

1 **Real-World Performance of Large Language Models in Emergency**

2 **Department Chest Pain Triage**

3 Xiangbin Meng^{1,2*}, Jia-ming Ji^{3*}, Xiangyu Yan^{4*}, Hua Xu³, Jun gao¹, Junhong Wang⁵, Jingjia
4 Wang¹, Xuliang Wang¹, Yuan-geng-shuo Wang¹, Wen Yao Wang¹, Jing Chen⁶, Kuo Zhang⁶, Da
5 Liu⁷, Zifeng Qiu⁸, Muzi Li⁹, Chunli Shao¹, Yaodong Yang^{3#}, Yi-Da Tang^{1,2#}

6 1.Department of Cardiology and Institute of Vascular Medicine, Peking University Third
7 Hospital.

8 2.State Key Laboratory of Vascular Homeostasis and Remodeling, Peking University.

9 3. Institute for Artificial Intelligence, Peking University.

10 4. Institute of Disaster and Emergency Medicine, Tianjin University.

11 5. Emergency department, Peking university third hospital.

12 6. Department of Cardiology, State Key Laboratory of Cardiovascular Disease, Fuwai
13 Hospital, National Center for Cardiovascular Diseases, Chinese Academy of Medical
14 Sciences and Peking Union Medical College.

15 7.Department of Cardiology, the First Hospital of Hebei Medicical University, Graduate
16 School of Hebei Medical University.

17 8. Peking University Health Science Center, Peking University First Hospital

18 9. Peking University Health Science Center, Peking University People's Hospital

19 * These authors contributed equally to this study.

20 **Abbreviated Title:** LLMs in Emergency Department Chest Pain Triage

21 **Keywords:** Large Language Models (LLMs), Real-World Performance, Emergency

22 Department Chest Pain Triage, Acute Coronary Syndrome (ACS), Diagnostic capabilities

1 **Word count:** 295 (abstract); 7848 (excluding abstract, figure legends, and references)

2 **Number of figures and tables:** 3 Figures, 1 Table (7 Supplementary Figures, 4

3 Supplementary Table)

4 **Conflict of interest statement**

5 The authors have no potential conflict of interest to declare.

6 **Corresponding to: Yi-Da Tang, MD, PhD and Yao-dong Yang, PhD**

7 **Yi-Da Tang, MD, PhD**, Department of Cardiology and Institute of Vascular Medicine, Peking

8 University Third Hospital. State Key Laboratory of Vascular Homeostasis and Remodeling,

9 Peking University, Beijing, China. No.49 Huayuanbei Road, Beijing 100191, China (Email:

10 tangyida@bjmu.edu.cn).

11 **Yao-dong Yang, PhD**, Institute for Artificial Intelligence, Peking University, Beijing, China.

12 No.5 Yi HeYuan Road, Beijing 100871, China (Email: yaodong.yang@pku.edu.cn).

13

14

15

16

17

18

19

20

21

22

Abstract

1
2 **Background:** Large Language Models (LLMs) are increasingly being explored for medical
3 applications, particularly in emergency triage where rapid and accurate decision-making is crucial.
4 This study evaluates the diagnostic performance of two prominent Chinese LLMs, "Tongyi
5 Qianwen" and "Lingyi Zhihui," alongside a newly developed model, MediGuide-14B, comparing
6 their effectiveness with human medical experts in emergency chest pain triage.

7 **Methods:** Conducted at Peking University Third Hospital's emergency centers from June 2021 to
8 May 2023, this retrospective study involved 11,428 patients with chest pain symptoms. Data were
9 extracted from electronic medical records, excluding diagnostic test results, and used to assess the
10 models and human experts in a double-blind setup. The models' performances were evaluated
11 based on their accuracy, sensitivity, and specificity in diagnosing Acute Coronary Syndrome
12 (ACS).

13 **Findings:** "Lingyi Zhihui" demonstrated a diagnostic accuracy of 76.40%, sensitivity of 90.99%,
14 and specificity of 70.15%. "Tongyi Qianwen" showed an accuracy of 61.11%, sensitivity of
15 91.67%, and specificity of 47.95%. MediGuide-14B outperformed these models with an accuracy
16 of 84.52%, showcasing high sensitivity and commendable specificity. Human experts achieved
17 higher accuracy (86.37%) and specificity (89.26%) but lower sensitivity compared to the LLMs.
18 The study also highlighted the potential of LLMs to provide rapid triage decisions, significantly
19 faster than human experts, though with varying degrees of reliability and completeness in their
20 recommendations.

21 **Interpretation:** The study confirms the potential of LLMs in enhancing emergency medical
22 diagnostics, particularly in settings with limited resources. MediGuide-14B, with its tailored

1 training for medical applications, demonstrates considerable promise for clinical integration.
2 However, the variability in performance underscores the need for further fine-tuning and
3 contextual adaptation to improve reliability and efficacy in medical applications. Future research
4 should focus on optimizing LLMs for specific medical tasks and integrating them with
5 conventional medical systems to leverage their full potential in real-world settings.

6

1 **Introduction**

2 Large Language Models (LLMs) are revolutionizing the medical field, particularly in accelerating
3 pre-hospital triage¹⁻¹⁰. These models leverage deep learning and natural language processing
4 technologies to capture patterns and relationships in text through multi-layer neural network
5 architectures, enabling efficient processing and precise understanding of vast medical data⁴⁻¹¹.
6 Trained on extensive text data, LLMs have acquired a wealth of vocabulary, grammar, semantics,
7 and a condensed vast knowledge system, allowing them to respond coherently and accurately to
8 various symptom descriptions, medical history information, and literature queries provided by
9 users¹²⁻¹⁵. Importantly, as the model parameters and training data volume increase, scientists have
10 observed significant "emergence abilities" in LLMs¹⁶⁻¹⁸. This ability is not only reflected in the
11 efficient processing of complex information but also in their superior performance in logical
12 reasoning and innovative thinking¹⁹. This gives LLMs unique advantages in simulating human
13 thought patterns, understanding, and applying knowledge, especially in the field of medical
14 diagnostic assistance^{20,21}.
15 Globally, especially in remote areas of developing and developed countries, there is a severe
16 shortage of primary healthcare resources, characterized by insufficient facilities, lack of
17 professional personnel, and financial constraints^{22,23}. Trained medical providers, including doctors,
18 nurses, and other community health workers, are scarce, making it difficult to provide high-quality
19 primary healthcare services, further exacerbating the imbalance of medical human resources
20 between urban and rural areas²⁴. Traditional clinical prediction models based on machine learning
21 or deep learning, despite having theoretical application prospects, are rarely deployed in actual
22 clinical practice. The main reason is that these models generally lack generalizability and cannot

1 effectively handle the complex and variable actual clinical data, and the required parameters are
2 often not easily obtained in clinical settings²⁵. In contrast, LLMs, with their strong information
3 integration and logical reasoning abilities, extensive knowledge reserves, and seamless integration
4 with human language, can overcome these challenges.

5 In medical diagnostic assistance scenarios, the value of LLMs is particularly prominent.
6 Traditional diagnostic models often struggle to cope with patients' complex symptom expressions,
7 intricate medical histories, and vast medical literature due to limited processing capacity or
8 insufficient knowledge coverage. LLMs, with their vast training data and deep learning
9 architectures, can quickly organize and integrate various clinical information, using their
10 embedded extensive medical knowledge base to accurately classify and analyze diseases^{8,9}.

11 Furthermore, LLMs' logical reasoning ability allows them to conduct in-depth analysis of complex
12 disease clues, construct disease progression path models, predict potential complications, and
13 provide doctors with detailed and in-depth diagnostic support.

14 However, a cautiously optimistic attitude should be maintained towards the application of LLMs
15 in the medical environment, fully recognizing their limitations. These limitations include but are
16 not limited to: potential bias in training data, which may lead to unfairness in model
17 decision-making; challenges in explaining complex medical details, which may affect the
18 understanding and trust of doctors and patients in model outputs; and the risk of misdiagnosis due
19 to over-reliance on technology without necessary human supervision. Therefore, while affirming
20 the transformative potential of LLMs in healthcare, it is also crucial to focus on and address these
21 challenges to ensure their application in medical diagnostics is both safe and effective.

22 A noteworthy challenge arises when applying LLMs to languages such as Chinese, Japanese,

1 Korean, Tamil, Hindi, Thai, and Vietnamese, which employ "non-segmented text" structures that
2 markedly differ from English in terms of grammar, syntax, and usage. Most research endeavors
3 have predominantly concentrated on English-centric models like ChatGPT, leaving a notable
4 research gap in evaluating the diagnostic proficiencies of large language models trained
5 specifically for non-English environments^{26,27}. Although English and 'non-segmented text'
6 languages share similar fundamental principles in the development of LLMs, they face distinct
7 technical and engineering challenges in practical applications due to differences in language
8 characteristics and available resources. This leads to variations in their implementation and
9 performance. We must consider the fundamental differences in grammatical structures, data
10 resources, vocabulary size and distribution, algorithmic implementation and optimization, as well
11 as the adaptability of technical architectures across different languages. This research gap impacts
12 the broad adoption of these models in diverse linguistic contexts and profoundly influences the
13 level of trust vested in LLMs.

14 The emergency medical setting is distinguished by its immediacy, intricate clinical presentations,
15 and the imperative need for prompt diagnosis²⁸. The environment of emergency room is often
16 dynamic, extremely busy, and high-pressure, requiring healthcare personnel to make rapid
17 decisions and handle multiple cases simultaneously²⁹. A study in 2019 found that the average wait
18 time for emergency department patients was approximately 40 minutes before being seen by a
19 physician, with doctors spending an average of 13-24 minutes per patient during the consultation
20 ³⁰⁻³². Emergency triage systems are used globally to assess patient severity and allocate
21 resources³³⁻³⁵. The US uses Emergency Severity Index (ESI), a 5-tier system. China uses
22 Emergency Triage Scale/Standard (ETS), a 4-tier system. ETS is like ESI, with levels 1&2 triaged

1 to resuscitation. Patients in levels 3&4 wait to see a physician. Though ETS is generally accurate,
2 some critical patients wait hours and misdiagnosis is a pronounced concern, as research
3 underscores a notably elevated misdiagnosis rate within the emergency room^{36,37}. This issue is
4 exacerbated, particularly for common symptoms associated with myocardial ischemia, which are
5 susceptible to oversight or misjudgment³⁸. The guidelines recommend reperfusion therapy within
6 12 hours of the onset of myocardial infarction^{39,40}, yet approximately 70% of acute myocardial
7 infarction patients succumb to the disease due to the missed opportunity for timely treatment⁴¹,
8 highlighting the risk of misdiagnoses leading to treatment delays. LLMs stand poised to bring
9 about significant transformations in specific medical contexts, notably expediting pre-hospital
10 triage procedures. We see potential in these models to facilitate rapid triage by assisting healthcare
11 providers in swiftly processing patient data and offering potential diagnoses rooted in symptoms,
12 medical histories, and pertinent literature.

13 The medical profession demands precise and dependable tools for informed decision-making.
14 While LLMs hold potential, they present difficulties in understanding context and obtaining
15 clarifications⁴²⁻⁴⁷. Addressing real-world medical issues requires handling multiple data modalities
16 and must also provide authenticity, authority, accessibility, safety, empathy, and a human touch⁴⁸.

17 Real-world medical problems often transcend the confines of multiple-choice tests and structured
18 tasks. The human or AI model approach to diagnostic interaction, whether single or multi-turn
19 dialogues and their ability to process various data modalities are equally important. Indeed,
20 evaluating these vast medical models might not be any simpler than developing them.

21 Standardized testing has largely evaluated these models' medical knowledge reserves and
22 diagnostic logical reasoning capabilities^{2,6,8,49}. However, medical issues in the real world often

1 surpass the scope of structured tasks and multiple-choice tests, exhibiting greater complexity and
2 uncertainty^{7,50}. Currently, there is a particular lack of systematic evaluation globally on the
3 effectiveness of large language models in real medical environments, especially based on rich,
4 diverse, and dynamically changing real medical data⁵¹.

5 In this study, we focused on evaluating two prominent Chinese language models, "Tongyi
6 Qianwen (通千)(V1.0.3)"and "Lingyi Zhihui (灵医智慧) (V2.2.0)"⁵²⁻⁵⁴, which are
7 developed based on a Transformer's autoregression framework, akin to other globally recognized
8 models like Meta's LLaMa Series, Google's PaLM and LaMDA, and OpenAI's ChatGPT series.
9 As general-purpose large models, they have not been specifically fine-tuned for the medical
10 domain and are currently offered as online services by their respective operators. One of the core
11 objectives of our research is to conduct a comprehensive evaluation of these two models using
12 case data from Chinese patients, focusing particularly on their performance in processing Chinese
13 medical contexts, given that they are primarily trained on Chinese datasets, although they also
14 incorporate a certain proportion of English data. A key component of the study is a comparative
15 analysis of the diagnostic accuracy of "Tongyi Qianwen" and "Lingyi Zhizhi" in handling complex
16 and urgent medical data, benchmarking their results against medical experts' judgments. Targeted
17 enhancements and optimizations were applied to the fine-tuning and alignment phases,
18 particularly for LLMs tasked with medical applications. Recognizing the performance variability
19 of LLMs underscores the importance of meticulously establishing benchmarks that are apt for
20 medical artificial intelligence. We further integrated the insights gained from the benchmark
21 conducted into developing a new model called **MediGuide-14B**.

22 **Methods**

1 **Study Design and Setting**

2 This retrospective study was conducted at the Emergency Chest Pain Centers of the Peking
3 University Third Hospital Group, involving five tertiary-level centers. The study received ethical
4 approval from the ethics committee of Peking University Third Hospital (M2023828), complying
5 with the Helsinki Declaration.

6 **Data Sources**

7 Chief complaints, current medical history, past medical history, family history, and personal
8 history were extracted from electronic medical records as unstructured text content. It is important
9 to note that diagnostic test results such as electrocardiograms, myocardial enzyme tests, and
10 echocardiograms were not included in the test dataset provided to the test model and the control
11 group. Confirmed diagnostic information and related evidence were only made available during
12 the subsequent double-blind evaluation phase to the expert review committee, which served as the
13 authoritative reference standard. This approach aimed to ensure that experts could conduct
14 accurate and fair comparative analyses of the diagnostic accuracy of each group by combining
15 comprehensive and detailed medical test data with the model's predictive results.

16 **Participants**

17 The inclusion criteria for this study were patients who visited the emergency chest pain center at
18 Peking University Third Hospital's five tertiary-level centers due to chest pain-related symptoms
19 between June 2021 and May 2023. The exclusion criteria were: 1. Cases with significant omissions
20 or incompleteness in medical records, such as missing key clinical assessment records or essential
21 auxiliary examination results, which are crucial for a comprehensive patient evaluation; 2. Cases
22 where the chief complaint information did not come directly from the patients themselves due to

1 unstable vital signs or other reasons but was obtained through the accounts of others, making it
2 difficult to accurately trace and complete. Given that such situations could affect the credibility of
3 the study results due to the one-sidedness of the information or decreased accuracy of symptom
4 description, these cases were excluded in the analysis process. Ultimately, 11,428 patients who
5 met the above inclusion and exclusion criteria were included in the study.

6 **Outcome**

7 The study employed repeated random sampling for case selection, using Python 3.8 to randomly
8 select 100 patient cases from the database, repeated 1000 times. This method ensured a diverse
9 and random sample, reducing potential bias. All medical records were anonymized to maintain
10 patient privacy, further minimizing selection bias and enhancing the generalizability of the
11 findings.

12 The primary outcome of this study was the accuracy of diagnosing Acute Coronary Syndrome
13 (ACS). The study used LLMs prompts that included patient demographics, clinical symptoms, and
14 medical history, which are commonly found and critically important in primary healthcare and
15 emergency scenarios. This approach mimics the situations where lab results are not yet available,
16 and doctors and medical professionals must rely solely on the patient's chief complaints and past
17 medical history to triage chest pain and provide rapid management. In the study, the
18 cardiovascular physicians' group also followed the same information dimensions and problem
19 structure as the LLMs for case analysis. These cardiology specialists, during the diagnostic
20 process, similarly needed to interpret each case based on the patient's age, gender, chief complaint,
21 present illness history, family history, and personal history, and make diagnostic and therapeutic
22 decisions accordingly.

1 **Study Size**

2 To ensure a statistical power of 0.8 and effectively compare the diagnostic accuracy between
3 human medical experts and Large Language Models (LLMs) optimized for the medical
4 environment, it is necessary to determine an appropriate sample size. Based on preliminary data
5 and relevant literature, the diagnostic accuracy of human experts usually falls within the range of
6 0.9 to 0.95, while the accuracy of general-purpose LLMs not specifically trained is between 0.7
7 and 0.8. LLMs that have been optimized for the medical field have improved their accuracy to
8 levels comparable to human experts, although there may still be slight differences between the
9 two⁵⁵. In this context, to ensure a statistical power of 0.8, sufficient to detect the potential small
10 difference in accuracy between human experts and optimized LLMs, we calculated that each
11 group needs approximately 3835 samples. This sample size ensures that even minor differences in
12 accuracy can be detected with an 80% probability in statistical tests.

13 Considering the study subjects, we selected five centers affiliated with Peking University Third
14 Hospital, each of which receives about 5000 patients with chest pain annually. Given that clinical
15 data often have high noise, sparsity, and heterogeneity, which not only increase the difficulty of
16 data analysis but also may affect the robustness of the research conclusions, we decided to set the
17 study period from June 2021 to May 2023, totaling two years. This time frame was chosen to
18 accumulate sufficient case data, overcome the inherent complexity of clinical data, and ensure the
19 effectiveness and reliability of the final statistical analysis.

20 Figure 1 illustrates the flow chart of the entire study. The sample size of 11,428 records was
21 determined based on the hospital's patient flow and data availability during the study period. This

1 size was deemed sufficient to provide a robust analysis of the LLMs' diagnostic performance.

2 **(Figure 1)**

3 **LLMs:**

4 "Tongyi Qianwen," developed by Alibaba Group, is a 100 billion-parameter model with a diverse
5 data foundation, including web texts and specialized literature, and utilizes advanced
6 reinforcement learning techniques like A2C/A3C and PPO, and Q-learning methods such as DQN
7 and C51^{56,57}.

8 "Lingyi Zhihui," created by Baidu Group for medical contexts, has 260 billion parameters. It
9 combines auto-regressive and auto-encoding frameworks, suited for natural language
10 understanding and generation, and supports zero-shot, few-shot learning, and detailed
11 fine-tuning^{58,59}.

12 **Diagnostic performance:**

13 The anonymized dataset, excluding lab test results, was analyzed by "Tongyi Qianwen," "Lingyi
14 Zhihui," and a group of human experts. The human experts were eight cardiovascular specialists
15 with over ten years of experience, certified by the Chinese Society of Cardiology and the Chinese
16 Medical Association. Both LLMs and human experts were given detailed patient information
17 encompassing age, gender, chief complaints, present illness, and medical history. In this study, all
18 participants underwent comprehensive training before commencement to ensure a thorough
19 understanding of the research process. During the actual testing phase, we observed and recorded
20 the time taken by each diagnostic group to complete a test unit containing 100 medical case
21 records.

22 **Prompt Engineering :**

1 For large language models, appropriate prompts are necessary to activate their respective
2 capabilities. (Supplementary Table S1)

3 The prompt is: "1. Based on the patient's basic information, chief complaint, symptoms, and
4 medical history, what do you consider to be the patient's diagnosis? 2. Considering your
5 considered diagnosis, what further tests and medical advice would you recommend? 3. Please
6 evaluate the risk of the patient's condition based on the above patient information".

7 Similarly, a group of human experts will respond to the same questions based on the patient's basic
8 information, chief complaint, symptoms, and medical history.

9 **Reference standards:**

10 The gold standard of final diagnosis was the consensus of an independent cardiology expert panel
11 with 20 years of clinical service experience following Chest Pain Management Guidelines^{38,60,61}.
12 This panel had full access to patient records and cardiovascular lab test results, ensuring a
13 comprehensive and authoritative standard for comparison. The definitive diagnosis for cases
14 where there was uncertainty was established using Voting Mechanisms. (**Figure 1**)

15 **Statistical analysis:**

16 The statistical analysis in this study was conducted using the SPSS 27.0 statistical package for
17 Windows, Python version 3.8, and R software version 4.2.2, provided by the R Foundation for
18 Statistical Computing. All continuous variables that adhered to a normal distribution were
19 represented as means with their 95% confidence intervals (CI). To identify initial differences in
20 baseline characteristics between treatment groups, bivariate analyses were performed utilizing
21 Student's t-test. For comparative analyses among multiple groups, a one-way ANOVA test was
22 employed. A p-value of less than 0.05 was designated as the threshold for statistical significance.

1 In assessing the differences in diagnostic efficacy for cardiovascular diseases among the study
2 groups, this research employed a comprehensive and multifaceted set of evaluation metrics,
3 including Accuracy, Sensitivity, Specificity, False Positive Rate (FPR), Recall, and confusion
4 matrix diagrams among other multidimensional indicators. These multidimensional indicators
5 together form a rigorous and comprehensive performance evaluation framework, aimed at
6 comprehensively comparing and assessing the strengths and weaknesses of each study group in
7 terms of diagnostic accuracy and effectiveness.

8 Before evaluating the diagnostic metrics of each group, the study initially assessed the distribution
9 characteristics of the results from 1000 test units in each group using probability density curves,
10 P-P plots, and Q-Q plots to conduct normality tests.

11 **The Development of MediGuide-14B**

12 MediGuide-14B is developed on the foundation of the Qwen-14B model, undergoing extensive
13 optimization and specialized transformation through a meticulously crafted medical data training
14 and tuning program. Qwen-14B boasts 14 billion model parameters, endowing it with powerful
15 learning capabilities, ample knowledge reserves, and robust logical reasoning. During its
16 foundational training, Qwen-14B assimilated knowledge from over three trillion tokens, spanning
17 Chinese, English, and various other languages, including specialized domains like programming
18 and mathematics. Qwen-14B excels in natural language understanding, mathematical
19 problem-solving, logical reasoning, and computer programming. This base model supports
20 comprehensive fine-tuning, allowing for deep and customized adjustments tailored to various
21 tasks and domains.

22 The special medical database built by the research team includes detailed medical records of

1 105,290 outpatients and inpatients, totaling 2 million pieces of professional medical data that have
2 been carefully cleaned and protected for privacy. The training of MediGuide-14B is completed on
3 a high-performance server equipped with A800 80G*8. Leveraging the power of the DeepSpeed
4 framework and the Transformer architecture, we have optimized MediGuide-14B for better
5 performance and efficiency. This is crucial in handling the complexities and nuances of medical
6 diagnostics. The integration of Cross-Entropy Loss function and Reinforcement Learning from
7 Human Feedback (RLHF) in the training process further refines the model's accuracy and
8 human-like understanding, addressing the high sensitivity yet lower specificity issue identified in
9 previous models.

10 **Role of the Funding Source:**

11 In the design of the study; collection, analysis, and interpretation of data; writing of the report; and
12 the decision to submit the paper for publication, the study sponsors had no involvement. All
13 responsibilities and decisions regarding the research were made independently by the authors.

14 **Results**

15 **Overview of Study Population**

16 In the study involving 11,428 individuals who presented with emergency chest pain, after initially
17 assessing 12,015 potential participants, 587 were excluded due to significant gaps in their medical
18 records or indirect patient complaints. The study group had an average age of 64.82 years, with a
19 broad age distribution from 15 to 101 years, highlighting a significant elderly presence, underlined
20 by a median age of 67 years. Men constituted 65.4% of the participants.

21 The average Body Mass Index (BMI) for the cohort was 25.41, with a standard deviation of 3.69.

22 Medical evaluations revealed an average systolic blood pressure of 124.28 mmHg, diastolic blood

1 pressure of 75.35 mmHg, and heart rate of 70.38 bpm. The patient histories showed varying
2 prevalences of conditions: 8.7% had chest pain, 15.3% experienced dyspnea or chest tightness,
3 and 3.7% had episodes of syncope. Additionally, there were notable rates of chronic conditions,
4 including diabetes (7.9%), hypertension (23.5%), and hyperlipidemia (17.3%). Lifestyle factors
5 were also recorded, with 18.9% having a smoking history and 15.3% with a history of alcohol
6 consumption. In terms of emergency severity, 2.8% of cases were classified as critical or severe,
7 while 14.7% were urgent, and the majority, 82.5%, were less urgent. The diversity of
8 cardiovascular conditions was evident in the primary diagnoses. (Supplementary Table S2)

9 The disease composition spectrum of 11,428 patients was analyzed based on the "primary
10 diagnosis" of discharge diagnosis. The most common cardiovascular issues were NSTEMI/UA
11 Unstable Angina (24.3%), followed by Stable Angina Pectoris (14.8%) and STMI (7.4%). Other
12 cardiovascular diagnoses included Chronic Coronary Syndrome, Aortic Dissection, and Acute
13 Pulmonary Embolism. Hypertensive emergencies varied in severity and risk, with a range of
14 stages and risks documented. Arrhythmias formed a significant part of the diagnoses, with
15 conditions like Paroxysmal and Persistent Atrial Fibrillation, Atrial Flutter, and
16 Wolff-Parkinson-White Syndrome being prevalent. Heart failure variants were also noted, along
17 with other cardiac conditions such as Old Myocardial Infarction and Aortic Valve Insufficiency.
18 This detailed assessment underscores the wide spectrum of cardiovascular diseases managed in the
19 emergency setting, reflecting the complexity and diversity of the patient population.
20 (Supplementary Table S3)

21 **Performance of LLM**

22 We assessed the normality of the LLM's performance distribution using kurtosis, skewness,

1 probability density curve, P-P diagram, and Q-Q diagram. (Supplementary Figure S1)

2 Our analysis confirmed a normal distribution without significant outliers. Regarding the diagnosis

3 efficacy, "Tongyi Qianwen" achieved an accuracy of 61.11% (95% CI:60.84%-61.29), with a

4 sensitivity of 91.67% (95% CI:91.37%-91.96%) and a specificity of 47.95% (95%

5 CI:47.65%-48.25%). "Lingyi Zhihui" demonstrated an accuracy of 76.40% (95%

6 CI:76.17%-76.63%), with a sensitivity of 90.99% (95% CI:90.67%-91.31%) and a specificity of

7 70.15% (95% CI:69.85%-70.44%). The human experts were asked to perform the diagnostic test

8 based on the same content fed to LLMs. A total of 8 physicians completed this task. Human

9 experts achieved a mean accuracy of 86.37% (95% CI:86.18%-86.55%), a sensitivity of 79.62%

10 (95% CI:79.20%-80.04%), and a specificity of 89.26% (95% CI:89.06%-89.46%) (**Table 1**) .

11 The language models " Tongyi Qianwen "and "Lingyi Zhihui" completed the task in 24.68 ± 2.23

12 and 28.75 ± 3.25 minutes, respectively. On the other hand, human physicians completed the task

13 within a range of 65.25 ± 10.45 minutes (**Figure 1a**).

14 We plotted the parameters of diagnostic performance using radar charts. The area under the curve

15 for "Lingyi Zhihui" was 8094.76 units, which was more significant than the 5597.88 units for

16 "Tongyi Qianwen". However, human experts had the best overall performance, totaling 9393.36

17 units (**Figure 2a**). The area for "Tongyi Qianwen" performance primarily spans over the

18 "Sensitivity" region but is relatively smaller in the "Specificity" and "Accuracy" regions. "Lingyi

19 Zhihui" had a larger area in all three parts compared to "Tongyi Qianwen", especially in the

20 "Specificity" domain. The area for human experts was substantial in both the "Specificity" and

21 "Accuracy" regions but slightly smaller in the "Sensitivity" domain.

22 Both "Tongyi Qianwen" and "Lingyi Zhihui" demonstrated a high level of sensitivity,

1 indicating their capability to detect the majority of true ACS cases. However, their specificity was
2 comparatively lower, implying the potential for misclassifying some non-ACS cases as ACS. High
3 sensitivity is pivotal in screening tools, as they aim to identify the most genuine cases, even at the
4 risk of producing some false positives. Ensuring the accurate detection of real diseases or
5 abnormalities is a critical attribute of screening tools. Consequently, LLMs are well-suited as
6 screening tools, particularly in life-threatening emergency scenarios. In such situations, the
7 primary goal during screening is to identify as many true cases as possible, minimizing the risk of
8 overlooking vital information. However, it's important to acknowledge that a trade-off exists
9 between high sensitivity and specificity, meaning that while a screening tool can capture most
10 genuine cases, it may also generate some false positives (false alarms), which must be carefully
11 considered. **(Figure 2 b. c)**

12 "Tongyi Qianwen" model achieved a true positive rate (TPR) of 91.67% and a concomitant
13 false positive rate (FPR) of 52.05%. The model's accuracy is 43.10%, consistent with its recall. On
14 the other hand, "Lingyi Zhihui" shows that its TPR and recall rate are both 90.99%, but its FPR is
15 significantly reduced to 29.85%, and its accuracy rate is 56.87%. In contrast, the human expert's
16 TPR was 79.62%, the FPR was reduced considerably to 10.74%, and the accuracy was 76.34%,
17 consistent with its recall rate **(Figure 2b)**.

18 Among the cases misdiagnosed as ACS by the "Tongyi Qianwen" test, approximately 7.34%
19 (95% CI: 7.07%-7.60%) were eventually diagnosed as aortic dissection, and 3.45% were
20 diagnosed as acute pulmonary embolism according the reference standard (95% CI:
21 3.26%-3.63%). The rest were other non-ACS diseases with chest pain manifestations. Among the
22 cases misdiagnosed as ACS by the "Lingyi Zhihui", the average proportion of cases that were

1 eventually diagnosed as acute aortic dissection was approximately 9.14% (95% CI: 8.83%-9.44%).
2 The average proportion of patients who were eventually diagnosed with acute pulmonary
3 embolism was 2.83% (95% CI: 2.63%-3.03%).

4 Although human experts presented higher accuracy, there were still some cases misdiagnosed.
5 Among the total cases misdiagnosed as ACS by human experts, acute aortic dissection accounted
6 for 5.27% (95% CI: 4.80%-5.74%) and acute pulmonary embolism accounted for 0.96% (95% CI:
7 0.74%-1.18%). The discrepancies among LLMs and human experts are statistically significant.
8 (Supplementary Figure S2)

9 **Advancements in Medical Large Language Models: The Performance of MedGuide-13B**

10 After evaluating the performance of various commercially available closed-source Large
11 Language Models in medical diagnostics, we enhanced their capabilities by improving model
12 architecture, refining algorithms, and boosting fine-tuning and alignment techniques to increase
13 accuracy and reduce misdiagnoses. From these comprehensive benchmarks, we distilled key
14 insights that provided a solid foundation for the development of new language models.
15 Consequently, we developed the MediGuide-14B model, which was derived by making precise
16 adjustments to the Qwen-14B base model. The Qwen-14B model, known for its strong natural
17 language understanding and problem-solving capabilities, served as an ideal starting point for the
18 development of MediGuide-14B.

19 In the development process of MediGuide-14B, we first meticulously analyzed the issues
20 encountered by existing commercial general-purpose large language models when executing
21 medical tasks and accordingly implemented a series of targeted parameter optimizations to
22 enhance their performance in the healthcare domain. In the initial phase, our focus was on

1 bolstering the model’s understanding of medical terminology. This involved expanding the
2 medical professional vocabulary database and refining the model’s processing mechanisms for
3 these terms. During fine-tuning, we incorporated multi-turn dialogue data derived from real-world
4 medical scenarios involving 300,000 patients, significantly enhancing the model’s professionalism
5 and accuracy within medical contexts.

6 Employing supervised fine-tuning (Supervised Fine-Tuning, SFT), the fine-tuned large model
7 showed a significant improvement in accuracy when dealing with professional medical texts
8 compared to the original model. Subsequently, during the alignment process of the large model,
9 we introduced reinforcement learning from human feedback technology (Reinforcement Learning
10 from Human Feedback, RLHF) to guide the output distribution of the large model. We solicited
11 feedback and optimization from medical experts on the model’s outputs, thereby creating a
12 contrastive dataset imbued with human preferences, ensuring that the decision-making process of
13 the large model not only fully leverages its reasoning capabilities but also aligns with the
14 judgment standards of medical professionals. A reward model (Reward Model, RM) was trained
15 on this dataset, and reinforcement learning techniques were used to conduct further fine-tuning
16 alignment.

17 The aligned model (aligned model) following this process demonstrated a substantial
18 enhancement in its generalization ability when handling actual medical data, closely adhering to
19 the practical needs of medicine and effectively improving the accuracy of complex case analysis.

20 Lastly, in the model’s inference process, we employed a chain-of-thought decomposition method
21 where complex medical scenario questions posed by users were finely dissected to accurately
22 capture key information. This helped the model better comprehend the core content and logical

1 structure of the problem, thereby enhancing both the accuracy and relevance of its responses. After
2 such granular decomposition, the model independently analyzed each sub-problem before
3 synthesizing answers from all sub-problems to form a comprehensive and logically coherent final
4 answer.

5 Through the above-mentioned parameter optimizations and adjustments tailored for medical tasks,
6 MediGuide-14B has achieved a 44% improvement in capability over its predecessors.
7 (Supplementary Figure S3)

8 To assess the efficacy of large language models in diagnosing cardiovascular diseases, we
9 constructed the CVIDB test set, comprising 1,233 single-choice questions and 203 multiple-choice
10 questions. This standardized and high-quality test set provides detailed explanations for each
11 question, offering insights into the reasoning behind the correct answers and enhancing learning
12 and understanding of complex topics. The test set covers various subtypes, developmental stages,
13 and related complications and comorbidities of cardiovascular diseases, effectively testing the
14 depth and breadth of large language models' understanding of the field. We have made this test set
15 publicly available on GitHub for researchers and developers to download free of charge. The
16 access link is: <https://github.com/mengxiangbin123/CVIDB.git>

17 After completing the training and development of the MediGuide large model, we conducted a
18 series of standardized assessments, including several important medical benchmark tests. These
19 tests are^{55,62,63}: USMLE, a repository of simulated questions for the United States Medical
20 Licensing Examination; MedMCQA, a large-scale medical multiple-choice question dataset
21 covering various disciplines, derived from medical entrance exams in India; CMC, a large-scale
22 multitask knowledge assessment benchmark focusing on Chinese medical knowledge; and

1 MCMLE, a simulation of the Chinese medical qualification exam; along with the cardiovascular
2 disease-specific benchmark test set, CVIDB. These resources aim to comprehensively evaluate the
3 performance and generalization ability of large language models in medical knowledge and
4 clinical decision-making skills. We compared the performance of MediGuide-14B (V5.0) with
5 other leading models in the industry, including ChatGPT-4, ChatGPT-3.5 Turbo, Tongyi Qianwen
6 (v1.0.3), Lingyi Zhihui (v2.2.0), LLaMA 2-14B, and Qianwen-14B -Base.

7 Focusing on the United States Medical Licensing Examination (USMLE), ChatGPT-4
8 demonstrated superior performance with a score of 80.28%, closely followed by MediGuide-14B
9 at 78.63%, while LLaMA 2-13B trailed significantly with only 35.04%. For the MedMCQA,
10 which consists of multiple-choice questions from Indian medical entrance exams, ChatGPT-4
11 again led with a score of 72.51%, although here, the performance differences among the newer
12 models were relatively narrower. In contrast, models like ChatGPT-3.5 Turbo and Qianwen-14B
13 -Base showed relatively lower scores, 56.25% and 42.86%, respectively. The Composite Medical
14 Content (CMC) dataset, which assesses the models' understanding of medical knowledge
15 specifically in the Chinese context, saw MediGuide-14B performing the best with a score of
16 77.56%. ChatGPT-4 and Lingyi Zhihui also showed strong results with scores above 73%.

17 Performance on the Medical Chinese Medical Licensing Examination (MCMLE) again
18 highlighted the effectiveness of ChatGPT-4 and MediGuide-14B, which scored 74.58% and
19 75.41%, respectively, demonstrating their robustness in handling questions related to the Chinese
20 Medical Licensing Examination. Lower-tier models, such as LLaMA 2-13B, had notably weaker
21 performance, indicating possible challenges in their language-specific medical knowledge. Lastly,
22 in the Cardiovascular Disease Intelligence Diagnostic Benchmark (CVIDB), MediGuide-14B

1 exhibited exceptional capability, scoring the highest at 80.85%, showing its potential utility in
2 applications focused on cardiovascular health. ChatGPT-4 remained consistent across all
3 benchmarks with scores generally above 75%, reinforcing its overall reliability in medical domain
4 question answering. (Supplementary Table S4).

5 MediGuide-14B underwent a thorough evaluation process like that of 'Tongyi Qianwen' and
6 'Lingyi Zhihui.' It was tested using 1000 test units, each consisting of 100 distinct real-world cases
7 sourced from actual medical scenarios. This rigorous testing framework provided a comprehensive
8 assessment of MediGuide-14B's performance in real-life conditions. The model achieved an
9 impressive accuracy rate of 84.52%. It demonstrated high sensitivity in correctly identifying
10 positive results and commendable specificity in correctly identifying negative results
11 (Supplementary Figure S4).

12 **Extended recommendations by LLMs**

13 For the reference standards, we invited a panel of four distinguished cardiovascular specialists,
14 each with over twenty years of clinical experience. To further evaluate the possibility of LLM's
15 role in emergency Chest Pain Triage, we asked them to arbitrarily evaluate the treatment
16 recommendations generated from the prompts. The evaluation was based on established guidelines
17 for diagnosing and treating chest pain and full access to an array of essential patient data: from
18 electrocardiograms (ECGs) and cardiac enzyme tests to echocardiograms, NT-proBNP evaluations,
19 and when indicated, results from coronary angiography.

20 As illustrated in Supplementary Figure S5, 3.32% (95% CI:3.22%-3.41%) of the
21 recommendations generated by "Tongyi Qianwen" were deemed unsuitable. This resulted in
22 significant omissions of critical content that could potentially endanger patients. However, 12.88%

1 (95% CI:12.71%-13.04%) of the recommendations were considered reasonable, albeit incomplete,
2 with no direct harm to patients. The remaining 83.81% (95% CI:83.64%-83.98%) of
3 recommendations were classified as comprehensive and appropriate.

4 Regarding the "Lingyi Zhihui" model, 3.40% (95% CI:3.30%-3.50%) of recommendations were
5 deemed inappropriate with inherent risks. 43.44% (95% CI:43.18%-43.69%) were considered
6 reasonable but not exhaustive, devoid of direct patient harm. Meanwhile, 53.16% (95%
7 CI:52.91%-53.41%) of recommendations were thoroughly comprehensive and relevant.

8 The diagnostic suggestions from human experts were deemed that 2.48% (95% CI:0.28%-4.68%)
9 were inappropriate and could potentially compromise timely patient treatment. Another 12.96%
10 (95% CI:7.69%-18.23%) were considered reasonable but not comprehensive, while a substantial
11 84.56% (95% CI:74.45%-94.67%) were acknowledged as fully comprehensive and appropriate.

12 The performance assessment of the MediGuide-14B group in terms of treatment recommendations
13 reveals a nuanced picture. A small fraction, specifically 2.90% (95% CI:1.86%-3.94%), of the
14 recommendations were categorized as unreasonable and risky, highlighting areas where caution is
15 necessary. On the other hand, 13.52% (95% CI:11.38%-15.64%) of the recommendations were
16 deemed reasonable, albeit incomplete, suggesting a foundation of sound medical guidance that
17 could benefit from further elaboration or additional information. The majority, 83.58% (95%
18 CI:81.28%-85.98%) of the recommendations from the MediGuide-14B group stood out as both
19 reasonable and comprehensive, indicating a high level of proficiency in providing well-rounded
20 and thorough treatment advice. (Supplementary Figure S5)

21 **The impact of employment status of patients on the ACS Diagnostic Accuracy of LLMs**

1 Next, we sought to evaluate the impact of employment status of patients on LLM diagnostic
2 accuracy. We hypothesized that employment status may be linked to inherent characteristics that
3 could influence the information extracted from a patient's chief complaint. Given that individuals
4 covered by the Urban and Rural Resident Basic Medical Insurance (URRBMI) are typically
5 unemployed or self-employed, while those covered by the Urban Employee Basic Medical
6 Insurance (UEBMI) are typically employed by institutions, we leveraged health insurance data
7 extracted from medical records to infer the employment status of patients^{64,65}. In the case of
8 patients under the URRBMI insurance plan, "Tongyi Qianwen" exhibited a diagnostic accuracy of
9 58.56%, while it demonstrated a slightly higher accuracy of 60.74% among patients under the
10 UEBMI plan ($P < 0.05$). Interestingly, the diagnostic accuracy of "Lingyi Zhihui" was also affected
11 by the employment status of the patients: an accuracy rate of 77.62% under URRBMI and 74.75%
12 under UEBMI ($P < 0.05$).

13 In the evaluation of MediGuide-14B's performance across different insurance types, the group
14 demonstrated notable results in the realm of Supplementary Diagnosis accuracy. For cases covered
15 under the URRBMI, the MediGuide-14B achieved a mean accuracy of 83.65%, under the UEBMI
16 category, the mean accuracy recorded was slightly higher, at 85.04% ($P = 0.046$). (Supplementary
17 Figure S6)

18 **The impact of patient's history on LLMs diagnosis efficacy**

19 Medical history is crucial in diagnosis, offering insights into a patient's past health and disease risk
20 factors. While human doctors can interpret this data based on experience, LLMs face a challenge
21 in doing so effectively. We sought to analyze how medical history affects LLMs' diagnosis
22 performance.

1 In the absence of medical history, "Tongyi Qianwen" initially demonstrated a mean accuracy of
2 72.66% (95% CI: 72.43%-72.90%), a sensitivity of 81.65% (95% CI: 81.25%-82.06%), and a
3 specificity of 68.81% (95% CI: 68.53%-69.09%). Upon the inclusion of past medical histories, the
4 accuracy decreased to 61.110% (95% CI: 60.84%-61.29%). However, sensitivity increased to
5 91.67% (95% CI: 91.37%-91.96%), while specificity decreased to 47.95% (95% CI:
6 47.65%-48.25%). Additional details can be found in Figure 3. In the case of "Lingyi Zhihui,"
7 without medical history, the diagnostic accuracy was 74.17% (95% CI: 73.94%-74.39%),
8 sensitivity stood at 89.06% (95% CI: 88.73%-89.40%), and specificity at 67.78% (95% CI:
9 67.49%-68.07%) as depicted in Figure 3. Subsequently, with the incorporation of a more
10 comprehensive dataset, "Lingyi Zhihui" achieved an accuracy of 76.40% (95% CI:
11 76.17%-76.63%), sensitivity of 90.99% (95% CI: 90.68%-91.31%), and specificity of 70.15% (95%
12 CI: 69.85%-70.44%).

13 Regarding the "Tongyi Qianwen" model, the initial treatment suggestions in the absence of
14 medical history were deemed inappropriate and potentially harmful in 5.80% of instances (95% CI:
15 5.68%-5.93%). Approximately 51.33% (95% CI: 51.11%-51.55%) were considered reasonable but
16 incomplete, while 42.87% (95% CI: 42.64%-43.09%) were assessed as both comprehensive and
17 suitable. In subsequent recommendations when medical history was provided, the figures shifted
18 to 3.32% (95% CI: 3.22%-3.41%) being inappropriate, 12.88% (95% CI: 12.71%-13.04%) being
19 reasonable but partial, and 83.81% (95% CI: 83.64%-83.98%) being thorough and appropriate
20 (Fig. 3). For the "Lingyi Zhihui" model, the recommendations based on prompts without medical
21 history were categorized as inappropriate in 6.50% (95% CI: 6.37%-6.63%) cases, reasonable but
22 lacking in 28.45% (95% CI: 28.21%-28.69%), and both comprehensive and fitting in 65.05% (95%

1 CI: 64.80%-65.31%). When medical history was provided, "Lingyi Zhihui" recommendations
2 shifted to 3.40% (95% CI: 3.30%-3.50%) being inappropriate, 43.44% (95% CI: 43.19%-43.68%)
3 being reasonable but not exhaustive, and 53.16% (95% CI: 52.91%-53.41%) being comprehensive
4 and relevant. (Figure 3)

5 The removal of past medical history significantly impacted "Tongyi Qianwen," notably increasing
6 specificity by 20.86% and accuracy by 11.55%, while sensitivity declined by 10.02%. Conversely,
7 "Lingyi Zhihui" exhibited relatively minor changes, with specificity decreasing by 2.37%,
8 accuracy by 2.23%, and sensitivity by 1.93% (Supplementary Figure S7). Following the omission
9 of medical history, "Tongyi Qianwen" reduced the rate of inappropriate treatment
10 recommendations from 5.80% to 3.32%, while comprehensive and appropriate recommendations
11 surged from 42.87% to 83.81%. In contrast, "Lingyi Zhihui" saw a decline in inappropriate
12 recommendations from 6.50% to 3.40%, but comprehensive and appropriate recommendations
13 decreased from 65.05% to 53.16%. This indicates that both models decreased the frequency of
14 inappropriate recommendations when excluding medical history, with "Tongyi Qianwen" notably
15 enhancing comprehensive and appropriate suggestions while "Lingyi Zhihui" experienced a
16 decline.

17 In assessing the impact of removing past medical history on MediGuide-14B's model metrics, a
18 series of changes were observed. The accuracy of the model experienced a decrease of 2.90%.
19 Additionally, there was a 4.26% reduction in sensitivity, indicating a diminished capacity of the
20 model to correctly identify positive cases. Finally, the model's specificity also decreased by 1.58%,
21 reflecting a slight reduction in its ability to accurately identify negative cases. (Supplementary
22 Figure S7)

1 **Discussion**

2 As large language models (LLMs) continue to advance and find widespread application, they have
3 demonstrated transformative potential in medical tasks. However, evaluating the capabilities of
4 large language models is a complex and challenging scientific issue. Currently, the assessment of
5 these models' medical knowledge and logical reasoning abilities primarily relies on standardized
6 tests. Yet, real-world medical tasks often exceed the scope of structured tasks, presenting a high
7 level of complexity and uncertainty. There is a global lack of systematic evaluation of large
8 language models' effectiveness in actual medical settings. Moreover, the pre-training data for the
9 world's major language models is predominantly in English. For instance, the recently released
10 Llama3 has about 5% of its corpus in non-English languages; ChatGPT-3.5 has approximately
11 0.09905% of its pre-training data in Chinese. Even models intended primarily for Chinese
12 contexts, such as Tongyi Qianwen, Wenxin Yiyen, and Baichuan, have only 15-30% of their
13 datasets in Chinese. Considering the structural, grammatical, and usage differences between
14 non-segmented and segmented texts, the composition of different language families in LLMs'
15 pre-training datasets might affect their performance in various linguistic environments, which is a
16 scientific question worthy of in-depth discussion.

17 This study aims to fill these gaps, focusing on the specific applications of large language models
18 in emergency triage or consultation scenarios. This study compares the diagnostic performance of
19 AI-driven models and human expertise in triaging emergency chest pain cases. While previous
20 research has primarily focused on English-based ChatGPT models, this study is pioneering in
21 evaluating two LLMs designed for "non-segmented text" environments.

22 Traditional machine learning systems (MLS) use specific structured data from the Electronic

1 Emergency Triage System (EETS) to enhance the identification of critically ill patients. These
2 MLS employ a predictive model primarily using the CatBoost Python package and provide
3 real-time explanations via the SHAP method to help medical staff understand why certain patients
4 may require immediate treatment. However, these systems have limitations, such as potential
5 overfitting issues, a lack of effective capture of complex nonlinear relationships, and challenges in
6 processing unstructured data. While models built on traditional machine learning or deep learning
7 perform well on specific datasets, they generally lack generalizability and often show reduced
8 predictive power when patient populations and samples are changed. This is why currently, there
9 are no truly integrated predictive models in medical systems worldwide, and a significant gap
10 exists between academic research on predictive models and their clinical applications²⁵. Models
11 like GPT-4 and similar large language models, with their strong capabilities in understanding and
12 generating natural language, logical reasoning, and knowledge storage, show significant
13 advantages in handling various data types, including unstructured, multimodal, and dynamic data.
14 It's worth noting that "Tongyi Qianwen (LLM)" and "Lingyi Zhihui (LLM)" exhibited high
15 sensitivity but lower specificity, particularly when compared to human experts. This raises
16 concerns about potential overdiagnosis by the AI models, which could result in unnecessary tests
17 and treatments. However, high sensitivity is crucial for screening tools, as they aim to capture
18 most of the true cases, even if it leads to some false positives. The high sensitivity of the AI
19 models suggests their suitability as initial diagnostic tools to ensure potential positive cases are not
20 missed. While both LLMs demonstrated the ability to provide relevant medical advice, there were
21 notable differences in the depth and validity of their recommendations. This underscores the need
22 to optimize LLMs for medical scenarios before deployment in healthcare settings. It is essential to

1 utilize more realistic medical data during training and fine-tune the models to align with the
2 nuances of medical treatment itself ^{66,67}.

3 This study found that LLMs struggled to significantly improve diagnostic accuracy and treatment
4 recommendations when incorporating patients' medical histories. This could be attributed to two
5 key factors: Current LLMs rely on computational power and probabilistic calculations rather than
6 a deep understanding of disease mechanisms. Second, there's a need for more advanced algorithms
7 that can better extract relevant clinical information while filtering out noise from historical data.
8 The removal of past medical history significantly impacted "Tongyi Qianwen," leading to a
9 substantial increase in specificity and accuracy while decreasing sensitivity. Meanwhile, "Lingyi
10 Zhihui" exhibited minor changes in diagnostic metrics. Additionally, both models altered the
11 frequency of inappropriate treatment recommendations when medical history was omitted, with
12 "Tongyi Qianwen" improving its comprehensive and appropriate suggestions, while "Lingyi
13 Zhihui" declined. When utilizing LLMs for diagnostic support, healthcare practitioners should
14 acknowledge variations in how these models handle intricate and diverse data⁶⁸. Factors such as
15 model comprehensiveness, accuracy, and potential sources of interference should be considered.
16 Clinical judgment, rooted in experience, should guide decision-making. Continuous monitoring
17 and performance optimization are crucial. LLMs offer promise as diagnostic aids, but healthcare
18 professionals must weigh multiple factors to ensure the delivery of precise and thorough
19 diagnostic recommendations to patients.

20 Our study shifted focus to MediGuide-14B, our proprietary open-source model. This model has
21 been specifically fine-tuned for medical applications, providing us with an opportunity to
22 scrutinize its real-world efficacy. In our preceding analyses, we observed notable variations in

1 sensitivity and specificity across different LLMs, underscoring the need for a detailed examination
2 of each model's strengths and weaknesses. MediGuide-14B stands out as a large language model
3 dedicated to the medical sector, boasting an advanced integration of domain-specific datasets and
4 bespoke training approaches to optimize its diagnostic capabilities.

5 Conducting an exhaustive evaluation of MediGuide-14B's performance is pivotal not only for
6 gauging the broader applicability of LLMs in healthcare but also for charting the course for their
7 future development and potential areas of application. By juxtaposing MediGuide-14B against
8 other leading models in the field, we aim to deliver a nuanced appraisal of the model's accuracy,
9 efficiency, and reliability in medical diagnostics. This comparative analysis is intended to furnish
10 diverse insights and formative experiences, contributing significantly to the ongoing discourse on
11 the role and impact of large language models in healthcare research.

12 Our study illustrates the marked variability in the performance of different large language models
13 when processing real-world medical scenario information. The objective of our research is not
14 solely to compare and rank these models but to emphasize the adaptability and potential of LLMs
15 in the medical field. This realization underscores the necessity of careful and precise
16 benchmarking tailored for medical AI applications. We have also discovered that specific
17 fine-tuning and alignment of LLMs significantly enhance their ability to perform specialized tasks
18 within the medical domain, even on smaller-scale models. This finding is particularly significant
19 for vertical sectors like healthcare, as it suggests the feasibility of training and deploying efficient
20 medical LLMs at a lower cost. Such advancements allow us to apply cutting-edge AI technology
21 more effectively in clinical settings, thereby improving the quality and efficiency of healthcare
22 services.

1 LLMs have the potential to reshape certain aspects of healthcare, particularly in the context of
2 rapid pre-hospital Chest Pain Triage. The integration of these models has the potential to
3 streamline triage procedures, facilitating timely interventions even before a patient arrives at the
4 hospital. This not only enhances the effectiveness and scope of diagnostic and therapeutic
5 interventions but also promises to improve the efficiency of medical infrastructure, reduce patient
6 waiting times, alleviate the burden on emergency medical personnel, and ultimately alleviate the
7 financial strain on patients and healthcare systems⁶⁹⁻⁷¹. In conclusion, the diagnostic capabilities
8 demonstrated by LLMs, as evidenced in this study, underscore their significance in advancing the
9 field of rapid triage. It is reasonable to anticipate that soon, these models will play a pivotal role in
10 enhancing healthcare delivery, ultimately benefiting both patients and healthcare systems.

11 In an era increasingly dominated by AI, medical practitioners, particularly the younger generation,
12 will inevitably encounter an expanding array of medical AI entities. How they utilize AI, discern
13 which AI tools best assist them, and identify the specific functions where AI can provide support,
14 necessitates robust benchmarking efforts. Such benchmarks offer crucial guidance to healthcare
15 professionals in navigating the AI landscape. Our study aims to initiate this journey in the field,
16 laying a foundational step that we believe will serve as a vital reference for future researchers.
17 This work is poised to propel the further advancement and application of LLMs in healthcare,
18 ultimately aiding medical professionals in harnessing AI's full potential for improved patient care
19 and healthcare delivery.

20 **Limitations**

21 The reliance on retrospective data may introduce inherent biases, potentially impacting the
22 generalizability of the results. Additionally, LLMs' diagnostic performance in ACS triaging

1 scenarios might not reflect their capabilities in other medical conditions. The specificity
2 challenges highlighted in the study emphasize the need for broader and more diverse training
3 datasets. Furthermore, the comparative analysis between AI models and human experts, though
4 illuminating, is based on a restricted set of parameters, potentially overlooking nuanced aspects of
5 clinical decision-making. The study underscores the necessity of comprehensive, prospective
6 research to validate the findings and address these limitations.

7 **Declarations:**

8 This manuscript was edited by LLM “Tongyi Qianwen” for its English language, but human
9 authors read and made the final version.

1 **References:**

- 2 1. Ayers JW, Zhu Z, Poliak A, et al. Evaluating Artificial Intelligence Responses to Public
3 Health Questions. *JAMA Netw Open* 2023; **6**(6): e2317517.
- 4 2. Kanjee Z, Crowe B, Rodman A. Accuracy of a Generative Artificial Intelligence Model in a
5 Complex Diagnostic Challenge. *Jama* 2023; **330**(1): 78-80.
- 6 3. Minssen T, Vayena E, Cohen IG. The Challenges for Regulating Medical Use of ChatGPT
7 and Other Large Language Models. *Jama* 2023.
- 8 4. Will ChatGPT transform healthcare? *Nat Med* 2023; **29**(3): 505-6.
- 9 5. Gilbert S, Harvey H, Melvin T, Vollebregt E, Wicks P. Large language model AI chatbots
10 require approval as medical devices. *Nat Med* 2023.
- 11 6. Thirunavukarasu AJ, Ting DSJ, Elangovan K, Gutierrez L, Tan TF, Ting DSW. Large
12 language models in medicine. *Nat Med* 2023.
- 13 7. Singhal K, Azizi S, Tu T, et al. Large language models encode clinical knowledge. *Nature*
14 2023.
- 15 8. Howard A, Hope W, Gerada A. ChatGPT and antimicrobial advice: the end of the
16 consulting infection doctor? *Lancet Infect Dis* 2023; **23**(4): 405-6.
- 17 9. Arora A, Arora A. The promise of large language models in health care. *The Lancet* 2023;
18 **401**(10377): 641.
- 19 10. Jiang LY, Liu XC, Nejatian NP, et al. Health system-scale language models are
20 all-purpose prediction engines. *Nature* 2023; **619**(7969): 357-62.
- 21 11. Thapa S, Adhikari S. ChatGPT, Bard, and Large Language Models for Biomedical
22 Research: Opportunities and Pitfalls. *Ann Biomed Eng* 2023.

- 1 12. Ouyang L, Wu J, Jiang X, et al. Training language models to follow instructions with
2 human feedback. *arXiv pre-print server*2022.
- 3 13. Wayne, Zhou K, Li J, et al. A Survey of Large Language Models. *arXiv pre-print server*
4 2023.
- 5 14. Sharma P, Parasa S. ChatGPT and large language models in gastroenterology. *Nat Rev*
6 *Gastroenterol Hepatol*2023.
- 7 15. Li R, Kumar A, Chen JH. How Chatbots and Large Language Model Artificial Intelligence
8 Systems Will Reshape Modern Medicine: Fountain of Creativity or Pandora's Box? *JAMA*
9 *Intern Med*2023; **183**(6): 596-7.
- 10 16. Wei J, Wang X, Schuurmans D, et al. Chain-of-thought prompting elicits reasoning in
11 large language models. *Advances in Neural Information Processing Systems* 2022; **35**:
12 24824-37.
- 13 17. Miller K, Gunn E, Cochran A, et al. Use of Large Language Models and Artificial
14 Intelligence Tools in Works Submitted to Journal of Clinical Oncology. *Journal of clinical*
15 *oncology : official journal of the American Society of Clinical Oncology*2023; **41**(19): 3480-1.
- 16 18. Ayers JW, Poliak A, Dredze M, et al. Comparing Physician and Artificial Intelligence
17 Chatbot Responses to Patient Questions Posted to a Public Social Media Forum. *JAMA Intern*
18 *Med*2023; **183**(6): 589-96.
- 19 19. Wei J, Tay Y, Bommasani R, et al. Emergent Abilities of Large Language Models. *arXiv*
20 *pre-print server*2022.
- 21 20. Azizi Z, Alipour P, Gomez S, et al. Evaluating Recommendations About Atrial Fibrillation
22 for Patients and Clinicians Obtained From Chat-Based Artificial Intelligence Algorithms.

- 1 *Circulation: Arrhythmia and Electrophysiology* 2023; e012015.
- 2 21. Yang X, Chen A, PourNejatian N, et al. A large language model for electronic health
3 records. *NPJ Digit Med* 2022; **5**(1): 194.
- 4 22. Sheikh K, Ghaffar A. PRIMASYS: a health policy and systems research approach for the
5 assessment of country primary health care systems. *Health Research Policy and Systems*
6 2021; **19**(1): 31.
- 7 23. Mehmood A, Rowther AA, Kobusingye O, Hyder AA. Assessment of pre-hospital
8 emergency medical services in low-income settings using a health systems approach.
9 *International Journal of Emergency Medicine* 2018; **11**(1): 53.
- 10 24. Gizaw Z, Astale T, Kassie GM. What improves access to primary healthcare services in
11 rural communities? A systematic review. *BMC Primary Care* 2022; **23**(1): 313.
- 12 25. Markowitz F. All models are wrong and yours are useless: making clinical prediction
13 models impactful for patients. *npj Precision Oncology* 2024; **8**(1): 54.
- 14 26. Yeo YH, Samaan JS, Ng WH, et al. GPT-4 outperforms ChatGPT in answering
15 non-English questions related to cirrhosis. 2023.
- 16 27. Fang C, Ling J, Zhou J, et al. How does ChatGPT4 perform on Non-English National
17 Medical Licensing Examination? An Evaluation in Chinese Language. 2023.
- 18 28. Bijani M, Abedi S, Karimi S, Tehranineshat B. Major challenges and barriers in clinical
19 decision-making as perceived by emergency medical services personnel: a qualitative content
20 analysis. *BMC Emergency Medicine* 2021; **21**(1): 11.
- 21 29. Becker TK, Gausche-Hill M, Aswegan AL, et al. Ethical challenges in Emergency Medical
22 Services: controversies and recommendations. *Prehosp Disaster Med* 2013; **28**(5): 488-97.

- 1 30. Pines JM, Mullins PM, Cooper JK, Feng LB, Roth KE. National trends in emergency
2 department use, care patterns, and quality of care of older adults in the United States. *Journal*
3 *of the American Geriatrics Society* 2013; **61**(1): 12-7.
- 4 31. Vainieri M, Panero C, Coletta L. Waiting times in emergency departments: a resource
5 allocation or an efficiency issue? *BMC Health Serv Res* 2020; **20**(1): 549.
- 6 32. Neprash HT, Everhart A, McAlpine D, Smith LB, Sheridan B, Cross DA. Measuring
7 Primary Care Exam Length Using Electronic Health Record Data. *Med Care* 2021; **59**(1): 62-6.
- 8 33. Zhiting G, Jingfen J, Shuihong C, Minfei Y, Yuwei W, Sa W. Reliability and validity of the
9 four-level Chinese emergency triage scale in mainland China: A multicenter assessment.
10 *International Journal of Nursing Studies* 2020; **101**: 103447.
- 11 34. Tanabe P, Gimbel R, Yarnold PR, Kyriacou DN, Adams JG. Reliability and validity of
12 scores on The Emergency Severity Index version 3. *Academic emergency medicine* 2004;
13 **11**(1): 59-65.
- 14 35. Taboulet P, Maillard-Acker C, Ranchon G, et al. Triage des patients à l'accueil d'une
15 structure d'urgences. Présentation de l'échelle de tri élaborée par la Société française de
16 médecine d'urgence: la French Emergency Nurses Classification in Hospital (FRENCH).
17 *Annales françaises de médecine d'urgence* 2019; **9**(1): 51-9.
- 18 36. Kachalia A, Gandhi TK, Puopolo AL, et al. Missed and delayed diagnoses in the
19 emergency department: a study of closed malpractice claims from 4 liability insurers. *Annals of*
20 *emergency medicine* 2007; **49**(2): 196-205.
- 21 37. Hussain F, Cooper A, Carson-Stevens A, et al. Diagnostic error in the emergency
22 department: learning from national patient safety incident report analysis. *BMC Emerg Med*

- 1 2019; **19**(1): 77.
- 2 38. Gulati M, Levy PD, Mukherjee D, et al. 2021 AHA/ACC/ASE/CHEST/SAEM/SCCT/SCMR
3 Guideline for the Evaluation and Diagnosis of Chest Pain: Executive Summary: A Report of the
4 American College of Cardiology/American Heart Association Joint Committee on Clinical
5 Practice Guidelines. *Journal of the American College of Cardiology* 2021; **78**(22): 2218-61.
- 6 39. Amsterdam EA, Wenger NK, Brindis RG, et al. 2014 AHA/ACC guideline for the
7 management of patients with non–ST-elevation acute coronary syndromes: a report of the
8 American College of Cardiology/American Heart Association Task Force on Practice
9 Guidelines. *Journal of the American College of Cardiology* 2014; **64**(24): e139-e228.
- 10 40. Lawton JS, Tamis-Holland JE, Bangalore S, et al. 2021 ACC/AHA/SCAI guideline for
11 coronary artery revascularization: executive summary: a report of the American College of
12 Cardiology/American Heart Association Joint Committee on Clinical Practice Guidelines.
13 *Circulation* 2022; **145**(3): e4-e17.
- 14 41. Li J, Li X, Wang Q, et al. ST-segment elevation myocardial infarction in China from 2001
15 to 2011 (the China PEACE-Retrospective Acute Myocardial Infarction Study): a retrospective
16 analysis of hospital data. *The Lancet* 2015; **385**(9966): 441-51.
- 17 42. Komorowski M, Del Pilar Arias López M, Chang AC. How could ChatGPT impact my
18 practice as an intensivist? An overview of potential applications, risks and limitations. *Intensive*
19 *Care Med* 2023; **49**(7): 844-7.
- 20 43. Madden MG, McNicholas BA, Laffey JG. Assessing the usefulness of a large language
21 model to query and summarize unstructured medical notes in intensive care. *Intensive Care*
22 *Med* 2023.

- 1 44. Patel SB, Lam K. ChatGPT: the future of discharge summaries? *Lancet Digit Health* 2023;
2 5(3): e107-e8.
- 3 45. Ali SR, Dobbs TD, Hutchings HA, Whitaker IS. Using ChatGPT to write patient clinic
4 letters. *Lancet Digit Health* 2023; 5(4): e179-e81.
- 5 46. van Heerden AC, Pozuelo JR, Kohrt BA. Global Mental Health Services and the Impact of
6 Artificial Intelligence-Powered Large Language Models. *JAMA psychiatry* 2023; 80(7): 662-4.
- 7 47. Kwok KO, Wei WI, Tsoi MTF, et al. How can we transform travel medicine by leveraging
8 on AI-powered search engines? *Journal of Travel Medicine* 2023; 30(4).
- 9 48. Liu F, Panagiotakos D. Real-world data: a brief review of the methods, applications,
10 challenges and opportunities. *BMC Med Res Methodol* 2022; 22(1): 287.
- 11 49. Li S. Exploring the clinical capabilities and limitations of ChatGPT: a cautionary tale for
12 medical applications. *Int J Surg* 2023.
- 13 50. Wang G, Yang G, Du Z, Fan L, Li X. ClinicalGPT: Large Language Models Finetuned with
14 Diverse Medical Data and Comprehensive Evaluation. *arXiv pre-print server* 2023.
- 15 51. Shea YF, Lee CMY, Ip WCT, Luk DWA, Wong SSW. Use of GPT-4 to Analyze Medical
16 Records of Patients With Extensive Investigations and Delayed Diagnosis. *JAMA Netw Open*
17 2023; 6(8): e2325000.
- 18 52. Zhao WX, Zhou K, Li J, et al. A survey of large language models. *arXiv preprint*
19 *arXiv:230318223* 2023.
- 20 53. Sun Y, Wang S, Feng S, et al. Ernie 3.0: Large-scale knowledge enhanced pre-training
21 for language understanding and generation. *arXiv preprint arXiv:210702137* 2021.
- 22 54. Wang S, Sun Y, Xiang Y, et al. Ernie 3.0 titan: Exploring larger-scale knowledge

- 1 enhanced pre-training for language understanding and generation. *arXiv preprint*
2 *arXiv:2112.12731* 2021.
- 3 55. Singhal K, Azizi S, Tu T, et al. Large language models encode clinical knowledge. *Nature*
4 2023; **620**(7972): 172-80.
- 5 56. Rudolph J, Tan S, Tan S. War of the chatbots: Bard, Bing Chat, ChatGPT, Ernie and
6 beyond. The new AI gold rush and its impact on higher education. *Journal of Applied Learning*
7 *and Teaching* 2023; **6**(1).
- 8 57. Chien AA, Lin L, Nguyen H, Rao V, Sharma T, Wijayawardana R. Reducing the Carbon
9 Impact of Generative AI Inference (today and in 2035). Proceedings of the 2nd Workshop on
10 Sustainable Computer Systems; 2023; 2023. p. 1-7.
- 11 58. Peng C, Yang X, Chen A, et al. A study of generative large language model for medical
12 research and healthcare. *npj Digital Medicine* 2023; **6**(1): 210.
- 13 59. Li X, Fan Y, Cheng S. AIGC In China: Current Developments And Future Outlook. *arXiv*
14 *preprint arXiv:2308.08451* 2023.
- 15 60. Collet JP, Thiele H, Barbato E, et al. 2020 ESC Guidelines for the management of acute
16 coronary syndromes in patients presenting without persistent ST-segment elevation. *Eur Heart*
17 *J* 2021; **42**(14): 1289-367.
- 18 61. Association EMBotCM, Association CPBotCHIEP. Expert consensus on emergency
19 diagnosis and treatment of acute chest pain. *Chinese Journal of Emergency Medicine* 2019;
20 **28**(4): 413-20.
- 21 62. Yang A, Xiao B, Wang B, et al. Baichuan 2: Open large-scale language models. *arXiv*
22 *preprint arXiv:2309.10305* 2023.

- 1 63. Mbakwe AB, Lourentzou I, Celi LA, Mechanic OJ, Dagan A. ChatGPT passing USMLE
2 shines a spotlight on the flaws of medical education. Public Library of Science San Francisco,
3 CA USA; 2023. p. e0000205.
- 4 64. Yu H. Universal health insurance coverage for 1.3 billion people: What accounts for
5 China's success? *Health policy* 2015; **119**(9): 1145-52.
- 6 65. He W. Does the immediate reimbursement of medical insurance reduce the
7 socioeconomic inequality in health among the floating population? Evidence from China.
8 *International Journal for Equity in Health* 2023; **22**(1): 1-14.
- 9 66. Su P, Vijay-Shanker K. Investigation of improving the pre-training and fine-tuning of BERT
10 model for biomedical relation extraction. *BMC Bioinformatics* 2022; **23**(1): 120.
- 11 67. Wu C, Zhang X, Zhang Y, Wang Y, Xie W. PMC-LLaMA: Further Finetuning LLaMA on
12 Medical Papers. *arXiv pre-print server* 2023.
- 13 68. Ferryman K, Mackintosh M, Ghassemi M. Considering Biased Data as Informative
14 Artifacts in AI-Assisted Health Care. *N Engl J Med* 2023; **389**(9): 833-8.
- 15 69. The , Lancet. AI in medicine: creating a safe and equitable future. *Lancet* 2023;
16 **402**(10401): 503.
- 17 70. Han T, Lisa, Papaioannou J-M, et al. MedAlpaca -- An Open-Source Collection of Medical
18 Conversational AI Models and Training Data. *arXiv pre-print server* 2023.
- 19 71. Yunxiang L, Zihan L, Kai Z, Ruilong D, You Z. ChatDoctor: A Medical Chat Model
20 Fine-tuned on LLaMA Model using Medical Domain Knowledge. *arXiv pre-print server* 2023.

21

Author Contributions:

The study was primarily designed by Yi-Da Tang, Xiangbin Meng, Xiangyu Yan, and Jun Gao.

Individual Contributions:

- Yi-Da Tang: Led the study design, contributed cardiovascular expertise, involved in data collection and analysis, reviewed and approved the final manuscript.
- Xiangbin Meng: Contributed cardiovascular expertise, involved in data collection and analysis, reviewed and approved the final manuscript.
- Xiangyu Yan: Provided domain-specific knowledge, involved in data collection, analysis, and interpretation, reviewed and approved the final manuscript.
- Jun Gao: Contributed cardiovascular expertise, involved in data collection and analysis, reviewed and approved the final manuscript.
- Jingjia Wang, Xuliang Wang, Yuan-geng-shuo Wang, Wenyao Wang, Chunli Shao: Contributed cardiovascular expertise, were involved in data collection and analysis, reviewed and approved the final manuscript.
- Junhong Wang, Yaodong Yang, Muhan Zhang, Xiaojuan Cui, Jing Chen, Kuo Zhang, Da Liu, Jia-ming Ji, Zifeng Qiu, Muzi Li: Provided domain-specific knowledge in data collection, analysis, and interpretation, reviewed and approved the final manuscript.

Data Verification: Yi-Da Tang, Xiangbin Meng and Jun Gao directly accessed and verified the underlying data reported in the manuscript, ensuring its integrity.

Full Access and Responsibility: All authors confirm that they had full access to all the data in the study and accept the responsibility for the decision to submit for publication.

Collaboration and Coauthorship: The authors acknowledge and appreciate the collaboration

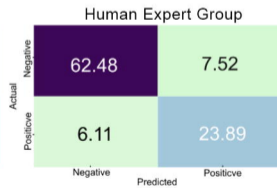
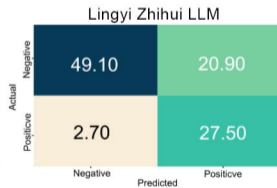
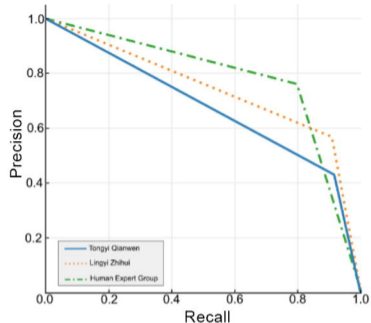
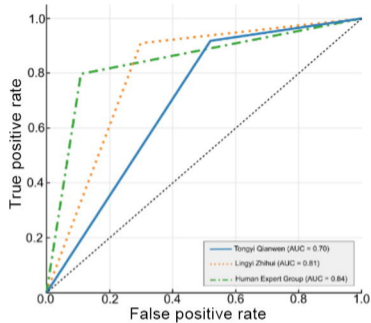
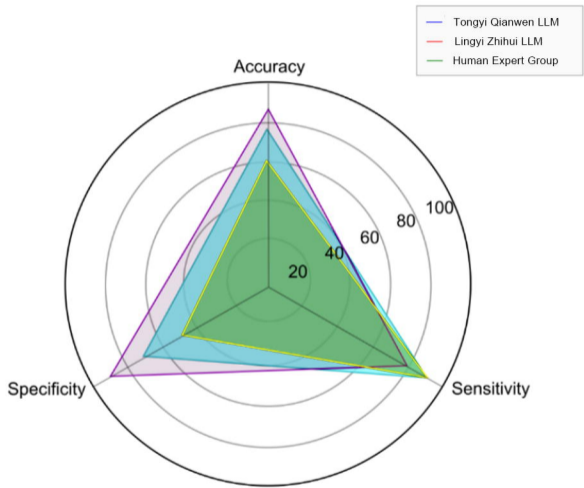
and coauthorship of colleagues in the locations where the research was conducted, reflecting the benefits of diversity in authorship in terms of background, career-stage, gender, geography, and race.

Acknowledgments and Personal Communications: All cited individuals in acknowledgments or personal communications have provided their written consent.

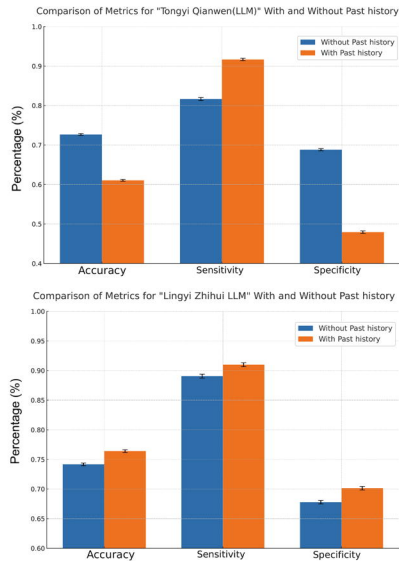
Author Statement Form: All authors have signed the author statement form, which will be uploaded with the submission.

Acknowledgments:

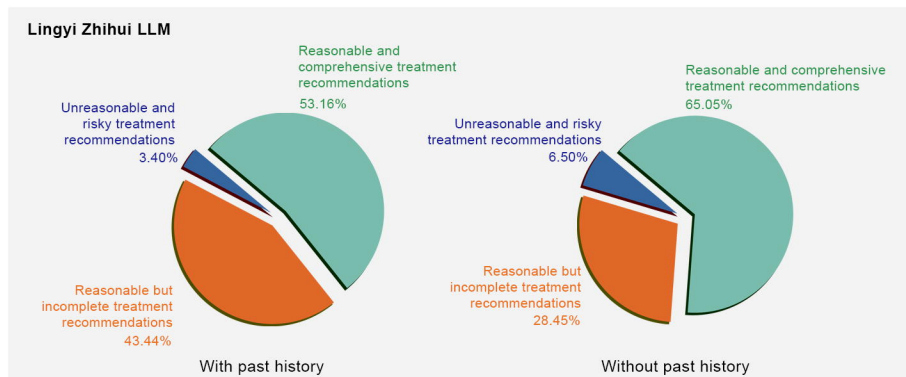
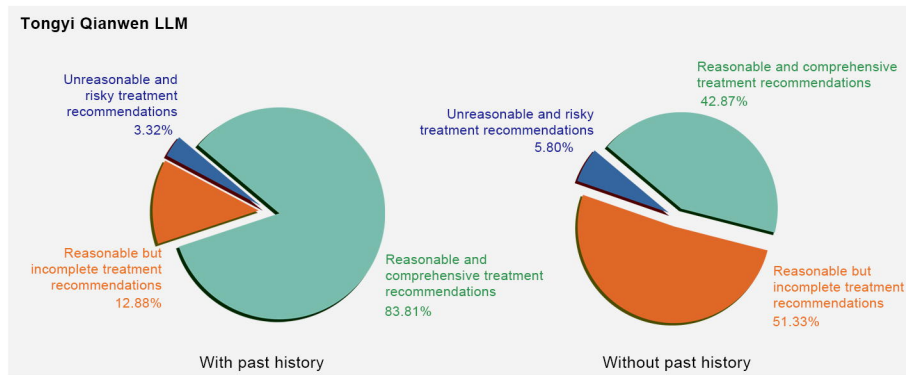
This study was funded by the National Key R&D Program of China (2020YFC2004705), National Natural Science Foundation of China (81825003, 91957123, 82270376), CAMS Innovation Fund for Medical Sciences (2022-I2M-C&T-B-119, 2021-I2M-5-003), Beijing Nova Program (Z201100006820002) from Beijing Municipal Science & Technology Commission, and CSC Special Fund for Clinical Research (CSCF2021A04).



a



b



c

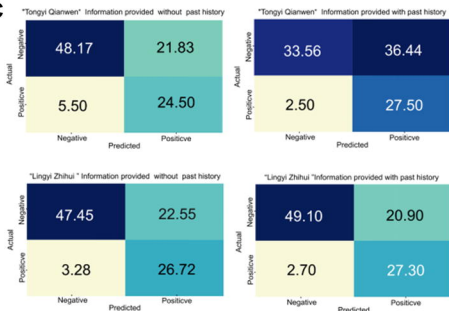
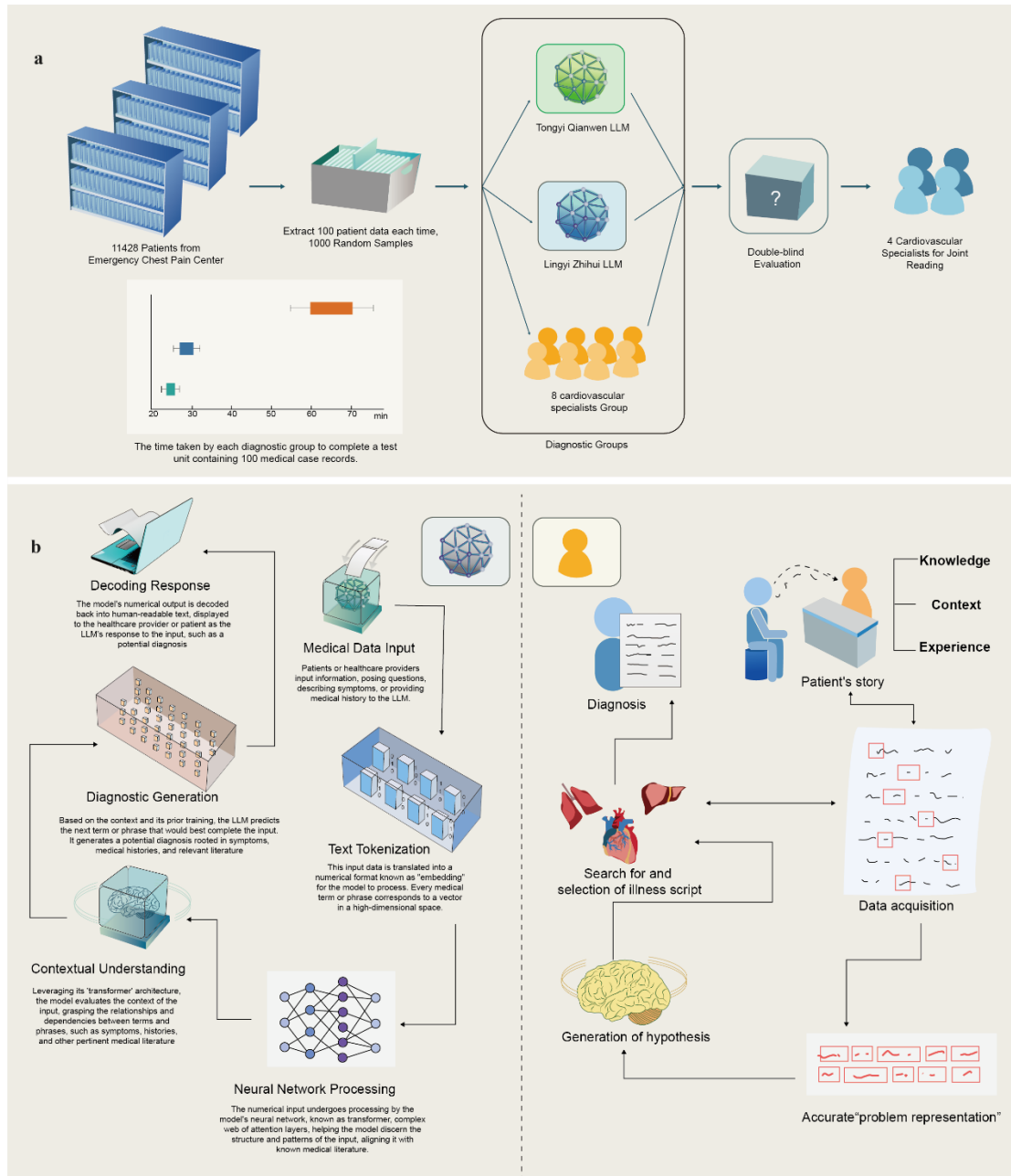


Table 1: Comparison of accuracy, sensitivity, and specificity in diagnosing emergency ACS between the two large language model diagnostic groups and the human expert diagnostic group.

Model/Groupe	Accuracy (95%CI)	Sensitivity (95%CI)	Specificity (95%CI)
TongyiQianwen (v1.0.3)	61.11% (95% CI:60.84%-61.29%)	91.67% (95%CI:91.37%-91.96%)	47.95% (95%CI:47.65%-48.25%)
Lingyi Zhihui (v2.2.0)	76.40% (95% CI:76.17%-76.63%)	90.99% (95%CI:90.67%-91.31%)	70.15% (95%CI:69.85%-70.44%)
Human Experts	86.37% (95%CI:86.18%-86.55%)	79.62% (95%CI:79.20%-80.04%)	89.26% (95%CI:89.06%-89.46%)

This table provides a detailed comparison of the key performance metrics when diagnosing emergency ACS among the two large language model diagnostic groups and the human expert diagnostic group.

Fig. 1: Research process diagram accompanied by a comparative illustration of diagnosis time between LLMs and human experts, along with a comparison of the diagnostic thought processes of LLMs and humans.



The research process diagram displays the primary steps and methodologies of this study. The comparative illustration showcases the time disparities between LLMs and human experts in completing diagnostic tasks. The thought process comparison further elucidates the cognitive and decision-making pathways employed by both during diagnosis.