVID-19 Mortality Prediction**

## Discussion

Our study reveals a notable performance gap between CML models and LLMs in predicting patient mortality via tabular data. RF and XGBoost emerged as the top CML performers, achieving over 80% accuracy and an F1 score of 0.86. In contrast, the best-performing LLM, GPT-4, achieved 62% accuracy and an F1 score of 0.43 in zero-shot classification. This disparity highlights the challenges LLMs face when dealing with purely tabular data. Notably, increasing our sample size from 5,000 patients in our previous study to 9,000 patients in this study significantly improved the performance of CML models. The AUC of RF improved from 0.82 to 0.94, underscoring the importance of large and diverse datasets in realizing the full potential of CMLs in medical tasks.

LLM performance heavily relies on the knowledge embedded within model weights, the complexity of input data, and the table-to-text transformation technique. Our approach, which uses a simple prompt and transformation to resonate with current clinical use, achieved results comparable to those of similar studies, with F1 scores of 0.50–0.60 across different medical tasks using LLMs such as the GPT-4 or GPT-3.5 (10,11,21). However, in line with many previous studies, we found that CMLs can outperform this zero-shot performance with even fewer than 100 training samples (10,21).

Given the performance gap between CMLs and LLMs, researchers have explored two main approaches for improving LLM performance: pipeline improvements and fine-tuning. Previous studies have shown that LLMs can close the gap in CML performance via pipeline improvements such as prompt engineering techniques (XAI4LLM), few-shot approaches (XAI4LLM, EHR-CoAgent, TabLLM), multiple runs of LLM to double-check results (EHR-CoAgent), the addition of a tree-based explainer alongside the LLM (XAILLM), or novel LLM-based text-to-table

transformation (Medi-TAB, TabLLM) (22–24). However, many of their evaluated tasks may not resonate with real-world use, as they have low-dimensional datasets (8-15 features) that do not reflect real-world complex medical data and limited sample sizes (<500 instances in rare classes) that restrict CMLs from reaching their maximum performance.

The alternative approach, fine-tuning or in-context learning, aims to modify the model weights to teach the model a new task, which has been evaluated on name entity recognition and text extraction (25,26). We validated this approach in our high-dimensional task, where fine-tuning Mistral increased the F1 score from 0.03 to 0.69, even with a resource-efficient QLoRa method. Our SHAP analysis provides initial evidence of an improved rationale after fine-tuning, as the top 10 features more closely align with XGBoost and clinician decision-making.

Despite these advancements, LLMs still face significant limitations that affect their applicability in medical settings. Their vulnerability to hallucination raises concerns about producing harmful information (27), whereas computational constraints impose token limits that can truncate responses and diminish interaction quality. (27,28). Data privacy is another crucial concern, particularly in medical contexts, as many powerful LLMs are proprietary or require cloud-based computations, increasing the risk of data leaks (29). Moreover, the cost of using LLM APIs for large clinical databases can disproportionately impact low- and middle-income communities (30). While open-source models present a more affordable alternative, they may not match the capabilities of proprietary models.

In light of these challenges, alternative approaches have emerged, including the use of small pretrained language models and rule-based systems. These offer resource-efficient alternatives to large LLMs. Previous studies have shown that rule-based and gradient boosting algorithms can achieve strong overall performance in specific tasks, such as extracting physical rehabilitation

exercise information from clinical notes (31,32). Additionally, fine-tuning pretrained BERT-like models has yielded promising results in some medical applications. However, our brief experiment with the zero-shot performance of pretrained models (BERT and ClinicalBERT) revealed their limitations, suggesting that further research is needed to optimize these approaches for complex medical tasks.

It is important to acknowledge several limitations of our study. While our fine-tuning method was resource efficient, it may not have been the optimal approach for achieving the highest performance. Our fine-tuned model was a small LLM with the lowest performance among our eight tested LLMs, indicating that fine-tuning larger and more accurate models could yield better results. Furthermore, our table-to-text transformation and prompts were designed to resonate with a medical user context, but more robust approaches (e.g., few-shot learning, advanced prompt engineering, and sophisticated transformation techniques) may achieve higher accuracies, especially in zero-shot classification (11,21). Although our sample size was substantial, the retrospective nature of our investigation necessitates prospective validation to confirm the generalizability of these findings. Moreover, since all participating hospitals were located within limited resources, potential disparities in healthcare access and quality might have influenced mortality rates compared with those in developed nations with robust medical systems.

Our findings highlight several critical areas for future research in the application of LLMs to medical data analysis. We propose the following research questions to advance the field:

- Does the LLM explanation of the prediction (death or survival) in human language align with the feature importance analysis? Can LLMs accurately explain their rationale?
- What would be the performance of fine-tuning pretrained models and large LLMs compared to small LLMs?

- Could we create a model to distinguish correct answers from incorrect answers via LLM output? How can we measure the certainty of the given answer?

## Conclusion

The efficacy of LLMs versus CML approaches in medical tasks appears to be contingent upon data dimensionality and data availability. In low-dimensional scenarios with limited samples, LLM-based methodologies may offer superior performance; however, as dimensionality increases and diverse sample sizes become available, CML techniques tend to outperform the zero-shot capabilities of LLMs. Notably, fine-tuning LLMs can substantially enhance their pattern recognition and logical processing, potentially achieving performance levels comparable to those of CMLs. The potential of LLMs to process both structured and unstructured data may outweigh marginally lower performance metrics than CMLs do. Ultimately, the choice between LLMs and CMLs should be guided by careful consideration of task complexity, data characteristics, and clinical context demands, with further research warranted to elucidate the precise conditions under which each methodology excels.

## Conflict of interest declaration

All the authors declare that they have no financial or nonfinancial conflicts of interest related to this work.

## Acknowledgments

We used ChatGPT 3.5 with the following prompt: “Is this paragraph grammatically correct and can you make it sound scientific? ” and grammar to improve the English language of the article.



Two authors, SAASN and MG, reviewed the suggestions and accepted relevant changes. All the authors are responsible for the validity of the final draft.

## **Funding**

No funding was available related to this work.

## **Authors' contributions**

MoGE: Conceptualization, Methodology, Programming, Investigation, Writing Original Draft; MaGE: Investigation, Methodology; ASN: Investigation, Methodology; HT: Writing Original Draft, Methodology; ZA: Investigation, Programming; AMK: Investigation; MS: Investigation; ZT: Investigation; FS: Investigation; MHB: Investigation; AF: Investigation; MT: Investigation; NG: Investigation; FH: Investigation; HA: Investigation; AA: Investigation; FA: Investigation; AS: Investigation; NA: Investigation; MAK: Investigation, Project Administration; HS: Investigation; AM: Investigation; SHZ: Investigation; OY: Investigation; RE: Investigation; MM: Investigation; DSN: Investigation; ALS: Investigation; HM: Investigation; SS: Investigation; ARS: Investigation; NG: Investigation; ET: Investigation, Validation; HH: Investigation; JSS: Reviewing and Editing the Manuscript; TS: Reviewing and Editing the Manuscript; AKS: Reviewing and Editing the Manuscript; ALS: Methodology, Validation, Reviewing and Editing the Manuscript; GN: Reviewing and Editing the Manuscript; IAD: Data Acquisition, Reviewing and Ed

## **Availability of Data**

The code and information for generating the output are available at <https://github.com/mohammad-gh009/Large-Language-Models-vs-Classical-Machine-Learning> and [https://github.com/Sdamirsa/Tehran\\_COVID\\_Cohort](https://github.com/Sdamirsa/Tehran_COVID_Cohort). In addition, the streamlet app with the

user interface and instructions for using it are available on the GitHub repo. Users can extract information via their OpenAI API and descriptions of the variables they want to extract. The anonymized MR enterography reports are available upon reasonable request containing the aim for use and ethical approval number from the corresponding author, ARR.

## **Abbreviations**

LLM: Large Language Model

CML: classical machine learning model

LR: Logistic regression

SVM: Support vector machine

DT: decision tree

KNN: k-nearest neighbor

RF: random forest

XGBoost: Extreme Gradient Boosting

MLP: multilayer perceptron

Zero-shot classification: ZSC

LASSO: least absolute shrinkage and selection operator

SMOTE: synthetic minority oversampling technique

QLoRA: quantized low-rank adaptation

MICE: Multiple Imputation by Chained Equations

ReLU: rectified linear unit

KBit: knowledge bit

CRP: C-reactive protein

LDH: lactate dehydrogenase

NLP: natural language processing

chain-of-thought (CoT)

## References

1. Karabacak M, Margetis K. Embracing Large Language Models for Medical Applications: Opportunities and Challenges. Cureus. 2023;
2. Dathathri S, Madotto A, Lan J, Hung J, Frank E, Molino P, et al. PLUG AND PLAY LANGUAGE MODELS: A SIMPLE APPROACH TO CONTROLLED TEXT GENERATION. In: 8th International Conference on Learning Representations, ICLR 2020. 2020.
3. Han JM, Babuschkin I, Edwards H, Neelakantan A, Xu T, Polu S, et al. Unsupervised neural machine translation with generative language models only. arXiv preprint arXiv:211005448. 2021;
4. Petroni F, Rocktäschel T, Lewis P, Bakhtin A, Wu Y, Miller AH, et al. Language models as knowledge bases? arXiv preprint arXiv:190901066. 2019;
5. Vaid A, Lampert J, Lee J, Sawant A, Apakama D, Sakhuja A, et al. Generative Large Language Models are autonomous practitioners of evidence-based medicine. arXiv preprint arXiv:240102851. 2024;

6. Zhang D, Yin C, Zeng J, Yuan X, Zhang P. Combining structured and unstructured data for predictive models: a deep learning approach. *BMC Med Inform Decis Mak*. 2020;20:1–11.
7. Sedlakova J, Daniore P, Horn Wintsch A, Wolf M, Stanikic M, Haag C, et al. Challenges and best practices for digital unstructured data enrichment in health research: A systematic narrative review. *PLOS Digital Health* [Internet]. 2023 Oct 11;2(10):e0000347-. Available from: <https://doi.org/10.1371/journal.pdig.0000347>
8. Zhou H, Gu B, Zou X, Li Y, Chen SS, Zhou P, et al. A survey of large language models in medicine: Progress, application, and challenge. *arXiv preprint arXiv:231105112*. 2023;
9. Wornow M, Xu Y, Thapa R, Patel B, Steinberg E, Fleming S, et al. The shaky foundations of large language models and foundation models for electronic health records. *NPJ Digit Med* [Internet]. 2023;6(1):135. Available from: <https://doi.org/10.1038/s41746-023-00879-8>
10. Hegselmann S, Buendia A, Lang H, Agrawal M, Jiang X, Sontag D. Tabllm: Few-shot classification of tabular data with large language models. In: *International Conference on Artificial Intelligence and Statistics*. PMLR; 2023. p. 5549–81.
11. Wang Z, Gao C, Xiao C, Sun J. MediTab: Scaling Medical Tabular Data Predictors via Data Consolidation, Enrichment, and Refinement. *arXiv preprint arXiv:230512081*. 2023;
12. Cui H, Shen Z, Zhang J, Shao H, Qin L, Ho JC, et al. LLMs-based Few-Shot Disease Predictions using EHR: A Novel Approach Combining Predictive Agent Reasoning and Critical Agent Instruction. *arXiv preprint arXiv:240315464*. 2024;

13. Nazary F, Deldjoo Y, Di Noia T, di Sciascio E. XAI4LLM. Let Machine Learning Models and LLMs Collaborate for Enhanced In-Context Learning in Healthcare. arXiv preprint arXiv:240506270. 2024;
14. Patel D, Timsina P, Gorenstein L, Glicksberg BS, Raut G, Cheetirala SN, et al. Comparative Analysis of a Large Language Model and Machine Learning Method for Prediction of Hospitalization from Nurse Triage Notes: Implications for Machine Learning-based Resource Management. medRxiv [Internet]. 2023 Jan 1;2023.08.07.23293699. Available from: <http://medrxiv.org/content/early/2023/08/10/2023.08.07.23293699.abstract>
15. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine learning in Python. the Journal of machine Learning research. 2011;12:2825–30.
16. Simpson L, Combettes PL, Müller CL. c-lasso--a Python package for constrained sparse and robust regression and classification. arXiv preprint arXiv:201100898. 2020;
17. Liu XY, Wu J, Zhou ZH. Exploratory Undersampling for Class-Imbalance Learning. IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics). 2009;39(2):539–50.
18. Devlin J, Chang MW, Lee K, Toutanova K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In: North American Chapter of the Association for Computational Linguistics [Internet]. 2019. Available from: <https://api.semanticscholar.org/CorpusID:52967399>
19. Wang G, Liu X, Ying Z, Yang G, Chen Z, Liu Z, et al. Optimized glycemic control of type 2 diabetes with reinforcement learning: a proof-of-concept trial. Nat Med [Internet]. 2023;29(10):2633–42. Available from: <https://doi.org/10.1038/s41591-023-02552-9>

20. Dettmers T, Pagnoni A, Holtzman A, Zettlemoyer L. Qlora: Efficient finetuning of quantized llms. *Adv Neural Inf Process Syst*. 2024;36.
21. Nazary F, Deldjoo Y, Di Noia T, di Sciascio E. XAI4LLM. Let Machine Learning Models and LLMs Collaborate for Enhanced In-Context Learning in Healthcare. *arXiv preprint arXiv:240506270*. 2024;
22. Kojima T, Gu SS, Reid M, Matsuo Y, Iwasawa Y. Large language models are zero-shot reasoners. *Adv Neural Inf Process Syst*. 2022;35:22199–213.
23. Brown T, Mann B, Ryder N, Subbiah M, Kaplan JD, Dhariwal P, et al. Language models are few-shot learners. *Adv Neural Inf Process Syst*. 2020;33:1877–901.
24. Han Z, Gao C, Liu J, Zhang SQ. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:240314608*. 2024;
25. Kanzawa J, Yasaka K, Fujita N, Fujiwara S, Abe O. Automated classification of brain MRI reports using fine-tuned large language models. *Neuroradiology* [Internet]. 2024; Available from: <https://doi.org/10.1007/s00234-024-03427-7>
26. Akbasli IT, Birbilen AZ, Teksam O. Human-Like Named Entity Recognition with Large Language Models in Unstructured Text-based Electronic Healthcare Records: An Evaluation Study. 2024;
27. O'Neill M, Connor M. Amplifying Limitations, Harms and Risks of Large Language Models. *arXiv preprint arXiv:230704821*. 2023;
28. Savage T, Wang J, Gallo R, Boukil A, Patel V, Ahmad Safavi-Naini SA, et al. Large Language Model Uncertainty Measurement and Calibration for Medical Diagnosis and

- Treatment. medRxiv [Internet]. 2024 Jan 1;2024.06.06.24308399. Available from: <http://medrxiv.org/content/early/2024/06/10/2024.06.06.24308399.abstract>
29. Wirth FN, Meurers T, Johns M, Prasser F. Privacy-preserving data sharing infrastructures for medical research: systematization and comparison. BMC Med Inform Decis Mak [Internet]. 2021;21(1):242. Available from: <https://doi.org/10.1186/s12911-021-01602-x>
  30. Gangavarapu A. Introducing L2M3, A Multilingual Medical Large Language Model to Advance Health Equity in Low-Resource Regions. arXiv preprint arXiv:240408705. 2024;
  31. Sivarajkumar S, Gao F, Denny P, Aldhahwani B, Visweswaran S, Bove A, et al. Mining Clinical Notes for Physical Rehabilitation Exercise Information: Natural Language Processing Algorithm Development and Validation Study. JMIR Med Inform [Internet]. 2024;12:e52289. Available from: <http://dx.doi.org/10.2196/52289>
  32. Chen S, Li Y, Lu S, Van H, Aerts HJWL, Savova GK, et al. Evaluating the ChatGPT family of models for biomedical reasoning and classification. J Am Med Inform Assoc [Internet]. 2024;31(4):940–8. Available from: <http://dx.doi.org/10.1093/jamia/ocad256>

This is a supplementary file to "**Large Language Models versus Classical Machine Learning: Performance on COVID-19 Mortality Prediction Using High-Dimensional Tabular Data**" by Mohammadreza Ghaffarzadeh-Esfahani, Mahdi Ghaffarzadeh-Esfahani, Arian Salahi-Niri, Hossein Toreyhi, Zahra Atf, Amirali Mohsenzadeh-Kermani, Mahshad Sarikhani, Zohreh Tajabadi, Fatemeh Shojaeian, Mohammad Hassan Bagheri, Aydin Feyzi, Mohammadamin Tarighatpayma, Narges Gazmeh, Fateme Heydari, Hossein Afshar, Amirreza Allahgholipour, Farid Alimardani, Ameneh Salehi, Naghmeh Asadimanesh, Mohammad Amin Khalafi, Hadis Shabanipour, Ali Moradi, Sajjad Hossein Zadeh, Omid Yazdani, Romina Esbati, Moozhan Maleki, Danial Samiei Nasr, Amirali Soheili, Hossein Majlesi, Saba Shahsavani, Alireza Soheilipour, Nooshin Goudarzi, Erfan Taherifard, Hamidreza Hatamabadi, Jamil S. Samaan, Thomas Savage, Ankit Sakhuja, Ali Soroush, Girish Nadkarni, Ilad Alavi Darazam, Mohamad Amin Pourhoseingholi, Seyed Amir Ahmad Safavi-Naini

\* Correspondence to: Seyed Amir Ahmad Safavi-Naini ([sdamirsa@ymail.com](mailto:sdamirsa@ymail.com)), Mohamad Amin Pourhoseingholi ([aminphg@gmail.com](mailto:aminphg@gmail.com)), and Ilad Alavi Darazam ([ilad13@yahoo.com](mailto:ilad13@yahoo.com))

## List of Supplementary Materials

- **Supplementary Section 1.** Details and hyperparameters of seven machine learning models
- **Supplementary Section 2.** Details of fine-tuning the LLM with the QLoRA
- **Supplementary Table S1.** Description of the raw prompts and improved prompts used during the experiments.
- **Supplementary Table S2.** Environments, model settings and model parameters for each LLM
- **Supplementary Figure S1.** Confusion matrices for **classical machine learning** models (CMLs) trained on the internal validation test set, including logistic regression (a), support vector machine (b), K-nearest neighbor (c), decision tree (d), random forest (e), XGBoost (f), and multilayer perceptron (MLP)
- **Supplementary Figure S2.** Confusion matrices for seven classical machine learning models tested on the external validation set: logistic regression (a), support vector machine (b), K-nearest neighbor (c), decision tree (d), random forest (e), XGBoost (f), and multilayer perceptron (MLP) (g)
- **Supplementary Figure S3.** Confusion matrices for LLM zero-shot tasks, including Mixtral-8x7B-Instruct-v0.1 (a), Llama-3-70B (b), Mistral-7B-Instruct (c), Llama-3-8B (d), gpt-4° (e), gpt-3.5-turbo (f), gpt-4-turbo (g), and gpt-4 (h)
- **Supplementary Figure S4.** Confusion matrices for fine-tuned mistral-7b-instruct-v0.2 on the internal validation test set (a) and external validation set (b).
- **Supplementary Figure S5.** Features with granular (left column) and global (right column) impacts on predictions among CMLs
- **Supplementary Figure S6.** Features granular (left column) and global (right column) impacts on predictions among LLMs



## **Supplementary Section 1. Details and hyperparameters of seven machine learning models**

### **Supplementary Section 1.1 Logistic Regression**

In the logistic regression model employed for this study, L2 regularization (penalty='l2') was applied with a convergence tolerance set to 0.0001, a regularization strength of 1.0, and the use of the limited-memory Broyden–Fletcher–Goldfarb–Shanno algorithm for optimization. The model was configured to fit the intercept term, with the intercept scaling set to 1 and no specific class weights assigned. The logistic regression utilized a maximum of 100 iterations for convergence, automatic detection of the multiclass scenario, 5-fold cross-validation, and a verbosity level set to 0 for minimal output during training. The random state of 42 was set for reproducibility.

### **Supplementary Section 1.2 Support Vector Machine**

An SVC model was employed in this study with a regularization parameter C set to 1.0, utilizing the radial basis function kernel ('rbf') with a default degree of 3. The gamma parameter was set to 'scale', the coefficient of the kernel function was set to 0.0, and shrinking during optimization was enabled. The model was configured without probability estimates, and the tolerance for convergence was set to 0.001. A cache size of 200 was allocated for optimization (cache\_size=200), with no specific class weights assigned (class\_weight=None). The model had no specified maximum iteration limit (max\_iter=-1) and employed the 'ovr' (one-vs-rest) strategy for decision function shape. Break ties were not considered in the decision function (break\_ties=False), 5-fold cross-validation was used, and a random state of 42 was set for reproducibility.

### **Supplementary Section 1.3 Decision Tree**

A DT classifier was employed in this study with the Gini impurity criterion ('gini'), utilizing the 'best' strategy for splitting nodes. The tree had no specified maximum depth (max\_depth=None), and a minimum of 2 samples were required to split an internal node (min\_samples\_split=2). The minimum number of samples required to be in a leaf node was set to 1 (min\_samples\_leaf=1), with no specified minimum weight fraction for a leaf node (min\_weight\_fraction\_leaf=0.0). The maximum number of features considered for splitting was not restricted (max\_features=None), 5-

fold cross-validation was used, and a random state of 42 was set for reproducibility. There were no limitations on the maximum number of leaf nodes (`max_leaf_nodes=None`), and the minimum impurity decrease for a split was set to 0.0 (`min_impurity_decrease=0.0`). No specific class weights were assigned (`class_weight=None`), and the cost complexity pruning alpha (`ccp_alpha`) parameter was set to 0.0. Additionally, monotonic constraints on features were not specified (`monotonic_cst=None`).

#### **Supplementary Section 1.4 K-nearest neighbor**

A KNN classifier was utilized in this study with the number of neighbors set to 5 (`n_neighbors=5`), employing uniform weights for neighbor contributions (`weights='uniform'`). The algorithm used for computing nearest neighbors was automatically determined (`algorithm='auto'`), with a leaf size of 30 (`leaf_size=30`). The Minkowski distance metric with a power parameter of 2 (`p=2`) was employed (`metric='minkowski'`). Additional metric parameters were not specified (`metric_params=None`), 5-fold cross-validation was used, and parallel processing was not utilized (`n_jobs=None`).

#### **Supplementary Section 1.5 Random forest**

An RF classifier was employed in this study with 100 estimators (`n_estimators=100`), utilizing the Gini impurity criterion (`criterion='gini'`). The DT in the RF had no specified maximum depth (`max_depth=None`), and a minimum of 2 samples was required to split an internal node (`min_samples_split=2`). The minimum number of samples required to be in a leaf node was set to 1 (`min_samples_leaf=1`), with no specified minimum weight fraction for a leaf node (`min_weight_fraction_leaf=0.0`). The square root of the total number of features was considered for splitting at each node (`max_features='sqrt'`). The maximum number of leaf nodes was unrestricted (`max_leaf_nodes=None`), and the minimum impurity decrease for a split was set to 0.0 (`min_impurity_decrease=0.0`). Bootstrap sampling was enabled (`bootstrap=True`), out-of-bag scoring was not utilized (`oob_score=False`), 5-fold cross-validation was used, and parallel processing was not employed (`n_jobs=None`). A random state of 42 was set for reproducibility during the model-building process (`random_state=42`). The model verbosity level was set to 0 (`verbose=0`), and warm starting was disabled (`warm_start=False`). No specific class weights were assigned (`class_weight=None`), and the cost complexity pruning alpha (`ccp_alpha`) parameter was set to 0.0. The maximum number of samples for bootstrapping was not specified

(max\_samples=None), and monotonic constraints on features were not specified (monotonic\_cst=None).

## **Supplementary Section 1.6 XGBoost**

The XGBoost classifier model was configured with default hyperparameters, including a base score of 0.5, utilizing a gradient boosting tree-based approach ('booster='gbtree'). The maximum depth of each tree was set to 3, and the learning rate was 0.1. The subsample ratio of training instances was set to 1, indicating the use of all training data. The objective function employed for binary classification was 'binary: logistic', and the number of boosting rounds (n\_estimators) was set to 100. The model was configured to use a single CPU core for parallelism (n\_jobs=1), 5-fold cross-validation was used, and the random seed number was set to 42 for reproducibility. Default values were applied for parameters such as gamma, min\_child\_weight, colsample\_bytree, reg\_alpha, reg\_lambda, scale\_pos\_weight, and verbosity, among others.

## **Supplementary Section 1.7 Multilayer Perceptron Neural Network**

A neural network model was employed for classification via the MLP classifier from Scikit-learn. Hyperparameter tuning was conducted through grid search optimization to identify the most effective configuration. The explored hyperparameters included variations in the hidden layer sizes, with options for architectures consisting of a single layer with 100 neurons or two layers with 50 neurons each. The rectified linear unit (ReLU) activation function was chosen, and the Adam solver was employed for optimization. The regularization term (alpha) was set to 0.0001 to control overfitting. The model's training was limited to a maximum of 200 iterations, and a random state of 42 was specified for reproducibility. Early stopping was enabled, with a validation fraction of 0.1 and a criterion of 10 consecutive iterations with no improvement (n\_iter\_no\_change), to halt the training process. The grid search was conducted via 5-fold cross-validation, and the best-performing neural network classifier was identified as the one with the optimal combination of hyperparameters. The resulting best classifier, determined through the grid search, was then further evaluated on the training data, and its performance metrics were used for subsequent analysis and interpretation.

## Supplementary Section 2. Details of fine-tuning the LLM with the QLoRA

QLoRA is a novel and efficient fine-tuning approach designed to reduce memory usage significantly. QLoRA introduces innovations such as a new 4-bit data type, double quantization to reduce the memory footprint, and paged optimizers to manage memory spikes. The approach demonstrates superior performance across various instruction datasets, model types (LLaMA, T5), and scales (e.g., 33B and 65B parameter models), showing state-of-the-art results through fine-tuning on a small high-quality dataset (20).

We enabled 4-bit loading, employed double quantization, utilized a quantization type of "nf4," and specified the computed data type as torch.bfloat16 for our model. We then created a tokenizer via the pretrained model and generated a causal language model, leveraging quantization configurations from the specified BitsAndBytesConfig ("bnb\_config"). We enabled gradient checkpointing for our model and prepared it for knowledge bit (KBit) training via the prepare\_model\_for\_kbit\_training function from the parameter-efficient fine-tuning (PEFT) library.

We utilized the PEFT library to create a LoraConfig object with specified parameters, including a 16-layer model with Lora attention, targeted projection modules, a dropout rate of 0.1, no bias, and a task type of 'CAUSAL\_LM'. We subsequently generated a PEFT model on the basis of this configuration.

We utilized the transformer library to establish a training pipeline for our language model. The pipeline involved initializing a trainer with our specified model, training dataset, and training arguments. The training configuration entailed setting the output directory to "logs" and conducting training for 4 epochs. We utilized a per-device batch size of 1, with gradient

accumulation steps of 4. The chosen optimizer was "paged\_adamw\_32bit", with a learning rate of  $2e-4$  and a weight decay of 0.001. We employed mixed-precision training with fp16 while disabling bfloat16. To control the maximum gradient norm, we set it to 0.3. The maximum number of training steps was not limited (-1), and we incorporated a warm-up ratio of 0.03 and enabled grouping by sequence length. The learning rate scheduler employed was "cosine". The training progress was reported to TensorBoard, and evaluation was performed at the end of each epoch. We disabled saving checkpoints during training by setting save\_steps to 0. Additionally, we set logging\_steps to 25 to log training progress at regular intervals. We specifically explored QLORA, combined with the bits and bytes library, to enhance language models while requiring fewer resources.

162 **Supplementary Table S1. Description of the raw prompts and improved prompts used during**  
163 **the experiments.**

Experiment	Raw Prompt	Improved Prompts	Improved prompt after 5 attempts
All LLMs including: Zero-shot classification including Mixtral-8x7B-Instruct-v0.1, Llama-3-70B, Mistral-7B-Instruct, Llama-3-8B, gpt-4o, gpt-3.5-turbo, gpt-4-turbo(g), gpt-4	"Does the patient survive or die based on the provided medical history? patient history is: "+["patient medical history"]	You're tasked with analyzing the present symptoms, past medical history,  laboratory data, age, and gender of COVID-19 patients to determine their outcome,  which is enclosed in square brackets. Your goal is to predict whether the patient will "survive" or "die" based on this information.	You're tasked with analyzing the present symptoms, past medical history,  laboratory data, age, and gender of COVID-19 patients to determine their outcome,  which is enclosed in square brackets. Your goal is to predict whether the patient will "survive" or "die" based on this information. Predict patient mortality in JUST ONE word and DO NOT answer vaguely.
Fine-tuned LLM (mistral-7b-instruct-v0.2)	(This prompt was not used in the fine-tuning task.)	You're tasked with analyzing the present symptoms, past medical history,  laboratory data, age, and gender of COVID-19 patients to determine their outcome,  which is enclosed in square brackets. Your goal is to predict whether the patient will "survive" or "die" based on this information.	(This prompt was not used in the fine-tuning task.)

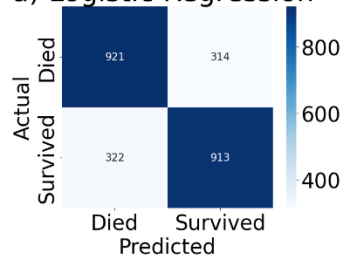
164

165 **Supplementary Table S2. Environments, model settings and model parameters for each LLM**

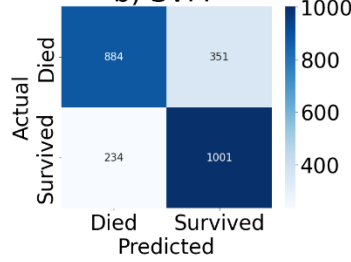
Model name in paper	Model names	Environment and model setting	Model parameter	Date/Source of use
Mixtral-8x7b	Mixtral-8x7b-Instruct-v0.1	Poe Web Interface	46.7 b	April 2024
Llama3-70b	Llama-3-70b	Poe Web Interface	70B	May 2024
Mistral-7b	Mistral-7b-Instruct	Poe Web Interface	7 B	April 2024
Llama3-8b	Llama-3-8b	Poe Web Interface	8 B	May 2024
GPT4o	gpt-4o-2024-05-13	OpenAI API: Temperature 1; max_tokens: 1024; seed: 123	175 B	June 2024
GPT3.5	gpt-3.5-turbo-0125	OpenAI API: Temperature 1; max_tokens: 1024; seed: 123	20 B	April 2024
GPT4T	gpt-4-turbo-2024-04-09	OpenAI API: Temperature 1; max_tokens: 1024; seed: 123	175 B	June 2024
GPT4	gpt-4-0613	OpenAI API: Temperature 1; max_tokens: 1024; seed: 123	175 B	June 2024
BERT (18)	bert-base-uncased	Transformers library (Python)	110 M	<a href="https://huggingface.co/google-bert/bert-base-uncased">huggingface.co/google-bert/bert-base-uncased</a>
ClinicalBERT <sup>a</sup> (19)	-	Transformers library (Python)	110 M	<a href="https://huggingface.co/medicalai/ClinicalBERT">huggingface.co/medicalai/ClinicalBERT</a>

166 Footnote: a: The ClinicalBERT model was developed using a large multicenter dataset comprising 1.2  
167 billion words covering a wide range of diseases. This base language model was further fine-tuned via  
168 electronic health records (EHRs) from over 3 million patient records to enhance its performance in clinical  
169 settings.

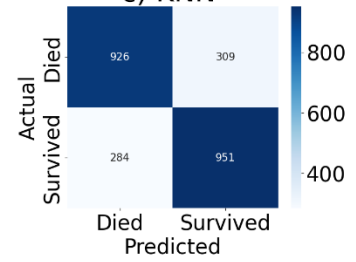
a) Logistic Regression



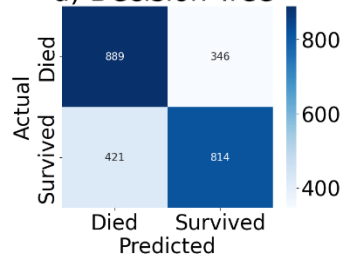
b) SVM



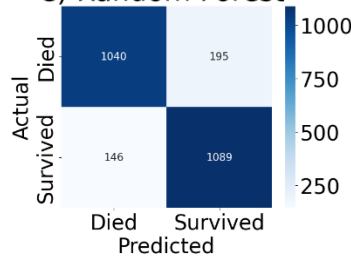
c) KNN



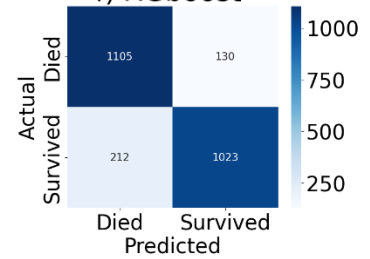
d) Decision Tree



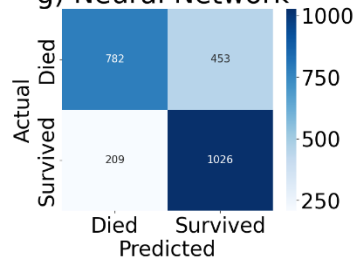
e) Random Forest



f) XGboost

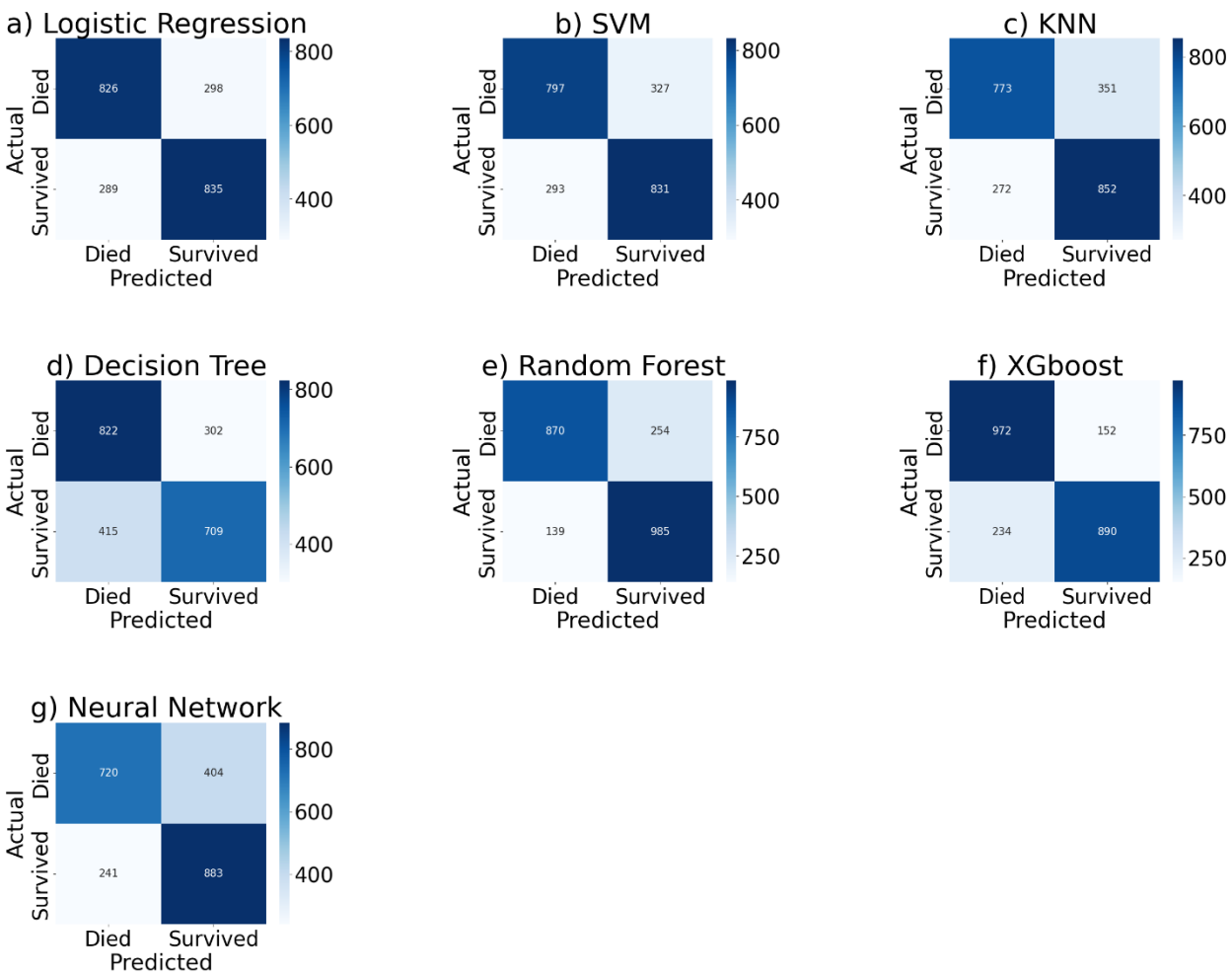


g) Neural Network



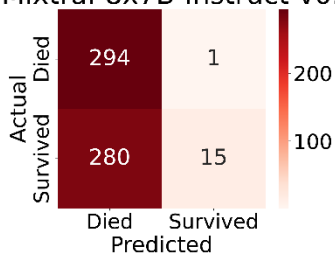
**Supplementary Figure S1. Confusion matrices for classical machine learning models (CMLs) trained on the internal validation test set, including logistic regression (a), support vector machine (b), K-nearest neighbor (c), decision tree (d), random forest (e), XGBoost (f), and multilayer perceptron (MLP)**



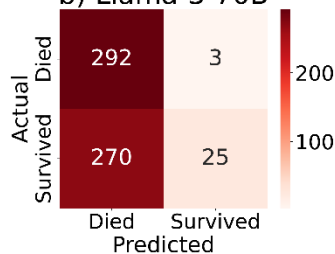


**Supplementary Figure S2. Confusion matrices for seven classical machine learning models tested on the external validation set: logistic regression (a), support vector machine (b), K-nearest neighbor (c), decision tree (d), random forest (e), XGBoost (f), and multilayer perceptron (MLP) (g)**

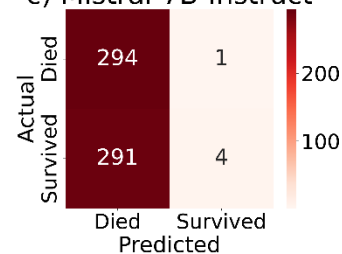
a) Mixtral-8x7B-Instruct-v0.1



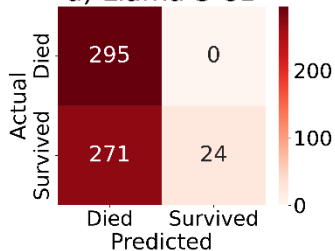
b) Llama-3-70B



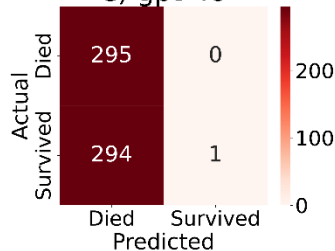
c) Mistral-7B-Instruct



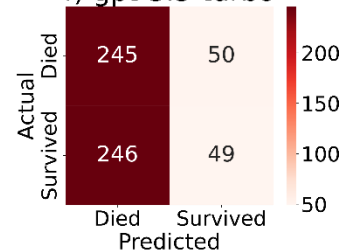
d) Llama-3-8B



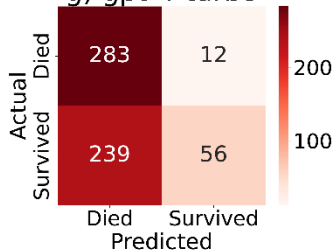
e) gpt-4o



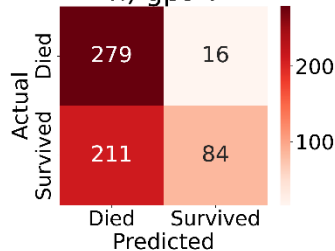
f) gpt-3.5-turbo



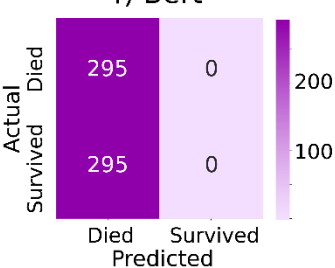
g) gpt-4-turbo



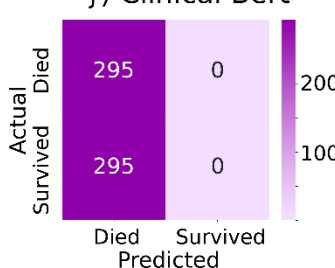
h) gpt-4



i) Bert

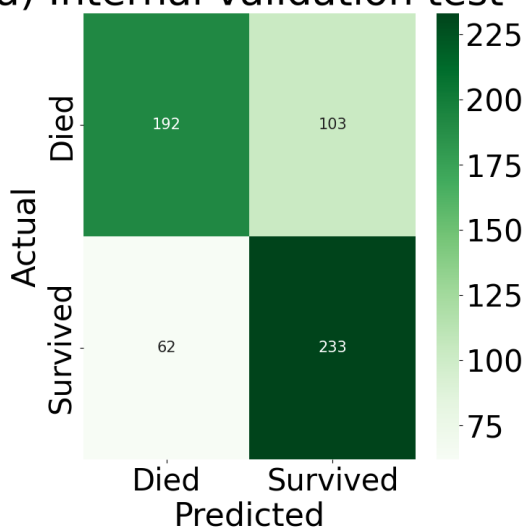


j) Clinical Bert

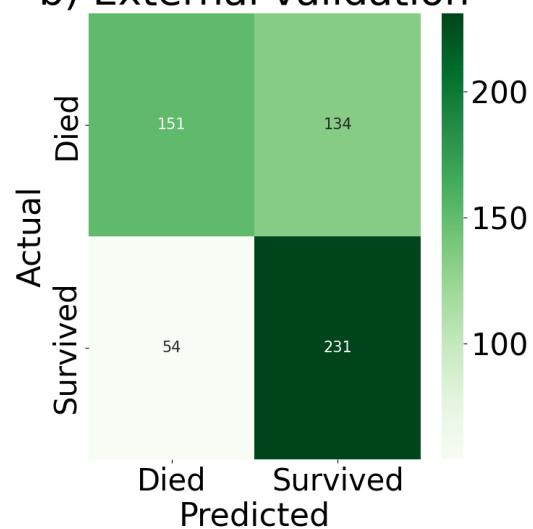


**Supplementary Figure S3. Confusion matrices for LLM zero-shot tasks, including Mixtral-8x7B-Instruct-v0.1 (a), Llama-3-70B (b), Mistral-7B-Instruct (c), Llama-3-8B (d), gpt-4o (e), gpt-3.5-turbo (f), gpt-4-turbo (g), and gpt-4 (h)**

a) Internal validation test

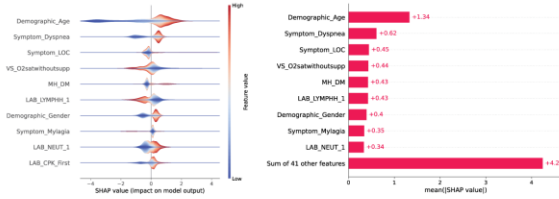


b) External validation

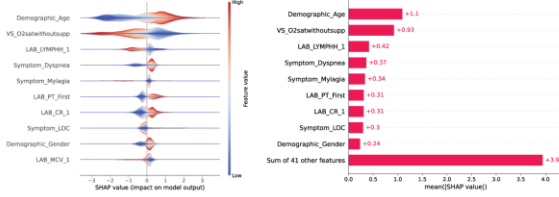


**Supplementary Figure S4. Confusion matrices for fine-tuned mistral-7b-instruct-v0.2 on the internal validation test set (a) and external validation set (b).**

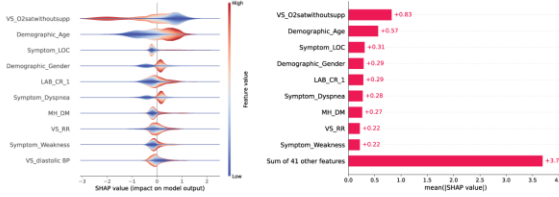
Logistic Regression



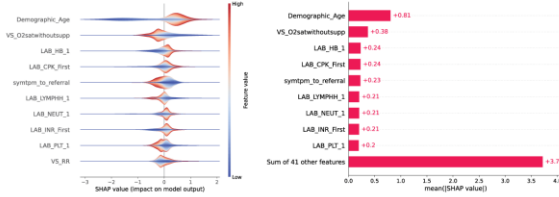
Support Vector Machine



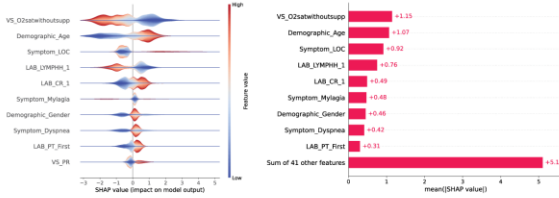
Decision Tree



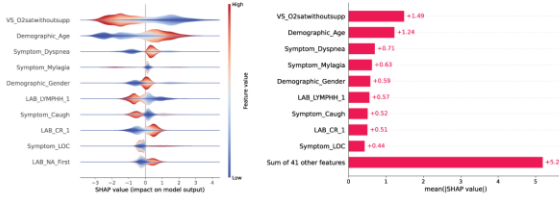
k-Nearest Neighbor



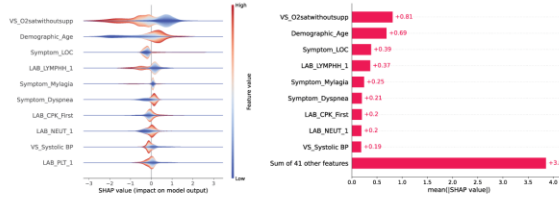
Random Forest



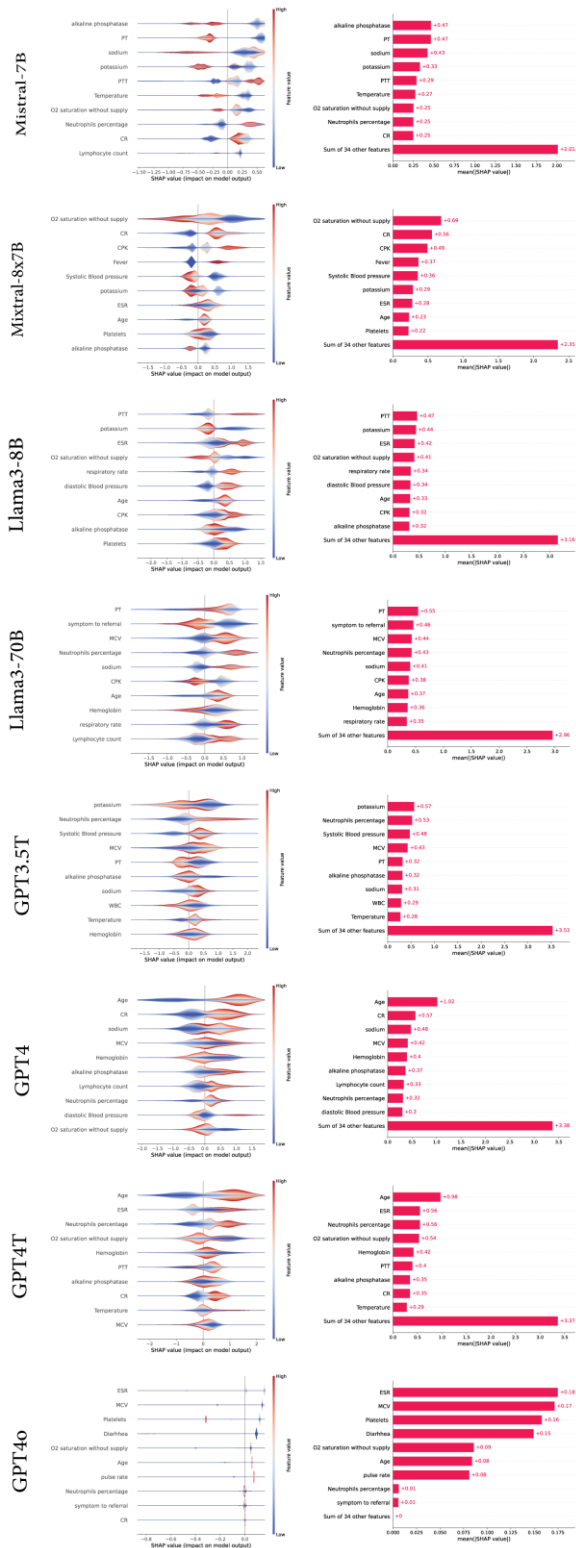
XGBoost



Multilayer Perceptron



**Supplementary Figure S5. Features with granular (left column) and global (right column) impacts on predictions among CMLs**



32

33 **Supplementary Figure S6. Features with granular (left column) and global (right column)**  
 34 **impacts on predictions among LLMs**