
Reliability in AI-Assisted Critical Care: Assessing Large Language Model Robustness and Instruction Following for Cardiac Arrest Identification

Anonymous Author(s)

Affiliation

Address

email

Ugurcan Vurgun

Sy Hwang

Sunil Thomas

Ashley Batugo

Ana Acevedo

Aarthi Kaviyarasu

Benjamin S. Abella

Oscar J.L. Mitchell

Danielle L. Mowery

Abstract

This study systematically evaluates the performance, robustness, and instruction-following capabilities of large language models (LLMs) in identifying in-hospital cardiac arrest (IHCA) events. We assessed 51 open-source LLMs—comprising 36 general-purpose models and 15 medical-specific models—against GPT-4o, serving as a benchmark. Our analysis focused on model accuracy, robustness, and adherence to clinical instructions, with robustness measured using confidence intervals derived from non-parametric bootstrapping across several runs. While GPT-4o set a high standard with consistent performance across metrics, several open-source models demonstrated competitive results (e.g., Mistral-Nemo-Instruct-2407: F1: 0.84 ± 0.05 , Balanced Accuracy: 0.84 ± 0.04), albeit with greater variability. Medical-specific models showed strong recall, but often exhibited wider confidence intervals, indicating potential challenges in maintaining consistent performance. Instruction-following evaluation revealed that some general-purpose models excelled in adhering to clinical guidelines (e.g., unsloth/Meta-Llama-3.1-8B-Instruct: 99.0%), while certain medical models struggled with consistency. Our findings emphasize the potential of LLMs in critical care settings, highlighting the need to balance accuracy and instruction-following capabilities for reliable clinical deployment.

1 Introduction

1.1 Background on IHCA and importance of accurate detection

In-hospital cardiac arrest (IHCA) is a leading cause of death and disability, with almost 300,000 patients experiencing IHCA annually in the United States alone [1, 2]. While short- and long-term survival after IHCA have both improved gradually over the past several decades, overall survival remains low, with 20-30% of patients surviving to hospital discharge [3, 4, 5]. Patients who do survive frequently experience significant physical disability and neurological injury [3, 6]. Multiple interventions have been shown to improve outcomes after IHCA, and high-quality post-arrest care, including coronary angiography, targeted temperature management, expert neuroprognostication, and the avoidance of premature withdrawal of life-sustaining therapies, have been associated with reductions in both morbidity and mortality after IHCA [7, 8]. However, given the complex, multi-disciplinary, and time-sensitive nature of post-arrest care, adherence to post-arrest care guidelines is low, and variable implementation likely contributes to wide inter-hospital variation in survival

32 following IHCA [9, 10, 11, 12]. Unsurprisingly, quality improvement (QI) aimed at optimizing
33 post-arrest care has been identified as a key intervention to improve survival from IHCA [1].

34 While individual care teams can readily identify IHCA at the bedside, notification of QI groups or
35 post-arrest care teams after a patient has been successfully resuscitated relies on human processes that
36 often fail, especially for IHCA events occurring at night or on weekends. Traditional retrospective
37 methods used to identify IHCA patients include administrative billing codes, such as the International
38 Classification of Diseases-9 (ICD-9), and/or the use of manually assembled QI or research registries.
39 These methods are limited for use in patient care and QI efforts as there are often delays and
40 inaccuracies in coding and documentation [13, 14]. For example, the use of ICD-9 codes to identify
41 cardiac arrest patients has been found to lack both sensitivity and specificity, with one large registry
42 study demonstrating that over 36% of IHCA events were not identified using a combination of billing
43 and procedure codes [13, 14]. Inaccurate or delayed identification of patients in the immediate
44 post-IHCA period can further hamper the implementation of efforts designed to improve the quality
45 of post-IHCA care and IHCA patient outcomes. As such, the accurate identification of patients after
46 IHCA represents a modifiable factor to improve QI processes and reduce morbidity and mortality
47 after IHCA.

48 1.2 Overview of LLMs and their potential in healthcare

49 Large Language Models (LLMs) have demonstrated significant potential in healthcare, particularly in
50 processing and understanding complex medical texts [15]. Recent advancements have led to LLMs
51 capable of performing a wide range of tasks, including medical diagnosis, treatment recommendation,
52 and clinical documentation analysis [16, 17]. In the context of IHCA detection, LLMs offer several
53 advantages: they can interpret unstructured clinical notes to extract relevant information [18], analyze
54 the broader context of a patient’s medical history [15], process large volumes of data rapidly for
55 real-time or near-real-time detection [16], apply consistent criteria across cases to reduce human
56 variability and bias [17], and continuously update with new data and guidelines to stay current with
57 medical knowledge.

58 However, the application of LLMs in critical care settings like IHCA detection also presents chal-
59 lenges, including the need for high accuracy, robustness to variations in clinical documentation, and
60 the ability to follow specific medical instructions consistently.

61 1.3 Study Objectives and Hypotheses

62 The primary goal of this study is to systematically evaluate the performance, robustness, and
63 instruction-following capabilities of various LLMs in the context of IHCA detection. Specifically, the
64 objectives are to:

- 65 1. Assess the accuracy of multiple open-source LLMs in detecting IHCA events.
- 66 2. Evaluate the robustness of these models by analyzing performance consistency across
67 different runs.
- 68 3. Investigate the ability of LLMs to follow instructions.
- 69 4. Compare the effectiveness of general-purpose LLMs with models specifically fine-tuned for
70 medical applications.
- 71 5. Benchmark the performance of open-source LLMs against GPT-4o.

72 The study is grounded in the following hypotheses: While we expect GPT-4o to exhibit superior
73 performance, we hypothesize that several open-source models will achieve competitive results. We
74 also anticipate that medical-specific models may deliver higher precision in IHCA detection, but
75 could face challenges in maintaining consistent adherence to guidelines.

76 2 Methods

77 This retrospective, observational study was reviewed and approved by the University of Pennsylvania
78 Institute Review Board.

79 2.1 Data Collection and Preprocessing

80 We collected data from the Hospital of the University of Pennsylvania, focusing on adult (≥ 18 years
81 old) inpatient encounters from June 2018 to March 2022. The dataset included patients with reported
82 clinical emergency or rapid response calls, identified from a quality improvement database hosted on
83 REDCap, an encrypted, HIPAA-compliant, and Penn Medicine-approved data management software.
84 We extracted discharge summaries, progress notes, and assessments from the electronic health record
85 (EHR) and de-identified them using an adapted version of PHilter [19].

86 Our preprocessing pipeline involved several key steps:

- 87 1. **Note type selection:** We prioritized discharge summaries, progress notes, plan of care, and
88 assessments based on their likelihood of containing relevant IHCA information.
- 89 2. **Note sampling:** For each patient, we selected up to four notes, prioritizing them in the enu-
90 merated order listed above. This approach ensured a balanced representation of information
91 while focusing on the most informative note types.
- 92 3. **Labeling:** Notes were manually reviewed by clinical experts with emergency medicine
93 experience to determine true IHCA labels. Positive IHCA was defined as the loss of pulses
94 followed by the delivery of chest compressions.

95 For both GPT-4o and the open-source models, we developed a specialized instruction prompt:

```
96 INSTRUCTION_PROMPT = """
97 Given the patient's notes summary, assess whether the patient had an
98 in-hospital cardiac arrest (IHCA) event. IHCA refers to a cardiac arrest
99 that occurs in an inpatient setting where the patient is admitted and
100 receiving care. Cardiac arrests occurring in outpatient settings, such as
101 the hospital's cafeteria or during a visit, do not count as IHCA. Cardiac
102 arrest refers to the sudden loss of cardiac function leading to a lack of
103 a pulse or signs requiring that the patient receives chest compressions,
104 or ACLS.
105
106 Steps to determine IHCA:
107 1. Identify if there is a mention of cardiac arrest.
108 2. Determine if the cardiac arrest occurred in an inpatient setting.
109 3. Confirm if the patient received chest compressions.
110
111 Respond with 'positive' or 'negative' and give your rationale.
112 """
```

113 This prompt structure was designed to provide clear, consistent instructions to both GPT-4o and the
114 open-source models, enabling fair comparison across different model architectures and sizes.

115 3 Methods

116 3.1 Model Selection and Implementation

117 In this study, we evaluated 52 LLMs, including 51 open-source LLMs and GPT-4o, which served as a
118 benchmark. The open-source models included both general-purpose (36 models) and medical-specific
119 (15 models) variants, enabling a comprehensive comparison of their capabilities.

120 3.2 Experimental Design

121 The experimental design employed a zero-shot learning approach. To enhance reasoning capabilities,
122 we used a chain-of-thought prompting strategy that provided structured instructions for identifying
123 IHCA events. Each model processed the same set of patient notes with the same prompt, allowing for
124 a fair comparison (The temperature parameter was 0.2 for all models).

3.3 Evaluation Metrics

Model performance was evaluated using several metrics, including F1 score to measure the balance between precision and recall, balanced accuracy to account for class imbalance, and both precision and recall to assess the models' ability to correctly identify IHCA cases. To assess robustness, we derived confidence intervals through a non-parametric bootstrapping method with 100,000 resamples, providing a quantification of the uncertainty in our performance estimates without assumptions about the underlying data distribution.

3.4 Instruction Following Assessment

The models' ability to follow clinical instructions was evaluated by analyzing their responses to the provided prompt. We calculated an instruction-following score based on whether the models' predictions began with "positive" or "negative" within the first five tokens of their responses. This score was determined by the proportion of correct responses relative to the total number of predictions, offering insights into the models' reliability in adhering to specific formats in clinical settings.

3.5 Model Performance Comparison

We evaluated the performance of 52 different language models, including GPT-4o, on IHCA detection.

3.5.1 Dataset and Testing Methodology

The testing methodology differed significantly between GPT-4o and the open-source models. GPT-4o was evaluated on a balanced sample dataset of 21 patients due to API cost limitations, while the open-source models underwent a progressive evaluation, starting with (randomly sampled) smaller subsets of 10 and 100 documents, moving to 1,000 documents, and finally testing on the full dataset for the most promising models (F1 Score > 50%). All models were implemented using the HuggingFace Transformers library and executed on a Microsoft Azure Databricks cluster equipped with NVIDIA A100 GPUs, with an estimated total of approximately 720 GPU hours dedicated to this study.

4 Results

4.0.1 Dataset Characteristics

Our study analyzed 2,674 medical notes from 684 unique patients, with an uneven distribution of IHCA cases: 63.43% positive and 36.57% negative, emphasizing the need for metrics like balanced accuracy and F1 score, which account for uneven distributions. Note lengths had a maximum of 27,671 tokens, with an average of 2,334 tokens and a median of 1,911 tokens. For consistent evaluation, all LLMs were provided access to the first 7,000 tokens of each document, ensuring fair comparison between models with limited context windows (1K - 10K tokens) and those with capacities over 128K tokens.

4.1 F1 Score

In IHCA detection, a high F1 score indicates that the model can accurately identify cardiac arrest cases without generating excessive false positives or missing true events. This balance is crucial in a clinical setting where both over-identification (leading to unnecessary interventions) and under-identification (missing critical events) can have serious consequences. GPT-4o demonstrated the highest F1 score among all models (shown in Figure 1, see Appendix for other key plots), achieving 0.90 ± 0.01 . Among the open-source models, several exhibited competitive performance:

- *mistralai/Mistral-Nemo-Instruct-2407*: 0.84 ± 0.05
- *princeton-nlp/gemma-2-9b-it-SimPO*: 0.82 ± 0.07
- *THUDM/glm-4-9b-chat*: 0.79 ± 0.00
- *meta-llama/Meta-Llama-3.1-8B-Instruct*: 0.79 ± 0.05

Medical models like *ruslanmv/Medical-Llama3-8B* and *aaditya/Llama3-OpenBioLLM-8B* also showed competitive F1 scores (0.78 ± 0.01 and 0.76 ± 0.02 , respectively).

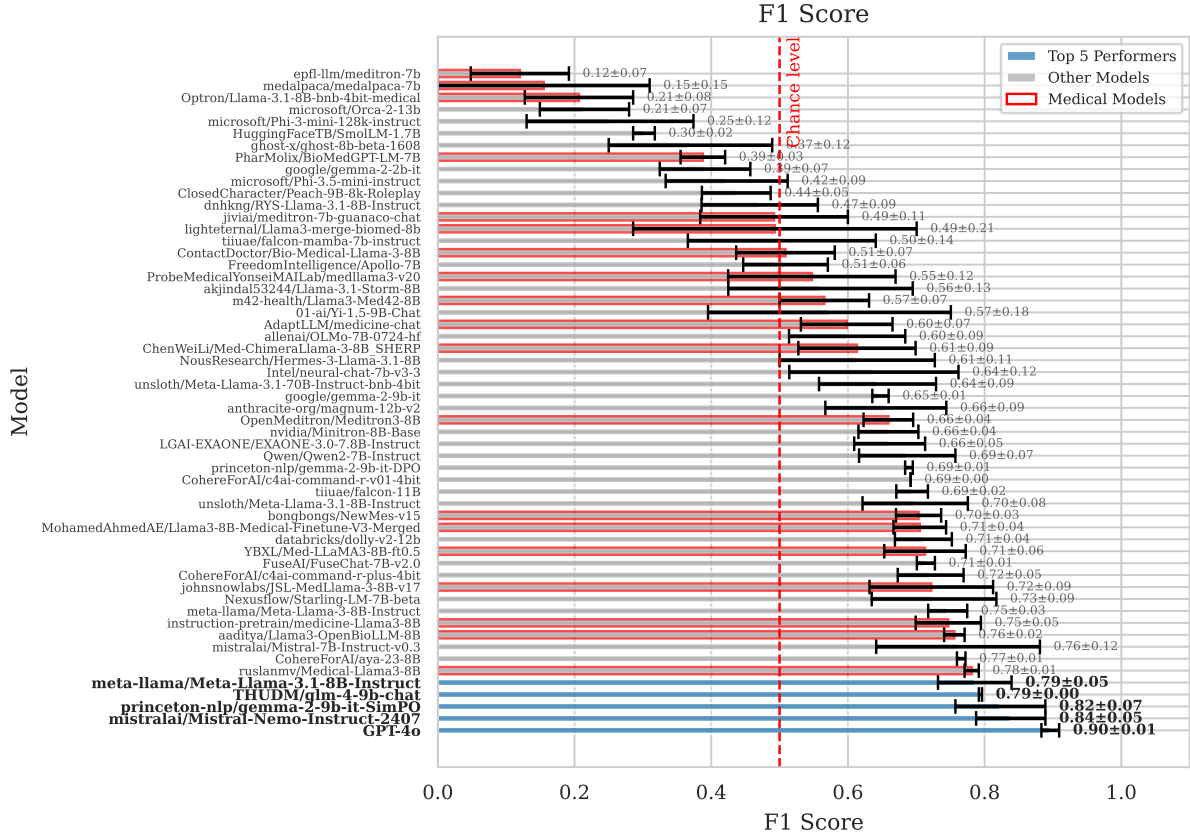


Figure 1: F1 score comparison of LLMs. This plot displays the F1 scores across models, providing insight into their precision and recall balance in IHCA detection. The confidence intervals were calculated with non-parametric bootstrapping across several runs of each model. The red dashed line represents chance level for the binary classification task.

4.2 Balanced Accuracy

Balanced accuracy is particularly important in our study due to the potential class imbalance in IHCA cases. GPT-4o led with a score of 0.84 ± 0.01 , followed by:

- princeton-nlp/gemma-2-9b-it-SimPO*: 0.84 ± 0.06
- mistralai/Mistral-Nemo-Instruct-2407*: 0.84 ± 0.04
- mistralai/Mistral-7B-Instruct-v0.3*: 0.76 ± 0.07
- unsloth/Meta-Llama-3.1-70B-Instruct-bnb-4bit*: 0.74 ± 0.05

Medical models such as *johnsnowlabs/JSL-MedLlama-3-8B-v17* (0.74 ± 0.03) exhibited slightly lower balanced accuracy, with similar confidence intervals.

4.3 Precision and Recall

In the context of IHCA detection:

- Precision (positive predictive value) indicates the proportion of correctly identified IHCA cases among all cases flagged as positive by the model. High precision minimizes false alarms, reducing unnecessary interventions and resource allocation.
- Recall (sensitivity) represents the proportion of actual IHCA cases correctly identified by the model. High recall is critical for ensuring that genuine IHCA events are not missed, as failing to detect a previous cardiac arrest could have severe consequences.

187 The trade-off between precision and recall is particularly important in IHCA detection, where the
188 cost of false negatives (missed IHCA events) may be considered higher than that of false positives.

189 4.3.1 Precision

190 Several open-source models demonstrated exceptionally high precision:

- 191 • *unsloth/Meta-Llama-3.1-70B-Instruct-bnb-4bit*: 0.99 ± 0.01
- 192 • *princeton-nlp/gemma-2-9b-it-SimPO*: 0.98 ± 0.02
- 193 • *microsoft/Phi-3.5-mini-instruct*: 0.97 ± 0.03
- 194 • *m42-health/Llama3-Med42-8B*: 0.97 ± 0.03
- 195 • *NousResearch/Hermes-3-Llama-3.1-8B*: 0.96 ± 0.04

196 Medical models like *m42-health/Llama3-Med42-8B* demonstrated strong precision, with narrow
197 confidence intervals similar to general-purpose models.

198 4.3.2 Recall

199 Models with high recall include:

- 200 • *GPT-4o*: 0.97 ± 0.03
- 201 • *CohereForAI/aya-23-8B*: 0.97 ± 0.02
- 202 • *THUDM/glm-4-9b-chat*: 0.95 ± 0.05
- 203 • *ruslanmv/Medical-Llama3-8B*: 0.94 ± 0.02
- 204 • *meta-llama/Meta-Llama-3-8B-Instruct*: 0.89 ± 0.10

205 Medical models like *ruslanmv/Medical-Llama3-8B* demonstrated strong recall, with confidence
206 intervals similar to general-purpose models.

207 4.4 Model Robustness and Variability

208 Robustness, as indicated by narrow confidence intervals, is crucial in clinical applications. A model
209 with consistent performance across different patient populations and clinical scenarios is more reliable
210 for real-world deployment. Wide confidence intervals suggest that a model’s performance may vary
211 significantly depending on the specific characteristics of the input data, which could be problematic
212 in diverse clinical settings.

213 The confidence intervals provide valuable insight into the robustness of the models. Narrow intervals,
214 as observed in *GPT-4o* (F1: 0.90 ± 0.01 , Balanced Accuracy: 0.90 ± 0.01), indicate consistent
215 performance across different samples, suggesting high reliability. *Mistral-Nemo-Instruct-2407* (F1:
216 0.84 ± 0.05 , Balanced Accuracy: 0.84 ± 0.04) shows slightly wider intervals, suggesting some
217 variability, but still relatively consistent performance. In contrast, wider intervals, such as those found
218 in *princeton-nlp/gemma-2-9b-it-SimPO* (F1: 0.82 ± 0.07 , Balanced Accuracy: 0.84 ± 0.06), indicate
219 greater variability and potential sensitivity to specific data samples.

220 Some medical models, like *MohamedAhmedAE/Llama3-8B-Medical-Finetune-V3-Merged* (F1: 0.71
221 ± 0.12), exhibited even wider confidence intervals, reflecting a lower degree of reliability in different
222 contexts. This suggests that these models may be less dependable when applied to varied datasets,
223 showing more significant performance fluctuations across different samples.

224 4.5 Instruction Following Performance

225 In a clinical context, the ability to follow specific instructions is paramount. This metric assesses
226 how well a model adheres to the given prompt structure, which is crucial for integration into existing
227 clinical workflows and for ensuring that the model’s outputs are in a format that is immediately useful
228 to healthcare providers, e.g., in our use case, consistently indicating “positive” or “negative” in its
229 response. Instruction following capabilities varied widely among the models (see Figure 2):

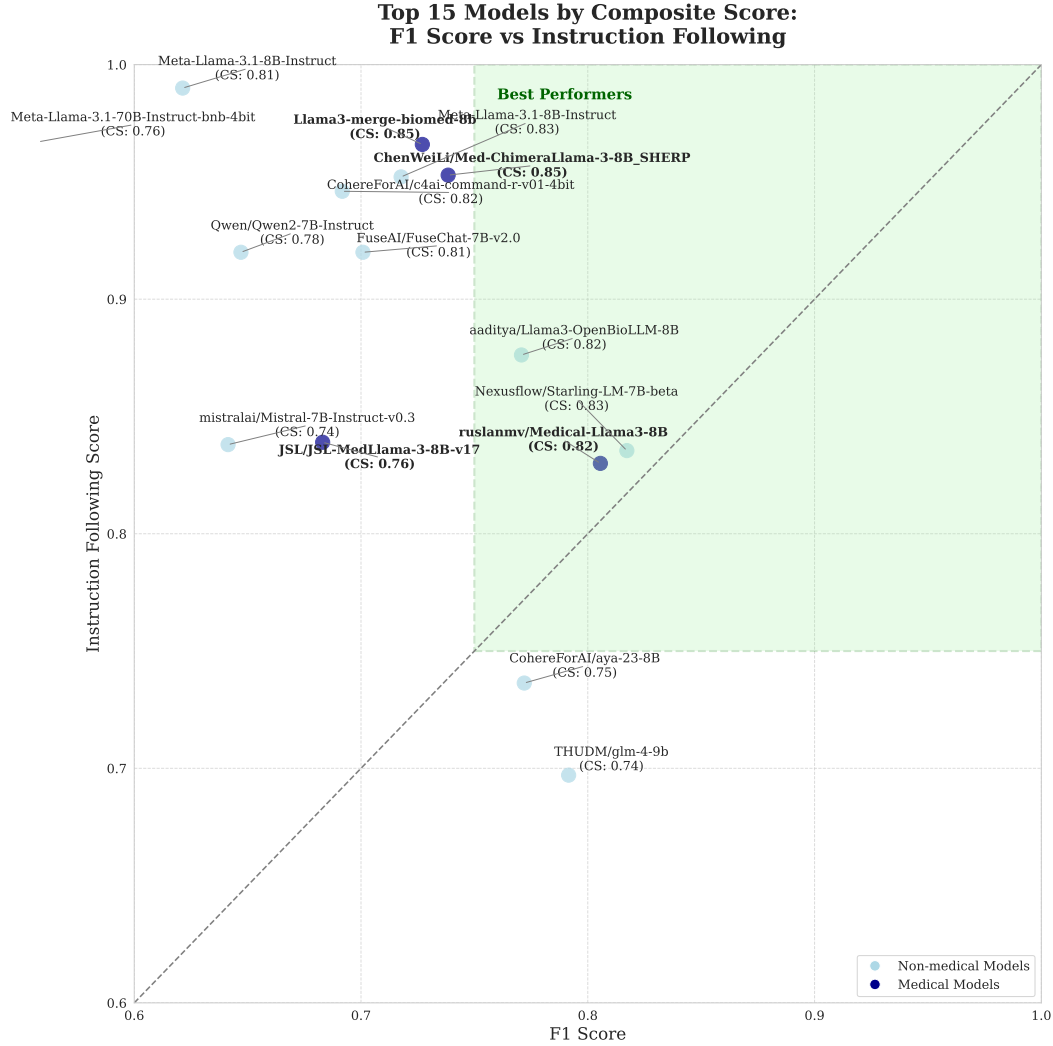


Figure 2: Top 15 models by composite score (CS). This plot illustrates the relationship between F1 score and instruction-following ability for the top 15 models, with the composite score representing the average of these two metrics. Medical model labels are shown in bold.

- General-purpose models like *unsloth/Meta-Llama-3.1-8B-Instruct* and *unsloth/Meta-Llama-3.1-70B-Instruct-bnb-4bit* excelled, achieving scores of 99.0% and 96.7%, respectively.
- Some medical models performed well, with *lighteternal/Llama3-merge-biomed-8b* and *ChenWeiLi/Med-ChimeraLlama-3-8B_SHERP* scoring 96.6% and 95.3%, respectively.
- However, other medical models struggled, with *MohamedAhmedAE/Llama3-8B-Medical-Finetune-V3-Merged* scoring as low as 32.6%.

This wide range in instruction following performance highlights the importance of this metric in assessing model reliability for clinical applications.

4.6 Overall Performance

GPT-4o demonstrated superior performance across all metrics, with high scores and low variability. Several open-source models, particularly larger variants like *Mistral-Nemo-Instruct-2407* and *princeton-nlp/gemma-2-9b-it-SimPO*, showed competitive performance, often achieving F1 scores and balanced accuracy within 10% of GPT-4o’s performance.

The variability in performance, as indicated by the confidence intervals, suggests that while some models (e.g., GPT-4o and *ruslanmv/Medical-Llama3-8B*) are highly consistent, others may be more sensitive to the specific characteristics of the input data. This variability is an important consideration for real-world deployment, where consistent performance across diverse patient populations is crucial.

Furthermore, the evaluation of instruction-following capabilities highlights the importance of model responsiveness in clinical contexts. While some medical models lag behind in instruction-following performance, surprisingly, top-performing general-purpose models demonstrate both strong task performance and reliable instruction adherence, making them potentially more versatile in real-world applications.

5 Discussion

5.1 Interpretation of Results

Our comprehensive evaluation of 52 large language models, including GPT-4o and 51 open-source alternatives (15 medical-specific and 36 general-purpose), reveals several key insights into their potential for IHCA detection:

- **Performance of Open-Source Models:** While GPT-4o demonstrated superior performance (F1: 0.90 ± 0.01 , Balanced Accuracy: 0.90 ± 0.01), several open-source models showed competitive results. Notably, *mistralai/Mistral-Nemo-Instruct-2407* (F1: 0.84 ± 0.05) and *princeton-nlp/gemma-2-9b-it-SimPO* (F1: 0.82 ± 0.07) achieved performance within 10% of GPT-4o, suggesting robust IHCA detection capabilities in accessible alternatives.
- **Impact of Medical Fine-tuning:** Medical-specific models showed mixed results. Some, like *ruslanmv/Medical-Llama3-8B* (Recall: 0.94 ± 0.02), demonstrated strong performance in certain metrics, but often exhibited wider confidence intervals compared to top general-purpose models. This suggests that fine-tuning on medical text can enhance domain-specific capabilities, but may not always translate to superior overall performance in IHCA detection.
- **Variability within Model Families:** We observed significant performance variability among models from the same family. Llama variants showed F1 scores ranging from 0.79 ± 0.05 (*meta-llama/Meta-Llama-3.1-8B-Instruct*) to 0.49 ± 0.21 (*lighteternal/Llama3-merge-biomed-8b*). Similarly, Gemma models ranged from 0.98 ± 0.02 to 0.75 ± 0.02 in precision. This highlights the impact of specific fine-tuning approaches and the importance of careful model selection.
- **Instruction Following Capabilities:** General-purpose models often outperformed medical-specific models in instruction following, with top performers like *unsloth/Meta-Llama-3.1-8B-Instruct* achieving near-perfect scores (99.0%). This suggests that structured prompt following may be more related to general language understanding than domain-specific training.
- **Robustness and Consistency:** Confidence intervals revealed significant differences in model robustness. Larger general-purpose models often demonstrated consistent performance across data samples, while others, including several medical models, showed wider variability. This shows the importance of considering both peak performance and consistency when evaluating models for clinical applications.

5.2 Clinical Implications

The high performance of several LLMs in IHCA detection suggests significant potential for improving the efficiency and accuracy of critical care event identification in clinical settings. Accurate detection of patients who have suffered IHCA during their hospitalization could greatly enhance QI efforts by ensuring more comprehensive capture of events for review and analysis. Additionally, high-precision models have the potential to reduce the substantial workload of reviewing individual patient records for IHCA for applications in both research and QI. The competitive performance of some open-source models also indicates that effective AI-assisted IHCA detection could be implemented more widely, even in resource-constrained healthcare settings.

Furthermore, models with high recall and strong instruction-following capabilities could be integrated into real-time monitoring systems, enabling faster responses to patients with prior IHCA events. How-

294 ever, the variability in performance and instruction-following capabilities among models illustrates
295 the importance of careful selection and rigorous validation before clinical deployment.

296 5.3 Limitations

297 Despite the promising results, our study has several limitations. The study was conducted on data
298 from a single institution, which may limit the generalizability of findings to other healthcare settings
299 with different patient populations or documentation practices. While we used a substantial dataset,
300 an even larger and more diverse dataset could provide more robust estimates of model performance
301 and generalizability. Although we evaluated a wide range of models, the rapidly evolving field
302 of LLMs means that newer models may have become available since the study’s completion. The
303 inherent biases in medical documentation practices could influence model performance and potentially
304 perpetuate existing healthcare disparities if not carefully addressed. Additionally, our evaluation
305 focused on IHCA detection from text alone, without incorporating other clinical data sources that
306 might be available in real-world settings.

307 5.4 Future Directions

308 We propose several directions for future research. Enhancing model robustness is crucial, with a
309 focus on developing techniques to improve the consistency of open-source models, particularly in
310 handling diverse patient populations and documentation styles. Explainable AI should be a priority,
311 focusing on improving the quality and consistency of model explanations to enhance interpretability,
312 facilitate chart review, and build clinical trust. Multi-modal integration presents an exciting avenue,
313 exploring the integration of text-based LLMs with other data modalities such as vital signs and lab
314 results for more comprehensive IHCA detection.

315 Prospective validation studies will be essential to evaluate the real-world impact of LLM-assisted prior
316 IHCA detection on patient outcomes and quality improvement processes. Ethical AI development
317 must be at the forefront, investigating methods to mitigate potential biases in model training and
318 deployment to ensure equitable performance across diverse patient groups. Lastly, optimizing
319 instruction following capabilities is crucial, with efforts directed towards developing specialized
320 training techniques to enhance the instruction-following capabilities of medical-specific models
321 without compromising their domain expertise.

322 6 Conclusion

323 This evaluation of 52 large language models for IHCA detection demonstrates the significant potential
324 of these AI technologies in this critical care application. While GPT-4o exhibited superior perfor-
325 mance and robustness, several open-source alternatives, both general-purpose and medical-specific,
326 showed promising results that approach clinical viability. The study revealed that fine-tuning on med-
327 ical text does not guarantee superior performance, with some general-purpose models outperforming
328 their medical counterparts in certain aspects, particularly in instruction following. The variability in
329 model performance, both within and across model families, emphasizes the complexity of applying
330 LLMs to specialized medical tasks. It highlights the need for careful model selection and rigorous
331 evaluation that considers not only peak performance but also consistency, robustness, and adherence
332 to clinical instructions.

333 As the field of AI in healthcare continues to evolve rapidly, our findings provide valuable insights
334 for researchers, clinicians, and health system leaders considering the implementation of LLMs for
335 critical care applications. The potential for improving IHCA detection and, consequently, patient
336 outcomes is significant. However, the path to clinical deployment requires addressing challenges in
337 model performance, consistency, and ethical considerations.

338 Future work should focus on enhancing the reliability and interpretability of these models, conducting
339 prospective validations in diverse clinical settings, and developing frameworks for responsible AI
340 integration in healthcare. With continued research and development, LLMs could become powerful
341 tools in the ongoing effort to improve the quality and efficiency of critical care, ultimately contributing
342 to better patient outcomes in IHCA and beyond.

References

- [1] Lars W. Andersen, Mads J. Holmberg, Katherine M. Berg, Michael W. Donnino, and Asger Granfeldt. In-Hospital Cardiac Arrest: A Review. *JAMA*, 321(12):1200–1210, 2019.
- [2] Lingling Wu, Bharat Narasimhan, Kirtipal Bhatia, Kam S. Ho, Chayakrit Krittanawong, Wilbert S. Aronow, Patrick Lam, Salim S. Virani, and Salpy V. Pamboukian. Temporal Trends in Characteristics and Outcomes Associated With In-Hospital Cardiac Arrest: A 20-Year Analysis (1999–2018). *Journal of the American Heart Association*, 10(23), 2021.
- [3] Saket Girotra, Brahmajee K. Nallamothu, John A. Spertus, Yan Li, Harlan M. Krumholz, and Paul S. Chan. Trends in Survival after In-Hospital Cardiac Arrest. *New England Journal of Medicine*, 367(20):1912–1920, 2012.
- [4] Dhaval Kolte, Sahil Khera, Wilbert S. Aronow, Chandrasekar Palaniswamy, Marjan Mujib, Chul Ahn, Sei Iwai, and Gregg C. Fonarow. Regional Variation in the Incidence and Outcomes of In-Hospital Cardiac Arrest in the United States. *Circulation*, 131(16), 2015.
- [5] Sebastian Wiberg, Mathias J. Holmberg, Michael W. Donnino, Jesper Kjaergaard, Christian Hassager, Lise Witten, Katherine M. Berg, Ari Moskowitz, Lars W. Andersen, Anne Grossestreuer, Dana Edelson, Joseph Ornato, Mary Ann Peberdy, Matthew Churpek, Michael Kurz, Monique Anderson Starks, Paul Chan, Saket Girotra, Sarah Perman, and Zachary Goldberger. Age-dependent trends in survival after adult in-hospital cardiac arrest. *Resuscitation*, 151:189–196, 2020.
- [6] Merel Schluep, Sanne E. Hoeks, Marleen Blans, Berber van den Bogaard, Annelies Koopman-van Gemert, Charlotte Kuijs, Christine Hukshorn, Nynke van der Meer, Matthijs Knook, Tim van Melsen, Rene Peters, Peter Perik, Kim Simons, Gerard Spijkers, Willem Vermeijden, Evelien J. Wils, Robert Jan Stolker, and Rik Endeman. Long-term survival and health-related quality of life after in-hospital cardiac arrest. *Resuscitation*, 167:297–306, 2021.
- [7] Ashish R. Panchal, Jason A. Bartos, José G. Cabañas, Michael W. Donnino, Ian R. Drennan, Karen G. Hirsch, and Peter J. Kudenchuk. Part 3: Adult Basic and Advanced Life Support: 2020 American Heart Association Guidelines for Cardiopulmonary Resuscitation and Emergency Cardiovascular Care. *Circulation*, 142(16_suppl_2), 2020.
- [8] Jacob C. Jentzer, Michael Scutella, Francis Pike, Justin Fitzgibbon, Nicole M. Krehel, Lindsey Kowalski, Clifton W. Callaway, Jon C. Rittenberger, Jason C. Reynolds, Gregory W. Barsness, and Cameron DeZfulian. Early coronary angiography and percutaneous coronary intervention are associated with improved outcomes after out of hospital cardiac arrest. *Resuscitation*, 123:15–21, 2018.
- [9] Raina M. Merchant, Robert A. Berg, Lin Yang, Lance B. Becker, Peter W. Groeneveld, and Paul S. Chan. Hospital Variation in Survival After In-hospital Cardiac Arrest. *Journal of the American Heart Association*, 3(1), 2013.
- [10] Rohan Khera, Andrew Humbert, Brian Leroux, Graham Nichol, Peter Kudenchuk, Damon Scales, Andrew Baker, Mike Austin, Craig D. Newgard, Ryan Radecki, Gary M. Vilke, Kelly N. Sawyer, George Sopko, Ahamed H. Idris, Henry Wang, Paul S. Chan, and Michael C. Kurz. Hospital Variation in the Utilization and Implementation of Targeted Temperature Management in Out-of-Hospital Cardiac Arrest. *Circulation: Cardiovascular Quality and Outcomes*, 11(11), 2018.
- [11] Ankur Gupta, Tharanga Perera, Ananda Ganesan, Trish Sullivan, Dennis H. Lau, Kevin C. Roberts-Thomson, Anthony G. Brooks, and Prashanthan Sanders. Complications of catheter ablation of atrial fibrillation: a systematic review. *Circulation: Arrhythmia and Electrophysiology*, 6(6):1082–1088, 2013.
- [12] Kjetil Sunde, Marius Pytte, D. Jacobsen, Arne Mangschau, Lars P. Jensen, Cathrine Smedsrud, Trond Draegni, and Petter A. Steen. Implementation of a standardised treatment protocol for post resuscitation care after out-of-hospital cardiac arrest. *Resuscitation*, 73(1):29–39, 2007.
- [13] Rohan Khera, John A. Spertus, Monique A. Starks, and et al. Administrative Codes for Capturing In-Hospital Cardiac Arrest. *JAMA Cardiology*, 2(11):1275–1277, 2017.

- 394 [14] Christopher DeZorzi, Brendan Boyle, Ahmed Qazi, Kanav Luthra, Rohan Khera, Paul S. Chan,
395 and Saket Girotra. Administrative billing codes for identifying patients with cardiac arrest.
396 *Journal of the American College of Cardiology*, 73(12):1598–1600, 2019.
- 397 [15] Laila Rasmy, Yang Xiang, Ziqian Xie, Cui Tao, and Degui Zhi. Med-BERT: Pretrained
398 contextualized embeddings on large-scale structured electronic health records for disease
399 prediction. *NPJ Digital Medicine*, 4(1):1–13, 2021.
- 400 [16] Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez,
401 Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature Medicine*,
402 29:1930–1940, 2023.
- 403 [17] Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben
404 Yan, Lifang He, Hao Peng, Jianxin Li, Jia Wu, Ziwei Liu, Pengtao Xie, Caiming Xiong, Jian
405 Pei, Philip S. Yu, and Lichao Sun. A comprehensive survey on pretrained foundation models: A
406 history from BERT to ChatGPT, 2023.
- 407 [18] Yikuan Li, Sheng Rao, José Roberto Ayala Solares, Abdelaali Hassaine, Rema Ramakrish-
408 nan, Dexter Canoy, Yajie Zhu, Kazem Rahimi, and Gholamreza Salimi-Khorshidi. BEHRT:
409 Transformer for electronic health records. *Scientific Reports*, 10(1):1–12, 2020.
- 410 [19] Beau Norgeot, Kathleen Muenzen, Thomas A Peterson, Xuancheng Fan, Benjamin S Glicksberg,
411 Gundolf Schenk, Eugenia Rutenberg, Boris Oskotsky, Marina Sirota, Jinoos Yazdany, et al.
412 Protected Health Information filter (Philter): Accurately and securely de-identifying free-text
413 clinical notes. *NPJ Digital Medicine*, 3(1):57, 2020.

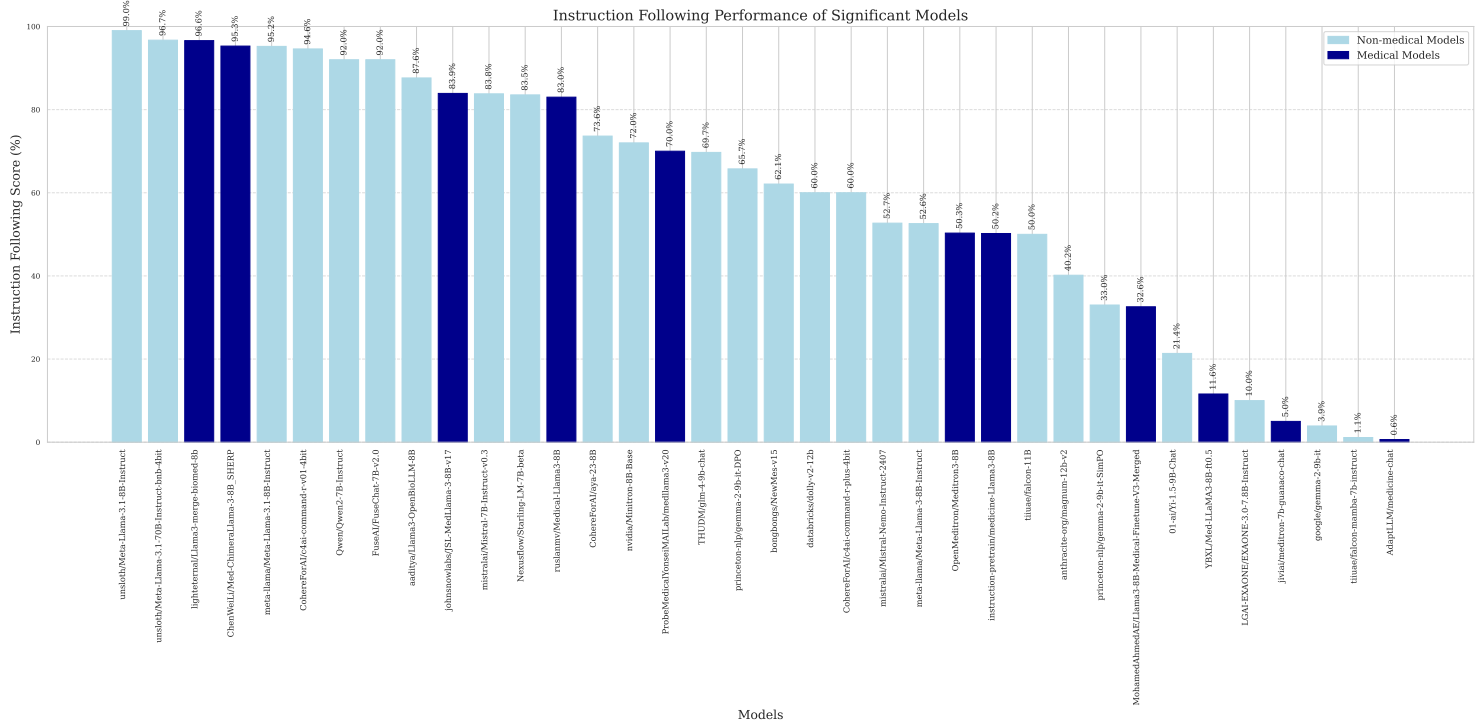


Figure A1: Instruction Following Performance of Significant Models. This bar chart displays the instruction following scores for models that showed statistically significant performance. The scores represent the percentage of responses that correctly followed the given instructions. Models are ordered by their instruction following ability, with medical models highlighted in dark blue.

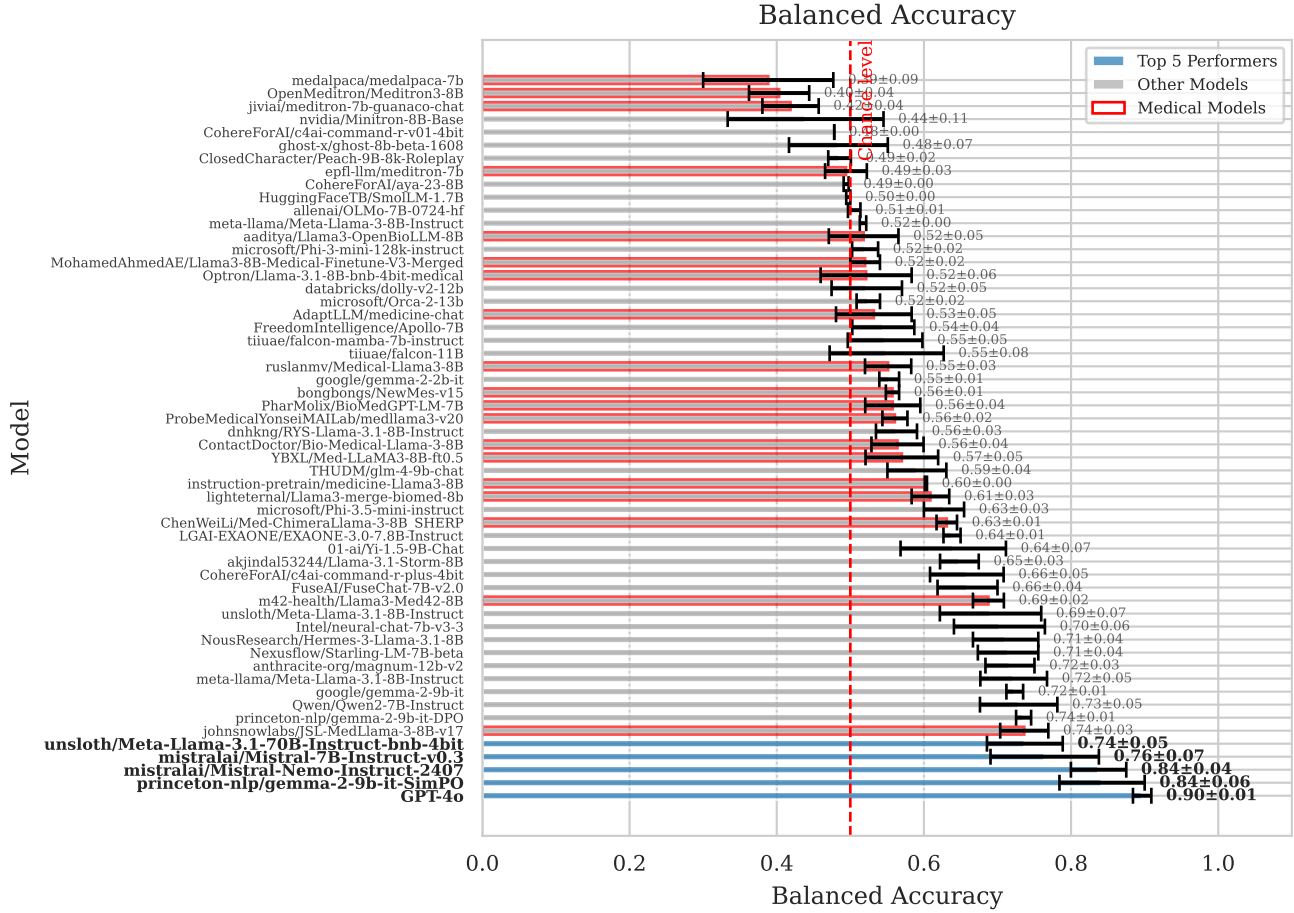


Figure A2: Balanced Accuracy Comparison of Models. This bar chart compares the balanced accuracy scores of various models, including both medical and non-medical models. The top 5 performers are highlighted, and medical models are distinguished by color. Error bars represent confidence intervals calculated through bootstrapping.

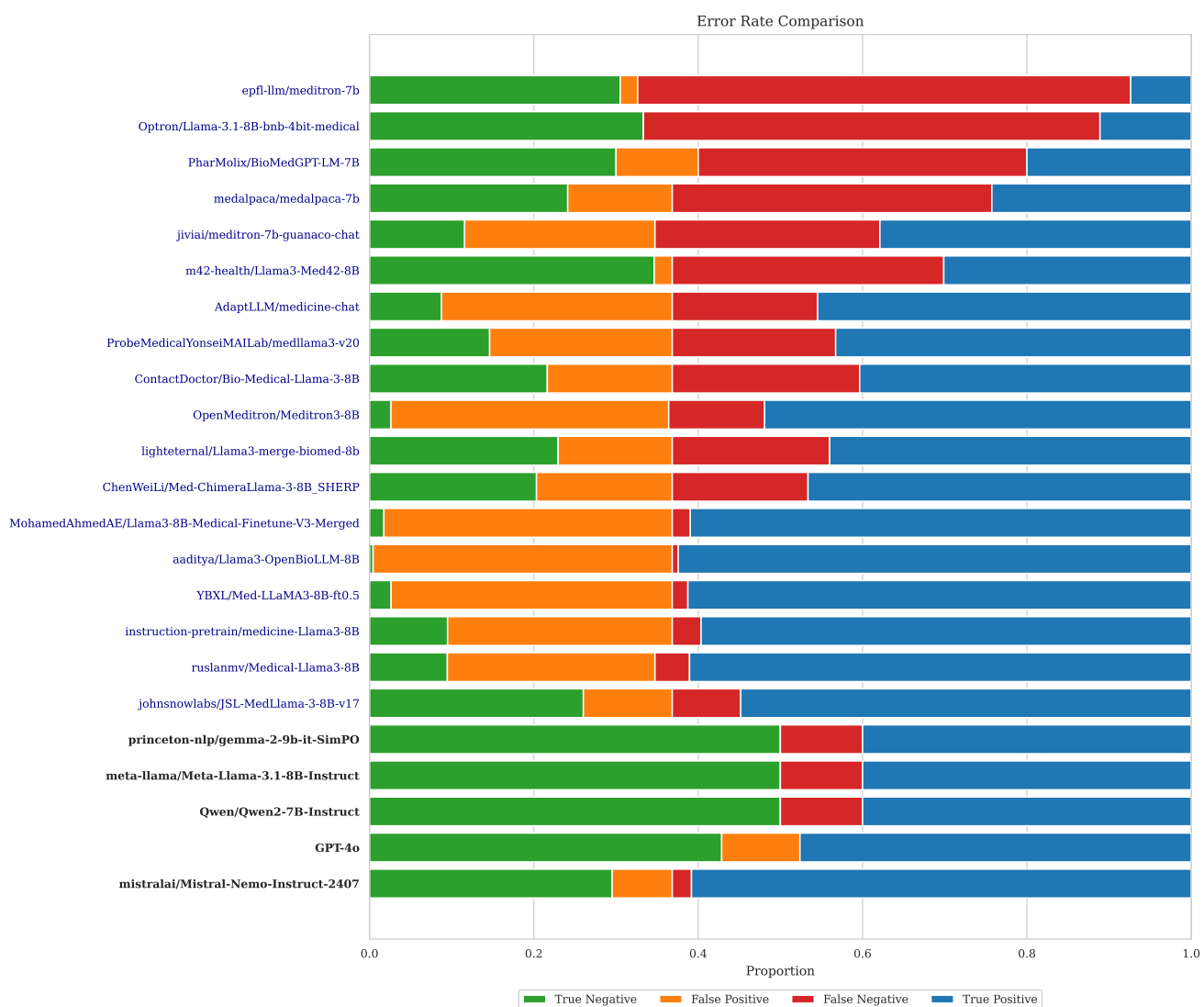


Figure A3: Error Rate Comparison of Top Models. This stacked bar chart illustrates the distribution of true negatives, false positives, false negatives, and true positives for the top-performing models. It provides a comprehensive view of each model's performance in terms of correct and incorrect predictions.

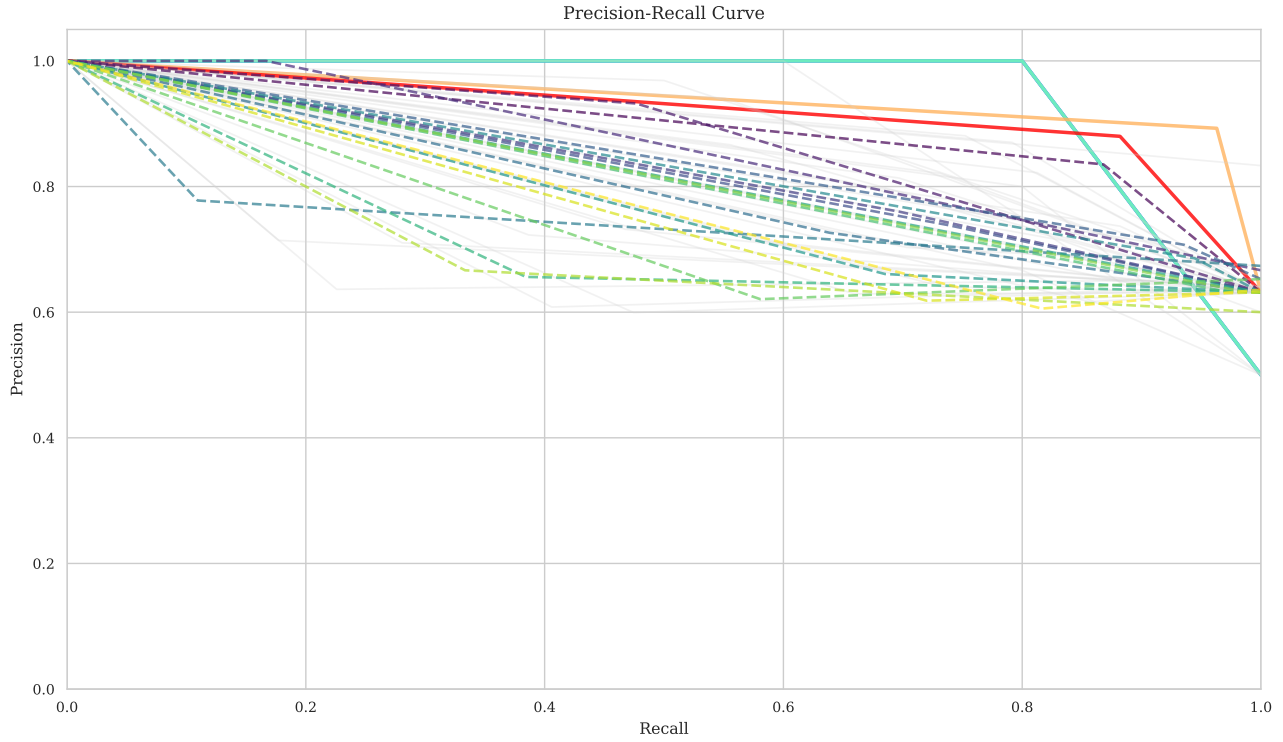
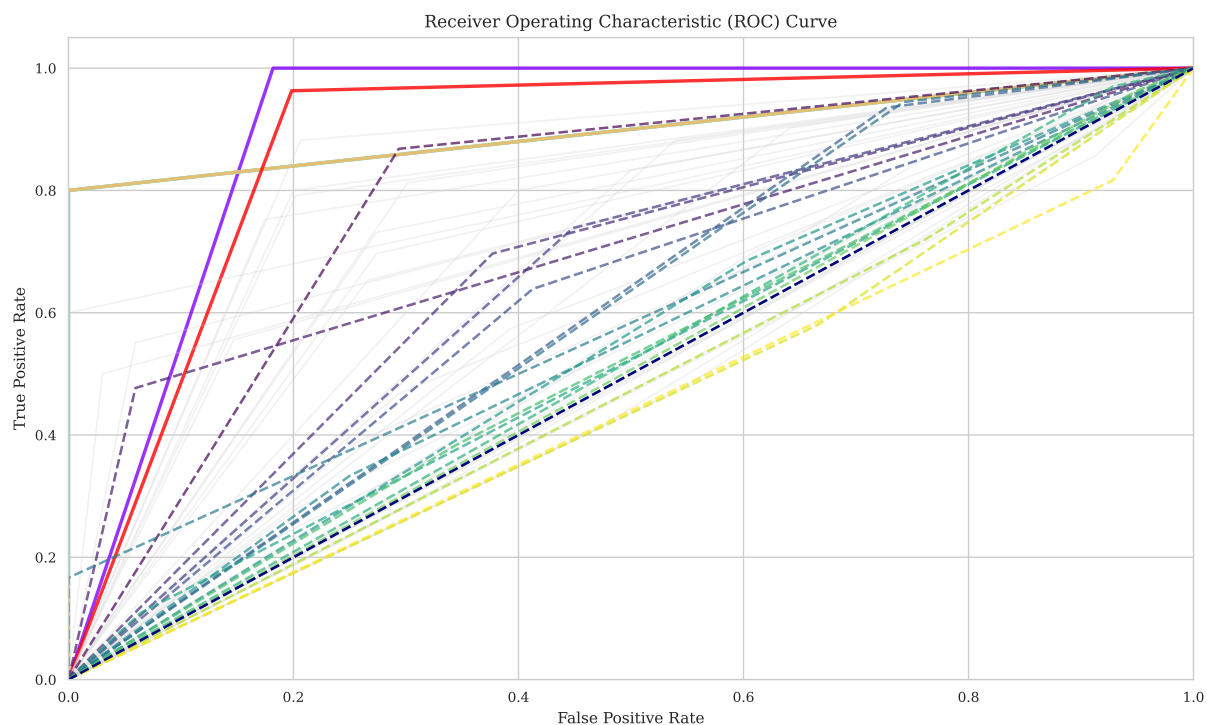


Figure A4: Precision-Recall Curves for Model Comparison. This plot shows the precision-recall curves for various models, allowing for a comparison of their performance across different classification thresholds. The area under each curve (AP) is indicated in the legend, with higher values suggesting better overall performance. Some models are not shown the legend due to space limitations.



GPT-4o (AUC=0.91)	ChenWeiLi/Med-ChimeraLlama-3-8B_SHERP (AUC=0.65)	YBXL/Med-LLaMA3-8B-ft0.5 (AUC=0.52)
Qwen/Qwen2-7B-Instruct (AUC=0.90)	ContactDoctor/Bio-Medical-Llama-3-8B (AUC=0.61)	medalpaca/medalpaca-7b (AUC=0.52)
meta-llama/Meta-Llama-3.1-8B-Instruct (AUC=0.90)	ruslanmv/Medical-Llama3-8B (AUC=0.60)	MohamedAhmedAE/Llama3-8B-Medical-Finetune-V3-Merged (AUC=0.51)
princeton-nlp/gemma-2-9b-it-SimPO (AUC=0.90)	instruction-pretrain/medicine-Llama3-8B (AUC=0.60)	aaditya/Llama3-OpenBioLLM-8B (AUC=0.50)
mistralai/Mistral-Nemo-Instruct-2407 (AUC=0.88)	Optron/Llama-3.1-8B-bnb-4bit-medical (AUC=0.58)	AdaptLLM/medicine-chat (AUC=0.48)
johnsnowlabs/SL-MedLlama-3-8B-v17 (AUC=0.79)	PharMolix/BioMedGPT-LM-7B (AUC=0.54)	jiviai/meditron-7b-guanaco-chat (AUC=0.46)
m42-health/Llama3-Med42-8B (AUC=0.71)	ProbeMedicalYonseiMAILab/medllama3-v20 (AUC=0.54)	OpenMeditron/Meditron3-8B (AUC=0.44)
lighteternal/Llama3-merge-biomed-8b (AUC=0.66)	epfl-llm/meditron-7b (AUC=0.52)	

Figure A5: Receiver Operating Characteristic (ROC) Curves for Model Comparison. This plot displays the ROC curves for various models, illustrating their ability to distinguish between positive and negative cases across different classification thresholds. The area under each curve (AUC) is provided in the legend, with higher values indicating better discriminative power. Some models are not shown on the legend due to space limitations.