

manuscript - Jurnal Dwi.docx

by hurum fatima

Submission date: 18-Dec-2025 08:27AM (UTC+0900)

Submission ID: 2847718419

File name: manuscript_-_Jurnal_Dwi.docx (31.25K)

Word count: 2289

Character count: 15733

1 **High Consistency, Limited Accuracy: Evaluating Large Language Models for Binary
2 Medical Diagnosis**

3 **Authors:** Dwi Anggriani¹, Syaiful Bachri Mustamin¹, Muhammad Atnang¹, Kartini Aprilia
4 Pratiwi Nuzry¹

5 **Affiliation:** ¹Department of Information ²Technology, Institut Sains Teknologi dan Kesehatan
6 'Aisyiyah Kendari, Kendari, Indonesia

7 **Corresponding Author:** Syaiful Bachri Mustamin, syaifulbachri@mail.ugm.ac.id

8 **Preprint:** medRxiv DOI: 10.64898/2025.12.08.25341823

9 **Abstract**

10 **Background:** Large Language Models (LLMs) have demonstrated impressive capabilities in
11 medical knowledge tasks, yet their reliability and consistency in clinical diagnosis remain
12 incompletely characterized.

13 **Objective:** To systematically evaluate the consistency and accuracy of state-of-the-art LLMs in
14 binary medical diagnosis, examining the relationship between reproducibility and diagnostic
15 performance.

16 **Methods:** We evaluated three LLMs (GPT-4o, Gemini-2.0-Flash, Qwen-Plus) on heart disease
17 diagnosis using 100 diverse clinical cases. Each model performed 4 independent assessments per
18 case (1,200 total predictions). We tested two prompt variations and measured intra-model
19 consistency, inter-model agreement, diagnostic accuracy, and prompt sensitivity.

20 **Results:** All models achieved exceptional intra-model consistency (99-100%), with Qwen
21 demonstrating perfect reproducibility. Inter-model agreement was similarly high (98-99%).
22 However, diagnostic accuracy remained at approximately 50%, equivalent to random guessing.
23 Models exhibited a strong bias toward positive diagnosis (49-51 false positives vs 0-1 false
24 negatives per 100 cases). Prompt variation had minimal impact (<3% change), and error patterns
25 were highly systematic, with all models making identical errors on 48-51% of cases.

26 **Conclusions:** Our findings reveal a critical dissociation between consistency and accuracy in
27 LLM medical diagnosis. While LLMs demonstrate remarkable reproducibility, their tendency

28 toward overdiagnosis and limited discriminative accuracy constrain direct clinical utility. Results
29 suggest LLMs may be better suited as supplementary decision-support tools rather than primary
30 diagnostic systems.

31 **Keywords:** Large Language Models, Medical Diagnosis, Consistency Analysis, Heart Disease,
32 Clinical Decision Support, AI Reliability, Reproducibility

33 **Introduction**

34 ⁴ Large Language Models (LLMs) have emerged as promising tools for clinical applications,
35 demonstrating impressive performance on medical licensing examinations and case analysis [1-
36 3]. However, their deployment in clinical settings raises critical questions about reliability and
37 consistency. While traditional diagnostic tools are expected to yield reproducible results, LLMs
38 employ stochastic generation that can lead to varying outputs [4], with documented tendencies
39 toward hallucinations and inconsistent reasoning [9,10]. Reproducibility challenges in AI
40 systems have been extensively documented [13], yet remain incompletely characterized for
41 medical LLMs.

42 Despite growing literature on LLM performance in medical question-answering [5,6], few
43 studies have systematically examined the relationship between consistency and accuracy in
44 diagnostic tasks. Most evaluations focus on single-run accuracy without assessing
45 reproducibility, and the influence of prompt engineering remains incompletely characterized
46 [7,11].

47 This study addresses these gaps through a ⁷ comprehensive evaluation of three state-of-the-art
48 LLMs on binary heart disease diagnosis. We aimed to: (1) quantify intra-model consistency
49 across repeated runs, (2) evaluate inter-model agreement, (3) measure diagnostic accuracy
50 relative to consistency, (4) assess prompt sensitivity, and (5) analyze error patterns to determine
51 if mistakes are random or systematic.

52 **Methods**

53 Dataset and Study Design

54 We utilized the UCI Heart Disease dataset [8], containing 303 patients with 13 clinical
55 parameters: demographics (age, sex), symptoms (chest pain type), vital signs (resting blood

56 pressure), laboratory values (cholesterol, fasting blood sugar), electrocardiography (resting
57 ECG), exercise testing (maximum heart rate, exercise-induced angina, ST depression, ST
58 segment slope), imaging (fluoroscopy vessel count), and thalassemia test results. Binary
59 outcomes indicated the presence or absence of significant coronary stenosis.

60 To ensure diverse representation, we performed k-means clustering (k=2) and selected 50 cases
61 from each cluster, yielding 100 test cases with balanced disease prevalence (51% positive, 49%
62 negative).

63 Models and Experimental Protocol

64 We evaluated three LLMs: GPT-4o (OpenAI), Gemini-2.0-Flash (Google), and Qwen-Plus
65 (Alibaba), accessed via APIs with temperature=0.7. Each model performed 4 independent
66 assessments per case, yielding 1,200 total predictions.

67 We tested two prompt variations: - **Prompt A (“Expert”)**: “You are Dr. CardioExpert, a
68 highly experienced cardiologist...” - **Prompt B (“Neutral”)**: “You are a medical AI
69 assistant trained to provide accurate and balanced diagnostic assessments...”

70 Both prompts provided identical clinical data and parameter definitions, requesting a binary
71 diagnosis (Yes/No) with a 2-3 sentence justification.

72 We implemented a SQLite-based checkpoint system enabling immediate data saving, automatic
73 duplicate prevention, and experiment resumption capability.

74 Outcome Measures

75 **Primary outcomes:** 1. **Intra-model consistency:** Proportion of runs with majority agreement
76 per case 2. **Diagnostic accuracy:** Using majority voting ($\geq 2/4$ runs), we calculated accuracy,
77 sensitivity, specificity, precision, and F1-score 3. **Inter-model agreement:** Pairwise agreement
78 and Cohen’s kappa between models

79 **Secondary outcomes:** 4. **Prompt sensitivity:** Proportion of cases with identical predictions
80 across prompts 5. **Error patterns:** Classification as all-correct, all-wrong, or mixed outcomes

81 Statistical analyses used Python with pandas, scikit-learn, and scipy. Significance was set at
82 $p < 0.05$.

83 **Results**

84 **Intra-Model Consistency: Exceptional Reproducibility**

85 All models demonstrated remarkably high consistency (Table 1, Figure 1). Qwen-Plus achieved
86 perfect consistency (100%) with the expert prompt, never varying across 4 independent runs.
87 GPT-4o and Gemini-2.0-Flash showed 99.0-99.5% average consistency. Notably, 96-100% of
88 cases achieved perfect agreement (4/4 identical predictions), and minimum consistency never fell
89 below 50%.

90 Table 1. Intra-Model Consistency

Model	Prompt	Avg Consistency	Min	Perfect (%)
GPT	Expert	99.25%	50%	98%
GPT	Neutral	99.00%	75%	96%
Gemini	Expert	99.50%	75%	98%
Gemini	Neutral	99.25%	50%	98%
Qwen	Expert	100.00%	100%	100%
Qwen	Neutral	99.75%	75%	99%

91 **Inter-Model Agreement: High Consensus**

92 Models showed 98-100% pairwise agreement, indicating remarkably similar reasoning patterns
93 (Table 2). Three-way agreement (all models concur) occurred in 98-99% of cases. Cohen's
94 kappa values near zero reflected extreme class imbalance (nearly all positive predictions) rather
95 than lack of agreement.

96 Table 2. Inter-Model Agreement

Prompt	GPT-Gemini	GPT-Qwen	Gemini-Qwen	All 3 Agree
Expert	98.0%	100.0%	98.0%	98%
Neutral	100.0%	99.0%	99.0%	99%

97 **Diagnostic Accuracy: Limited Despite High Consistency**

98 Diagnostic accuracy approximated random guessing (48-51%) despite 99-100% consistency
99 (Table 3). Models achieved perfect or near-perfect recall (98-100%) but extremely poor
100 specificity (~0-2%), generating 49-51 false positives versus 0-1 false negatives. This created a
101 consistency-accuracy gap of approximately 50 percentage points.

102 Table 3. Diagnostic Performance

Model	Prompt	Accuracy	Precision	Recall	F1	FP	FN
GPT	Expert	51.0%	51.0%	100%	67.6%	49	0
GPT	Neutral	49.0%	49.0%	100%	65.8%	51	0
Gemini	Expert	51.0%	51.0%	98%	67.1%	48	1
Gemini	Neutral	49.0%	49.0%	100%	65.8%	51	0
Qwen	Expert	51.0%	51.0%	100%	67.6%	49	0
Qwen	Neutral	48.0%	48.5%	98%	64.9%	51	1

103 Representative confusion matrices (Figure 2) showed models predicted “disease present” for
104 nearly all cases, with true negatives ≈ 0.

105 **Prompt Sensitivity: Minimal Impact**

106 Changing from expert to neutral prompt had a minimal effect (Table 4, Figure 3). GPT-4o
107 showed zero sensitivity (100% identical predictions), while Gemini and Qwen changed only 1-3
108 predictions (1-3% of cases). This suggests diagnostic behavior is deeply encoded rather than
109 easily modifiable through prompting.

110 Table 4. Prompt Robustness

Model	Agreement	Changes
GPT	100%	0/100
Gemini	98%	2/100
Qwen	99%	1/100

111 **Error Pattern Analysis: Systematic, Not Random**

112 Errors were highly systematic rather than random (Table 5). In 98-99% of cases, all three models
113 either succeeded together or failed together. Only 1-2% showed model disagreement, indicating
114 shared reasoning patterns or biases.

115 Table 5. Error Consistency

Pattern	Expert	Neutral
All correct	50%	48%
All wrong	48%	51%
Mixed	2%	1%

116 Qualitative analysis revealed models consistently cited elevated cholesterol, abnormal ECG
117 findings, or exercise abnormalities as disease evidence, even when ground truth indicated
118 absence of significant stenosis, suggesting risk factor conflation with diagnostic criteria.

119 **Discussion**

120 **Principal Findings**

121 This study demonstrates a critical dissociation between consistency and accuracy in LLM
122 medical diagnosis: exceptional reproducibility (99-100%) coexists with chance-level accuracy
123 (~50%). This 50-percentage-point gap represents a fundamental challenge for clinical
124 deployment.

125 **The Consistency-Accuracy Paradox**

126 High consistency indicates LLMs reliably apply learned reasoning patterns they are
127 systematically biased rather than randomly erring. This “consistent wrongness” is arguably more
128 concerning than random errors, suggesting fundamental limitations in medical reasoning
129 capabilities [10] rather than simple uncertainty.

130 Several mechanisms may explain this paradox:

131 **Medical Conservatism Bias:** LLMs trained on medical text may have learned that missing
132 disease (false negative) carries greater consequences than over-diagnosis (false positive),
133 encoding a “better safe than sorry” heuristic consistently applied.

134 **Risk Factor Conflation:** Qualitative analysis suggests models conflate risk factors with
135 diagnostic findings. Elevated cholesterol increases disease risk but doesn’t constitute diagnostic
136 evidence of current stenosis. LLMs may struggle to distinguish “high-risk patient” from
137 “disease-positive patient.”

138 **Lack of Discriminative Training:** Unlike supervised models trained explicitly on diagnostic
139 labels, LLMs learn from general medical text, emphasizing disease description more than
140 differential diagnosis, leaving them poorly calibrated for binary classification [12]. This
141 calibration deficit manifests as systematic over-prediction despite high confidence.

142 **Prompt Insensitivity: Deep-Rooted Behavior**

143 The minimal prompt impact (GPT: 0%, others: 1-2%) was unexpected. Refraining from “expert
144 cardiologist” to “neutral assessor” should have reduced conservatism, yet predictions remained
145 nearly identical. This suggests diagnostic behavior is deeply encoded in model weights rather
146 than modifiable through surface-level prompting [11], with important implications for prompt
147 engineering strategies. While prompt patterns can enhance certain LLM behaviors, our findings
148 indicate diagnostic reasoning may be resistant to prompt-level interventions [15].

149 **Inter-Model Agreement: Shared Limitations**

150 The 98-99% inter-model agreement despite different architectures and training procedures
151 suggests observed limitations reflect fundamental challenges in applying LLMs to medical
152 diagnosis rather than model-specific artifacts. Possible explanations include similar training data
153 sources, convergent learning of medical conservatism, shared limitations in processing structured
154 numerical data, and common challenges in threshold-based classification.

155 **Clinical Implications**

156 Current LLMs are not ready for primary diagnostic applications requiring binary classification.
157 The ~50% accuracy is unacceptable clinically and could lead to harmful over-diagnosis,
158 unnecessary testing, and patient anxiety [14]. In a 50% prevalence scenario, deploying these

159 models would result in 98-100% of disease cases correctly identified but only 0-2% of healthy
160 cases correctly identified, causing approximately 50% unnecessary downstream testing. These
161 findings underscore the importance of rigorous evaluation before clinical deployment, as ethical
162 considerations demand careful assessment of potential harms [14].

163 Despite primary diagnosis limitations, LLMs' high consistency and strong negative predictive
164 value suggest potential roles as: (1) second opinion tools where reproducibility builds physician
165 confidence, (2) triage assistants suitable for initial screening where high sensitivity is prioritized,
166 (3) medical education tools providing consistent feedback, and (4) research tools for hypothesis
167 generation.

168 **Technical Implications**

169 Results suggest general-purpose LLMs lack discriminative capabilities for diagnostic
170 classification. Future development should consider: supervised fine-tuning on labeled diagnostic
171 datasets, reinforcement learning from physician-verified diagnoses, calibration techniques for
172 binary classification thresholds, and hybrid architectures combining LLM reasoning with
173 specialized classifiers.

174 **Limitations**

175 Key limitations include: single condition (heart disease may not generalize), binary classification
176 (real diagnosis often involves multi-class assessment), dataset age (1980s diagnostic criteria may
177 differ from current standards), limited sample size (100 cases), structured input only (missing
178 narrative information), three models tested (limited sampling of LLM landscape), API-only
179 access (preventing internal mechanism analysis), and single temperature setting (0.7).

180 **Future Directions**

181 Important future work includes mechanistic studies examining which parameters LLMs
182 prioritize, improvement strategies testing fine-tuning and ensemble approaches, broader
183 evaluations across diverse diagnostic tasks, comparison with human physicians for baseline
184 performance, and theoretical development of consistency-accuracy frameworks.

185 **Conclusions**

186 This study provides rigorous evidence that LLMs achieve exceptional consistency (99-100%) but
187 limited accuracy (~50%) in binary medical diagnosis. This consistency-accuracy dissociation
188 represents a fundamental challenge for clinical deployment. Our findings indicate that high
189 consistency does not guarantee accuracy, diagnostic behavior is resistant to prompt engineering,
190 errors are systematic rather than random, and LLMs show strong positive diagnosis bias.

191 Current general-purpose LLMs are better suited as supplementary decision support rather than
192 primary diagnostic systems. Their exceptional reproducibility is clinically valuable, but limited
193 discriminative ability necessitates human oversight. Future development should prioritize
194 supervised fine-tuning, improved structured data processing, and calibration techniques to
195 address systematic biases.

196 This work contributes to a nuanced understanding of LLM capabilities and limitations in
197 healthcare, informing responsible development and deployment of AI-assisted clinical decision
198 support systems.

199 **Acknowledgments**

200 We thank Institut Sains Teknologi dan Kesehatan 'Aisyiyah Kendari for institutional support.
201 We acknowledge OpenAI, Google, and Alibaba Cloud for providing API access to GPT-4o,
202 Gemini-2.0-Flash, and Qwen-Plus, respectively, through standard commercial services. We
203 thank the UCI Machine Learning Repository and the original data collectors (Janosi, Steinbrunn,
204 Pfisterer, and Detrano) for making the Heart Disease dataset publicly available.

205 **References**

- 206 1. Thirunavukarasu AJ, Ting DSJ, Elangovan K, et al. Large language models in medicine.
207 Nat Med. 2023;29(8):1930-1940. doi:10.1038/s41591-023-02448-8
- 208 2. Lee P, Bubeck S, Petro J. Benefits, limits, and risks of GPT-4 as an AI chatbot for
209 medicine. N Engl J Med. 2023;388(13):1233-1239. doi:10.1056/NEJMsr2214184
- 210 3. Singhal K, Azizi S, Tu T, et al. Large language models encode clinical knowledge.
211 Nature. 2023;620(7972):172-180. doi:10.1038/s41586-023-06291-2

- 212 4. Kung TH, Cheatham M, Medenilla A, et al. Performance of ChatGPT on USMLE:
213 Potential for AI-assisted medical education using large language models. PLOS Digit
214 Health. 2023;2(2):e0000198. doi:10.1371/journal.pdig.0000198
- 215 5. Nori H, Lee YT, Zhang S, et al. Can Generalist Foundation Models Outcompete Special-
216 Purpose Tuning? Case Study in Medicine. arXiv:2311.16452. 2023.
- 217 6. Wang S, Zhao Z, Ouyang X, et al. ChatCAD: Interactive Computer-Aided Diagnosis on
218 Medical Images using Large Language Models. arXiv:2302.07257. 2023.
- 219 7. Moor M, Banerjee O, Abad ZSH, et al. Foundation models for generalist medical
220 artificial intelligence. Nature. 2023;616(7956):259-265. doi:10.1038/s41586-023-05881-
221 4
- 222 8. Janosi A, Steinbrunn W, Pfisterer M, Detrano R. Heart Disease [Dataset]. UCI Machine
223 Learning Repository. 1988. doi:10.24432/C52P4X
- 224 9. Ji Z, Lee N, Frieske R, et al. Survey of hallucination in natural language generation.
225 ACM Comput Surv. 2023;55(12):1-38. doi:10.1145/3571730
- 226 10. Liévin V, Hother CE, Motzfeldt AG, Winther O. Can large language models reason about
227 medical questions? Patterns. 2024;5(3):100943. doi:10.1016/j.patter.2024.100943
- 228 11. White J, Fu Q, Hays S, et al. A prompt pattern catalog to enhance prompt engineering
229 with ChatGPT. arXiv:2302.11382. 2023.
- 230 12. Jiang LY, Liu XC, Nejatbakhsh N, et al. Health system-scale language models are all-
231 purpose prediction engines. Nature. 2023;619(7969):357-362. doi:10.1038/s41586-023-
232 06160-y
- 233 13. McDermott MBA, Wang S, Marinsek N, et al. Reproducibility in machine learning for
234 health research. Sci Transl Med. 2021;13(586):eabb1655.
235 doi:10.1126/scitranslmed.abb1655
- 236 14. Chen IY, Pierson E, Rose S, et al. Ethical machine learning in healthcare. Annu Rev
237 Biomed Data Sci. 2021;4:123-144. doi:10.1146/annurev-biodataisci-092820-114757
- 238 15. Savage T, Nayak A, Gallo R, et al. Diagnostic reasoning prompts reveal the potential for
239 large language model interpretability in medicine. NPJ Digit Med. 2024;7(1):20.
240 doi:10.1038/s41746-024-01010-1



PRIMARY SOURCES

1	Submitted to University of Ulster Student Paper	1 %
2	www.researchsquare.com Internet Source	1 %
3	jurnal.istekaisiyah.id Internet Source	1 %
4	staging-medinform.jmir.org Internet Source	<1 %
5	www.frontiersin.org Internet Source	<1 %
6	Ruslan Ruslan, Meilani Lutfiati Andarini, La Ode Agus Salim, Andi Musdalifah, Maulidiyah Maulidiyah, Muhammad Nurdin. "Synthesis and characterization of electrode Ag-S-TiO ₂ /Ti for enhanced photocatalytic degradation of methylene blue", AIP Publishing, 2023 Publication	<1 %
7	diposit.ub.edu Internet Source	<1 %
8	du.diva-portal.org Internet Source	<1 %