

Benchmarking LLMs on Authentic Cases from Medical Journals

Wanlong Liu^{1,2*}, Junying Chen^{2*}, Yunjin Yang^{2*}, Prayag Tiwari³,
Wenyu Chen^{1†}, Benyou Wang^{2†}

¹University of Electronic Science and Technology of China

²The Chinese University of Hong Kong, Shenzhen

³Halmstad University, Sweden

liuwanlong@std.uestc.edu.cn, yunjinyang69@gmail.com, prayag.tiwari@hh.se,
cwy@uestc.edu.cn, wangbenyou@cuhk.edu.cn

Abstract

In recent years, large language models (LLMs) have demonstrated remarkable capabilities in the medical domain. However, existing medical benchmarks suffer from performance saturation and are predominantly derived from medical exam questions, which fail to reflect the complexity of real-world clinical scenarios. To bridge this gap, we introduce **ClinBench**, a challenging benchmark based on authentic clinical cases sourced from authoritative medical journals. Each question retains the complete patient information and clinical test results from the original case, effectively simulating real-world clinical practice. Additionally, we implement a rigorous human review process involving medical experts to ensure the quality and reliability of the benchmark. ClinBench supports both **textual and multimodal** evaluation formats, covering 11 medical specialties with over 2,000 questions, including a dedicated rare disease track, providing a comprehensive resource for assessing the medical reasoning capabilities of LLMs. We evaluate the performance of over 20 open-source and proprietary LLMs and benchmark them against human medical experts. Our findings reveal that human experts still retain an advantage within their specialized fields, while LLMs demonstrate superior overall performance on a broader range of medical specialties.

1 Introduction

Recent advancements in large language models (LLMs) have demonstrated a remarkable ability to understand and generate medical content, marking significant progress in the medical field (Thirunavukarasu et al., 2023; Liévin et al., 2024; Clusmann et al., 2023; Chen et al., 2024a). Their impressive performance underscores their potential to approach expert-level intelligence.

With the rapid advancement of medical LLMs, existing medical benchmarks lack sufficient challenge and face the issue of performance saturation. For instance, powerful LLMs such as GPT-5 (OpenAI, 2024), Gemini-2.5-Pro (Guo et al., 2025) have achieved accuracy approaching 90% on widely used medical benchmarks like MedQA (Jin et al., 2021), MedMCQA (Pal et al., 2022) and PubMedQA (Jin et al., 2019). To address this limitation, recent works (Zuo et al., 2025; Tang et al., 2025; McDuff et al., 2025) have attempted to introduce more challenging benchmarks by incorporating advanced, expert-level examinations, such as medical licensing tests.

However, **these benchmarks remain predominantly exam-oriented and fail to reflect realistic clinical scenarios**. Real-world clinical scenarios require physicians to integrate medical knowledge with practical experience to navigate complex and multifaceted situations, taking into account patient symptoms, medical history, imaging findings, and various diagnostic test results. In contrast, exam questions typically assess isolated pieces of medical knowledge without adequately simulating complex clinical contexts. Consequently, these questions do not sufficiently capture the intricacies and nuances of real-world clinical scenarios.

To address these limitations, we introduce ClinBench: a challenging and real-world medical benchmark for future medical LLMs. Our ClinBench has four key features: (1) **Real-world clinical scenarios**: The questions are sourced from **authoritative medical journals and based directly on real-world clinical cases**. Each question retains the complete patient information and clinical test results from the original case, effectively simulating real-world clinical practice. (2) **High challenge**: The questions are highly challenging, even for experienced physicians. Each question corresponds to a realistic clinical task, requiring specialized medical knowledge, analytical skills, medical image inter-

*These authors contributed equally to this work.

†Corresponding authors: Wenyu Chen and Benyou Wang.

Dataset	# Size	# Avg Lens	# Avg Option Num	# Saturation	Real-Med Scenario	Rare-Diseases
MedQA (USMLE) (Jin et al., 2021)	1,273	116.6	4	✓	✗	✗
PubMedQA (PQA-L) (Jin et al., 2019)	1,000	14.4	3	✓	✓	✗
MedMCQA (Pal et al., 2022)	4,183	12.8	4	✓	✗	✗
MMLU (Med) (Hendrycks et al.)	1,089	100.1	4	✓	✗	✗
MMLU-Pro (Med) (Wang et al., 2024)	586	166.6	10	✗	✗	✗
MedXpertQA _{Text} (Zuo et al., 2025)	2,450	257.4	8.7	✗	✗	✗
HealthBench (Arora et al., 2025)	5000	-	-	✗	✓	✗
ClinBench _{Text} (Ours)	2,014	462.1	8.4	✗	✓	✓

Table 1: Comparison with existing textual medical benchmarks. **# Saturation** indicates whether the dataset suffers from performance saturation. **Real-Med Scenario** denotes whether the questions reflect real-world clinical settings. **Rare-Diseases** indicates whether the benchmark includes the evaluation of rare diseases.

Dataset	# Size	# Images	# Image Rate	# Avg Lens	# Saturation	Real-Med Scenarios	Rare-Diseases
PMC-VQA (Zhang et al., 2023b)	33,430	29,021	0.9	61.8	✓	✗	✗
OmniMedVQA (Hu et al., 2024)	127,995	118,010	0.9	42.4	✓	✗	✗
GMAI-MMBench (Ye et al., 2024)	21,281	21,180	1.0	49.9	✓	✗	✗
MMM (H & M) (Yue et al., 2024)	1,752	1,994	1.1	83.6	✓	✗	✗
MMM-Pro (H & M) (Yue et al., 2024)	346	431	1.3	107.1	✗	✗	✗
MedXpertQA _{MM} (Zuo et al., 2025)	2,000	2,852	1.4	149.4	✗	✗	✗
ClinBench _{MM} (Ours)	2,014	4,978	2.5	421.7	✗	✓	✓

Table 2: Comparison with existing multimodal medical benchmarks. **# Image Rate** refers to the average number of images included per question. **# Saturation**, **Real-Med Scenario**, and **Rare-Diseases** are consistent with Table 1.

pretation, and extensive clinical experience. (3) **Quality assurance:** For each question, the stem is derived from authoritative medical cases. The golden answer is provided by an expert panel, and each question is thoroughly reviewed and validated by human experts. (4) **Comprehensive Evaluation:** ClinBench provides both textual and multimodal versions, covering 11 medical specialties and encompassing more than 2,000 questions. Additionally, the inclusion of a dedicated rare-disease track further enhances its clinical comprehensiveness, offering a more comprehensive evaluation.

We evaluate over 20 LLMs, including both open-source and proprietary LLMs. Additionally, we engage attending-level human medical experts to answer ClinBench questions, facilitating a comparison between human experts and LLMs. Our key contributions are summarized as follows:

- We propose **ClinBench**, the first medical multiple-choice benchmark focusing on realistic clinical scenarios. ClinBench has both textual and multimodal versions, with questions derived from authentic clinical cases, closely simulating the real-world scenarios.
- ClinBench is built upon authoritative medical journals with rigorous quality assurance processes. Comprehensive human checks and data leakage risk assessments are conducted to

ensure the reliability and quality of the questions.

- We evaluate ClinBench across more than 20 LLMs, providing a comprehensive assessment of the current medical capabilities of existing medical LLMs. Furthermore, through the comparison between human experts and LLMs, we find that human experts still retain an advantage within their specialized fields, while LLMs demonstrate superior overall performance across a broader range of medical specialties.

2 Comparison with Existing Benchmarks

Statistic Comparison. As shown in Tables 1 and 2, traditional text medical benchmarks like MedQA and PubMedQA have short questions with limited options, lacking the challenge of complex, specialized medical tasks. Additionally, MMLU (Hendrycks et al.), MMLU-pro (Wang et al., 2024) and MedXpertQA (Zuo et al., 2025) datasets, mostly sourced from educational exams, fail to accurately represent real clinical tasks. In contrast, our ClinBench_{Text} includes longer, more complex questions with multiple options, all derived from authoritative case journals, offering a better reflection of real clinical scenarios. Moreover, ClinBench_{MM} incorporates more images per question compared to existing multimodal medi-

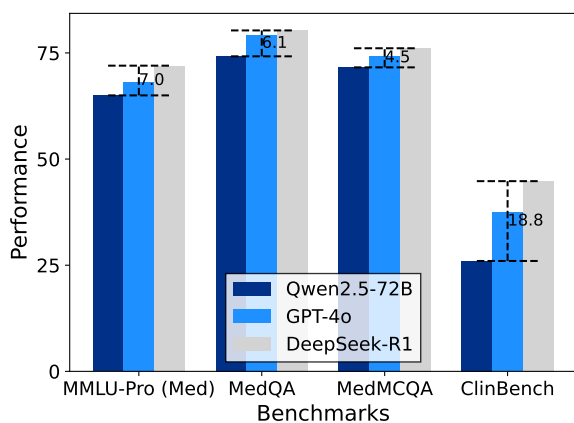


Figure 1: Performance gap among models on different benchmarks. The three exam-oriented benchmarks fail to effectively distinguish capability differences among models, while ClinBench_{Text} provides clearer gap.

cal benchmarks, reflecting the complexity of real-world multimodal medical scenarios. The recent emergence of HealthBench (Arora et al., 2025) has introduced rubric-based evaluation for open-ended medical tasks, where large language models are required to score responses according to predefined criteria. While such a rubric benchmark is suitable for open-ended assessments, it adds additional complexity to the evaluation process. By contrast, our ClinBench adopts a multiple-choice format, enabling more straightforward, consistent, and reliable evaluation.

Discriminative Comparison. As LLMs continue to advance, existing medical benchmarks struggle to effectively evaluate the performance gap among models. As illustrated in Figure 1, compared to benchmarks such as MMLU-Pro and MedQA, our ClinBench demonstrates greater discriminative power, enabling clearer distinctions between model performances. This discriminative nature makes ClinBench as a more suitable benchmark for evaluating and guiding the future development of medical LLMs.

3 The ClinBench Benchmark

3.1 Overview

ClinBench consists of two versions: textual and multimodal. The textual version, ClinBench_{Text}, contains 2,014 multiple-choice questions covering 11 distinct medical specialties. The multimodal version, ClinBench_{MM}, is based on ClinBench_{Text} but partially replaces textual information in question stems with medical images.

The questions in ClinBench are sourced from authoritative medical journals available on PubMed Central¹, a reputable repository of English-language medical case reports. We construct our benchmark based on medical journals for three primary reasons: (1) Medical journals offer comprehensive patient information and detailed clinical test results, thereby closely simulating the real-world diagnostic process. (2) The diagnoses and treatment plans reported in these journals are reviewed and validated by expert medical panels, ensuring the reliability and authority of the reference answers. (3) All patient information in these journals is rigorously anonymized, effectively addressing privacy and ethical concerns. Figure 2 provides two illustrative example from ClinBench_{Text} and ClinBench_{MM}, demonstrating a challenging question that integrates detailed patient information with medical imaging data.

3.2 The Construction of ClinBench_{Text}

Figure 3 illustrates the construction process of ClinBench_{Text}, which consists of three steps: data preprocessing, question stem construction, and candidate options construction.

Data Preprocess. We first convert medical journal PDFs into text format using the MinerU tool². Then, we apply a three-step filtering pipeline to the medical journals: (1) **Filtering for Diagnostic Cases:** We first select medical cases from PubMed Central (32M) using keyword tags. As a result, we obtain the full text of approximately 40K English medical case journals. (2) **Filtering for Complete Cases:** We then apply rule-based filtering to exclude incomplete case reports, retaining only those that contain essential sections: patient information, clinical test results, diagnostic conclusions, and treatment plans. Additionally, we discard cases that lack medical images in the patient information and clinical test results, ensuring that each question includes medical images. After this step, we obtain around 4K high-quality, complete medical cases. (3) **Removing Duplicates:** Finally, we eliminate duplicate or highly similar cases to maintain the dataset’s diversity and quality. A more detailed process is shown in Appendix C.

Question Stem Construction. ClinBench focuses on diagnostic problems because they represent a core, high-stakes clinical challenge that is

¹<https://pmc.ncbi.nlm.nih.gov/>

²<https://mineru.net/>

ClinBench_{Text}

1. General Information
Patient: Male, 66 years old
History: Long-term smoking history
Previous Findings: Two years ago, hospital examination showed multiple linear opacities and small nodules in both lungs.
Current Findings: This year's examination revealed the lesion in the left upper lobe's lingular segment increased ...;

2. Medical Examination Results
Imaging Examination
CT Findings: The CT scan showed multiple smallcystic lucencies with clear borders in both lungs. Multiple solid small nodules with clear borders were noted in the upper lobes of both lung and the left lower lobe. The two larger nodules were located in the left upper lobe and the apical segment of the right upper lobe, measuring approximately 6 mm × 5 mm and 5 mm × 4 mm, ...

Histological Examination
Left Upper Lobe Lingular Segment Larger Nodule:

- Tumor adjacent to the pleura, not penetrating it.
- Tumor cells arranged in bundles or woven patterns, demarcated from surrounding lung tissue.
- Tumor cells were spindle-shaped and relatively uniform in size. Most areas showed no significant atypia, but some showed increased activity and nuclear atypia, with occasional mitosis. Some tumor stroma showed myxoid degeneration.

... (Some content is omitted due to space limitation)

Immunohistochemistry
BCL-2: Negative; **CALP:** Negative; **CD34:** Positive (in vessels); **KL-67:** Positive (5%~10%); **Vimentin:** Positive **CD99, CKP, CR, DESMIN, EMA, HMW, S-100, SMA, TTF-1, WT1, CD117, Dog-1, CD68, HMB45:** Negative **β-Catenin:** Positive (nuclear); **STAT-6, P53 (mutant expression), NF (some cells):** Mixed results; **D-PAS:** Positive

Molecular Pathology
EWSR1-CREB1 FISH Detection: ... (Some content is omitted due to space limitation)

Gene	Mutation Site	Mutation Frequency (%)	Detection Method
HIPK2-YAP1	H1-Y3 fusion	-	RNA-based NGS
TP53	NM_000546.5	33.51	DNA-based NGS
ARID1B	NM_020732.3	19.33	DNA-based NGS

Table 1 NGS Mutation Details

Which of the following diseases is the most likely primary diagnosis for the patient? Please choose the most probable disease.
A: Malignant fibrous histiocytoma
B: Malignant mesothelioma
C: Synovial sarcoma
D: Solitary fibrous tumor of the pleura
 ...
H: Sarcomatoid carcinoma
I: Pulmonary fibromatosis with HIPK2-YAP1 fusion mutation
J: Solitary fibrous tumor

ClinBench_{MM}

1. General Information
Patient: Male, 66 years old
History: Long-term smoking history
Previous Findings: Two years ago, hospital examination showed multiple linear opacities and small nodules in both lungs.
Current Findings: This year's examination revealed the lesion in the left upper lobe's lingular segment increased ...;

2. Medical Examination Results
Imaging Examination
CT Findings: The CT Findings are shown in Figure 1((Lung CT)).

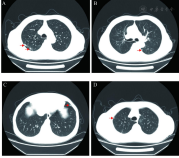


Figure 1: [Lung CT]

Histological Examination
Left Upper Lobe Lingular Segment Larger Nodule:
 Histological images of the larger nodule in the left upper lobe lingular segment are shown in **Figure 2**[Histology] (Sub-figures A-D)

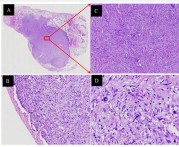


Figure 2: [Histology]

Immunohistochemistry
BCL-2: Negative; **CALP:** Negative; **CD34:** Positive (in vessels); **KL-67:** Positive (5%~10%); **Vimentin:** Positive **CD99, CKP, CR, DESMIN, EMA, HMW, S-100, SMA, TTF-1, WT1, CD117, Dog-1, CD68, HMB45:** Negative **β-Catenin:** Positive (nuclear); **STAT-6, P53 (mutant expression), NF (some cells):** Mixed results; **D-PAS:** Positive

Molecular Pathology
EWSR1-CREB1 FISH Detection:
 ... (Some content is omitted due to space limitation)

Gene	Mutation Site	Mutation Frequency (%)	Detection Method
HIPK2-YAP1	H1-Y3 fusion	-	RNA-based NGS
TP53	NM_000546.5	33.51	DNA-based NGS
ARID1B	NM_020732.3	19.33	DNA-based NGS

Table 1 NGS Mutation Details

Which of the following diseases is the most likely primary diagnosis for the patient? Please choose the most probable disease.
A: Malignant fibrous histiocytoma
B: Malignant mesothelioma
C: Synovial sarcoma
D: Solitary fibrous tumor of the pleura
 ...
H: Sarcomatoid carcinoma
I: Pulmonary fibromatosis with HIPK2-YAP1 fusion mutation
J: Solitary fibrous tumor

Figure 2: Two demos of ClinBench_{Text} and ClinBench_{MM}, respectively. More cases are shown in Appendix F.

amenable to objective, verifiable evaluation. This stands in stark contrast to tasks like treatment planning, whose inherent flexibility and open-ended nature render their assessment less standardized. To construct the questions, we recruited 30 undergraduate medical students who followed our detailed annotation guidelines (see Appendix C). For each case, they extracted the *Patient Information* and *Clinical Test Results* sections from the source material to form the question stem. Finally, every question in the dataset is uniformly framed as: “What is the most likely primary diagnosis for this patient?”

Candidate Options Construction. The construction of candidate options follows four main sub-steps: (1) Medical undergraduate students first extract the final diagnosis from the original journals’ *Diagnosis* sections as the correct answer. Potential alternative diagnoses mentioned in the discussions serve as distractors.³ This step yields one correct option and five distractors. (2) Next, the constructed question stem is submitted to GPT-5, which is instructed to generate the five most plausible diagnoses, ensuring that each option is clear and

³If insufficient alternative diagnoses are mentioned, students are instructed to propose additional plausible alternatives.

precise. (3) Finally, GPT-5 merges the candidate options from the previous two steps, eliminating ambiguous or duplicate entries, thus producing a final, concise set consisting of one correct option and multiple high-quality distractors.

Specialty Categorization. Finally, we classify all questions into 11 medical specialties based on the diseases indicated by the correct answers. The detailed categorization is provided in the Appendix C.3.

3.3 The Construction of ClinBench_{MM}

In real-world clinical scenarios, physicians rely not only on patient textual descriptions but also on medical images from clinical examinations to make informed decisions. To better reflect this multimodal scenario of medical diagnosis, we construct a multimodal version of the ClinBench benchmark dataset. Specifically, we employ undergraduate students to curate ClinBench_{MM} based on questions from ClinBench_{Text} by removing textual content that describes the associated medical images.

To ensure annotation consistency and quality, we provide the annotators with detailed guidelines (see Appendix C) that instruct them on how to identify and remove image-referential text from the question stems. If a question contains a low-quality

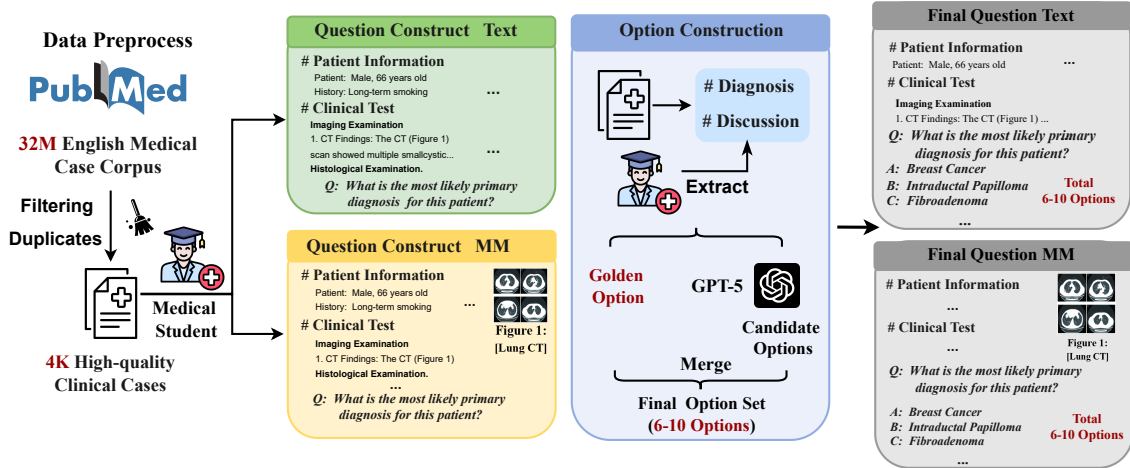


Figure 3: The pipeline for constructing ClinBench. ClinBench_{MM} is built upon ClinBench_{Text}, with the difference being that the question stems in the multimodal version have the textual descriptions of the associated medical images removed.

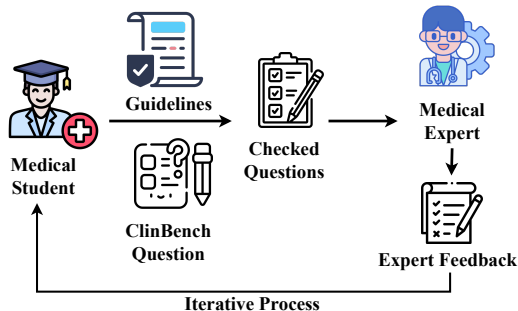


Figure 4: The human check pipeline for ClinBench.

image or one that lacks clinically relevant diagnostic information, it is excluded from the dataset. Following this procedure, we obtain over 2,000 multimodal questions, forming the ClinBench_{MM} benchmark. Compared to ClinBench_{Text}, the question stems in ClinBench_{MM} exclude descriptions related to medical images. ClinBench_{MM} challenges models to accurately extract and reason over visual content from medical images, providing a rigorous assessment of their multimodal understanding and diagnostic capabilities.

3.4 Human Expert Check

We summarize the potential issues that may arise during the question construction process:

(1) **Incorrect Question Stem:** The constructed question stem may **omit critical medical information** present in the original case report. Additionally, when constructing ClinBench_{MM} questions, students are required to manually remove textual descriptions of image-related content from the stem.

This process may further result in the inadvertent loss of important information essential for accurately understanding or answering the question.

(2) **Inappropriate Candidate Options:** During the option merging process, some candidate options may **overlap with the correct answer** or be ambiguously phrased, making it difficult to ensure a clear and unambiguous set of choices.

Two-level Human Check. To ensure the quality of the ClinBench benchmark, we implement a rigorous two-level human check process, and the process is shown in Figure 4. First, medical undergraduate students review the question stems and options to ensure clarity and accuracy. Then, we invite experienced practicing physicians (e.g., attending doctors) to conduct a sampling-based inspection. If any quality issues are identified, they are systematically summarized and fed back to medical students for targeted revision. This iterative process maintains the overall reliability and quality of the benchmark. More detailed process is shown in Appendix C.

4 Experiments

4.1 Experimental Setup

We evaluate all models under a zero-shot prompt setting. Models with fewer than 32B parameters are evaluated locally using 8 * A800 GPUs. During evaluation, we set the temperature to $t = 0.6$ and report the average results over three independent runs. For models larger than 32B parameters, we use the official APIs for evaluation. Detailed

Model	GH	Surg	Neuro	Oph/ENT	DI	Resp	Dent	OG	Ortho	Cardio	Ped	Avg
General Large Language Models												
GPT-4o-2024-11-20	35.7	<u>42.4</u>	37.3	29.8	34.0	40.6	38.3	31.7	35.4	39.8	39.0	37.4
Deepseek-V3.1	44.9	45.5	<u>45.2</u>	39.7	<u>40.9</u>	<u>42.4</u>	55.0	36.6	43.4	<u>46.4</u>	<u>44.1</u>	<u>44.3</u>
Grok-4	38.5	38.4	40.8	<u>33.1</u>	35.5	35.9	<u>40.0</u>	42.7	<u>48.5</u>	40.4	28.8	38.7
Phi-4	35.2	29.3	37.0	16.5	30.5	27.1	35.0	30.5	30.3	33.8	30.5	32.3
Llama-3.1-8B-Instruct	20.8	23.2	27.4	19.8	23.7	18.8	23.3	25.6	16.2	27.3	33.9	23.7
Llama-3.1-70B-Instruct	30.0	29.3	34.3	29.8	30.1	28.2	33.3	28.1	29.3	32.0	28.8	30.9
Qwen-2.5-32B-Instruct	27.1	24.2	31.9	22.3	18.7	25.9	23.3	23.2	27.3	29.0	32.2	26.9
Gemma-3-12B-it	12.4	25.3	21.9	18.2	16.3	15.9	25.0	23.2	22.0	19.2	15.3	17.4
Claude-Opus-4-20250514	<u>43.6</u>	<u>42.3</u>	46.6	<u>32.2</u>	45.8	43.5	<u>33.3</u>	<u>40.2</u>	50.5	47.3	45.8	46.3
Qwen3-32B	30.8	37.3	33.6	28.1	35.5	35.6	28.9	<u>36.7</u>	32.9	34.3	35.0	35.6
Qwen3-8B	25.8	28.3	28.8	27.2	22.1	21.2	25.0	22.0	22.2	29.0	18.6	25.7
General Large Reasoning Models												
GPT-5-0807	<u>52.1</u>	53.5	<u>57.9</u>	<u>44.6</u>	<u>50.7</u>	<u>45.3</u>	41.7	<u>51.2</u>	50.5	<u>52.0</u>	<u>47.5</u>	<u>51.5</u>
Gemini-2.5-pro-Thinking	56.3	<u>49.5</u>	54.1	47.9	54.7	47.1	<u>43.3</u>	53.7	<u>46.7</u>	53.3	49.2	52.6
Deepseek-R1	47.4	46.5	46.9	40.5	40.9	41.8	50.0	34.2	43.4	47.0	42.4	44.8
o3-mini-2025-01-31	36.6	37.4	44.5	33.1	37.4	35.9	31.7	40.4	39.5	39.0	38.8	38.8
o4-mini-2025-04-16	45.2	46.5	51.4	41.3	41.9	42.9	38.3	46.5	45.8	40.7	45.4	45.4
OpenAI-o1-mini	36.2	33.3	34.9	25.6	30.1	31.2	30.0	35.4	32.3	34.1	33.9	33.5
Qwen-3-235b-a22b	40.2	31.3	44.2	30.6	38.9	36.5	<u>43.3</u>	31.7	34.3	27.3	44.1	36.7
Llama-4-maverick	46.2	40.4	46.9	36.4	48.3	40.6	40.0	39.0	44.4	44.6	42.4	44.4
R1-Distill-Qwen-32B	34.0	36.4	38.4	28.1	36.0	31.2	26.7	29.3	26.3	33.5	35.6	33.9
Medical Large Language Models												
Llama-3-8B-UltraMedical	13.4	18.3	19.0	14.9	19.8	15.9	15.0	19.5	14.1	19.2	17.8	17.1
MedGemma-27B-it	19.1	36.7	31.2	35.8	17.9	16.5	18.0	18.1	21.9	11.7	18.8	19.1
Llama-3-70B-UltraMedical	28.5	33.0	34.2	31.7	33.6	31.8	36.7	34.2	31.3	32.0	31.4	31.1
Llama-3-OpenBioLLM-8B	17.1	23.2	24.7	19.8	23.7	20.6	26.7	23.2	20.2	23.7	23.7	22.8
Llama-3-OpenBioLLM-70B	38.0	33.3	37.0	33.1	32.0	38.8	36.7	<u>40.2</u>	37.4	38.9	39.0	37.2
Baichuan-M2-32B	39.9	30.3	37.7	32.2	31.0	30.6	31.7	36.6	28.3	42.0	38.9	37.5
HuatuogPT-o1-7B	24.3	25.3	27.7	21.5	26.6	21.2	26.7	23.2	26.3	34.4	20.3	26.5
HuatuogPT-o1-70B	<u>40.0</u>	<u>37.4</u>	<u>39.2</u>	<u>38.7</u>	<u>38.9</u>	<u>42.4</u>	<u>40.5</u>	32.0	<u>38.0</u>	43.7	<u>41.4</u>	<u>39.2</u>
Lingshu-32B	49.0	43.6	43.8	40.8	39.6	44.4	46.9	40.2	41.4	<u>43.5</u>	43.8	43.4
Human Experts	78.0	28.0	24.0	22.0	-	80.0	12.0	72.0	16.0	90.0	84.0	50.6

Table 3: Performance (accuracy) of various models across medical specialties on ClinBench_{Text}. **Bold** indicates the best performance, and underlined indicates the second best. Specialty abbreviations: GH (Gastroenterology and Hematology), Surg (Surgery), Neuro (Neurosciences), Oph/ENT (Ophthalmology and ENT), DI (Dermatology and Immunology), Resp (Respiratory and Thoracic Medicine), Dent (Dentistry), OG (Obstetrics and Gynecology), Ortho (Orthopedics), Cardio (Cardiovascular and Internal Medicine), Ped (Pediatrics). **Green** indicates questions assessed by human experts within their own specialty, while **red** denotes questions out of their specialty.

evaluation prompts are provided in Appendix D.1.

Evaluated Models We conduct evaluations on a wide range of LLMs and large multimodal models (LMMs) using ClinBench_{Text} and ClinBench_{MM}, respectively. Our benchmark includes both *proprietary* and *open-source models*, and additionally covers advanced large *reasoning* medical models, with a focus on capturing the latest advancements in medical reasoning capabilities. Detailed information of models is shown in Appendix D.2.

4.2 Main Results

Tables 3 and 4 present the main results on ClinBench_{Text} and ClinBench_{MM}, respectively.

Comparison of LLMs. (1) Overall, even the most advanced LLMs, such as Deepseek-R1 and OpenAI-o4-mini, achieve no more than 45% ac-

curacy, indicating relatively poor performance on ClinBench. This underscores the benchmark’s difficulty and its challenge to state-of-the-art LLMs. (2) Generally, large reasoning models outperform non-reasoning counterparts. For instance, Deepseek-R1 and R1-Distill-Qwen-32B show measurable gains over Deepseek-V3.1 and Qwen2.5-32B-Instruct, highlighting the benefits of test-time scaling in complex clinical contexts.

Take-away 1 *Large reasoning models consistently outperform non-reasoning counterparts, underscoring the value of test-time scaling in complex clinical tasks.*

Comparison of LMMs. (1) Consistent with the results observed on ClinBench_{Text}, current LMMs demonstrate relatively low performance on ClinBench_{MM}. (2) On the other hand, LMMs achieve approximately 30 points (which is close to

Model	GH	Surg	Neuro	Oph/ENT	DI	Resp	Dent	OG	Ortho	Cardio	Ped	Avg
Doubao-1.5-Vision-Pro-32k	29.2	34.0	34.0	33.3	33.7	28.2	20.7	34.2	33.7	36.0	28.8	32.4
GPT-4o	31.7	37.1	32.9	32.5	<u>39.6</u>	<u>38.8</u>	<u>37.9</u>	38.5	<u>41.8</u>	37.2	33.9	36.5
Gemini-2.0-Flash	27.7	34.0	39.6	28.3	28.7	31.2	27.6	31.7	29.6	33.5	28.8	31.6
Gemini-2.5-Flash	<u>36.8</u>	<u>39.2</u>	<u>40.5</u>	<u>35.8</u>	38.9	39.5	35.2	<u>39.8</u>	38.6	38.7	<u>37.5</u>	<u>38.1</u>
Gemma-3-27B-It	22.7	27.8	28.5	20.8	24.8	24.7	19.0	25.6	29.6	26.6	27.1	25.3
Internvl3-14B-It	18.0	33.3	40.0	28.6	20.0	23.1	20.0	30.0	26.7	<u>41.9</u>	35.0	29.1
Llama-3.2-11B-Vision-It	20.8	33.3	25.8	21.4	24.0	22.4	26.3	31.8	12.2	21.0	20.0	23.3
Qwen2.5-VL-72B-It	24.6	35.3	29.7	19.2	29.0	25.5	18.2	27.3	34.8	30.9	22.2	27.8
Qwen3-VL-235B-A22B	40.2	41.5	45.8	37.5	43.2	42.6	38.5	44.8	42.3	42.5	39.7	41.1
MedGemma-27B-It	24.7	23.7	<u>30.6</u>	34.8	23.6	25.3	20.0	25.0	<u>34.2</u>	32.2	36.8	27.6
HealthGPT-14B	18.3	21.5	19.8	18.6	20.4	19.7	22.1	19.2	20.8	21.6	18.5	20.0
HuatuoGPT-Vision-7B	18.6	33.0	24.7	20.8	17.7	18.5	17.5	17.5	28.6	20.3	16.9	21.0
HuatuoGPT-Vision-34B	24.6	28.4	27.2	26.5	25.7	24.8	23.5	24.3	28.7	26.9	25.6	26.0
Lingshu-7B	<u>27.8</u>	<u>30.6</u>	29.4	28.2	31.5	<u>30.1</u>	<u>28.7</u>	<u>29.8</u>	<u>27.3</u>	<u>33.2</u>	28.4	<u>29.5</u>
Lingshu-32B	34.2	39.5	36.8	<u>32.4</u>	35.6	33.7	31.8	34.2	37.9	35.8	<u>33.5</u>	35.6
Human Experts	80.0	30.0	28.0	22.0	-	80.0	12.0	74.0	16.0	90.0	84.0	51.6

Table 4: Performance of various models across medical specialties on ClinBench_{MM}. **Bold** indicates the best performance, and underlined indicates the second best.

the average performance on ClinBench_{Text}.) even though the key information of diagnostic images is removed in the question stem. This indicates that LLMs are able to effectively utilize medical image information, which contributes to the resolution of complex clinical problems.

4.3 LLM v.s. Human Experts

In this section, we aim to assess whether the capabilities of LLMs in complex medical scenarios, such as diagnostic reasoning, have reached the level of human experts. To this end, we compare the performance of LLMs with that of human experts on the ClinBench_{MM} dataset.

Setting. We invite experienced physicians from 5 different medical specialties, all of whom hold at least attending-level positions and have extensive clinical diagnostic experience. The detailed information is shown in the Appendix C.5. For the evaluation, we randomly select 50 questions from each medical specialty, totaling 500 questions, and assign each expert questions from two specialties: **one within their own area of expertise and another outside of their specialization.** During the answering process, experts are permitted to consult relevant medical literature and textbooks; however, the use of AI-assisted tools is strictly prohibited.

Medical Insights. From Table 3, we have the following observations: Human experts significantly outperform the strongest current LLMs, Deepseek-R1 and GPT-5, within their own specialized fields. On the other hand, experts perform poorly on questions outside their areas of expertise, whereas LLMs demonstrate relatively stable performance

across all medical specialties, highlighting their stronger generalization capabilities.

4.4 Rare Disease Track

Rare diseases have long posed significant challenges to the medical community (Schieppati et al., 2008; Stoller, 2018), primarily due to limited clinical data, insufficient diagnostic knowledge, and a lack of effective treatments. In this work, we include a dedicated rare-disease subset within our dataset. Medical students carefully select 129 rare-disease cases from various medical specialties, strictly adhering to internationally recognized rare-disease catalogs⁴. This rare disease subset provides a specialized evaluation track for LLMs, which is beneficial for advancing LLMs to overcome the challenges of rare diseases in human medicine.

As shown in Table 5, we observe that the performance of LLMs on rare diseases is significantly lower than on non-rare diseases, highlighting the challenge that rare diseases pose to current LLMs. Additionally, medical LLMs such as Lingshu-32B and HuatuoGPT-o1-7B exhibit relatively smaller performance gaps between rare and non-rare diseases. This may be attributed to their training on more medical texts, including materials related to rare diseases, enabling them to achieve better diagnostic capabilities in rare diseases.

4.5 Reliability

In clinical settings, particularly for diagnosis, the consequences of a single erroneous model response

⁴<https://www.who.int/standards/classifications/frequently-asked-questions/rare-diseases>

Model	Rare Acc	Non-Rare Acc
GPT-5	36.1 ^{-15.4}	51.5
GPT-4o	23.1 ^{-14.3}	37.4
Deepseek-R1	30.3 ^{-14.5}	44.8
Qwen3-8B	11.1 ^{-14.6}	25.7
Llama-3-8B-UltraMedical	9.5 ^{-7.6}	17.1
Huatuogpt-o1-7B	19.5 ^{-7.0}	26.5
Lingshu-32B	36.5 ^{-6.9}	43.4

Table 5: Comparison of model accuracy (%) for rare and non-rare diseases on ClinBench_{text}.

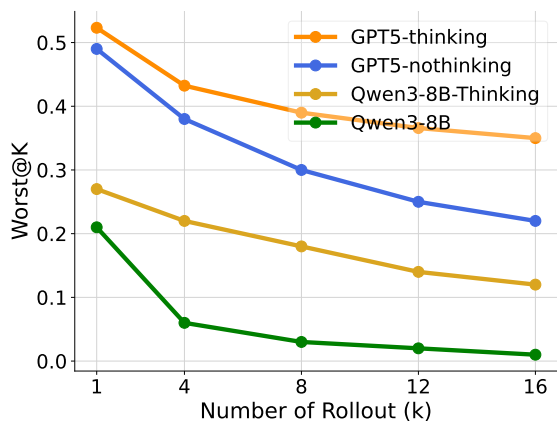


Figure 5: Plot of worst-case ClinBench scores for GPT5 and Qwen3-8B models across different rollouts (k).

can outweigh the benefits of numerous correct ones. To rigorously quantify model reliability under such demanding conditions, we evaluate the *worst@k* performance on the ClinBench benchmark. We define this metric as the minimum score achieved across k independent model responses. A steeper decline in the *worst@k* score as k increases signifies greater performance volatility and, consequently, lower model reliability.

We investigate this by testing both a closed-source model (GPT-5 (OpenAI, 2024)) and an open-source model (Qwen3-8B (Yang et al., 2025)), comparing their performance with and without reasoning-augmented prompting. As depicted in Figure 5, our analysis of *worst@16* performance reveals two key findings. First, for all values of k from 1 to 16, the reasoning-augmented models consistently outperform their non-reasoning models. Second, and more importantly, the reasoning models exhibit significantly greater stability, as evidenced by a much gentler degradation in their *worst@k* scores. This suggests that incorporating reasoning not only improves overall accuracy but also substantially enhances the model’s robustness and trustworthiness for clinical applications.

Take-away 2 *Large reasoning models consistently improve overall accuracy and enhance the model’s robustness and trustworthiness for clinical applications.*

4.6 Error Analysis

In main results, we found that specialized medical LLMs still cannot surpass some powerful general-purpose LLMs. To further identify the problems of current medical LLMs and understand LLMs’ failure modes on ClinBench, we categorized errors into 5 types (Key Information Omission, Clinical Context Misunderstanding, Medical Knowledge Deficiency, Priority Judgment Error, and Temporal Sequence Confusion) and analyzed GPT-5, Deepseek-V3.1, and two medical LLMs Huatuogpt-O1-70B and Baichuan-M2-32B. The detailed results and experimental settings are shown in Appendix D.4.

These results reveal distinct error distribution patterns across models. Medical LLMs (Huatuogpt-o1-70B and Baichuan-M2-32B) exhibit notably high rates of Priority Judgment Error and Temporal Sequence Confusion, indicating stronger medical knowledge but weaker clinical reasoning capabilities. In contrast, general-purpose LLMs (Deepseek-V3.1 and GPT-5) show higher incidence of Medical Knowledge Deficiency but demonstrate better contextual understanding and comprehensive reasoning abilities.

Take-away 3 *Specialized medical LLMs still underperform powerful general LLMs on complex clinical reasoning tasks (ClinBench), primarily because they lack robust contextual understanding and comprehensive reasoning abilities despite having rich domain knowledge.*

5 Conclusion

In this paper, we introduce ClinBench, a medical benchmark specifically designed to simulate real clinical scenarios. This challenging benchmark originates from authoritative medical cases and incorporates detailed patient information and clinical findings, offering a more realistic assessment of LLMs’ medical reasoning. Our comparative analysis of over 20 LLMs against medical experts demonstrates the continued strength of human specialists within their domains, while also highlighting the impressive ability of LLMs to generalize across a wider range of medical knowledge, suggesting their potential to complement and enhance

clinical expertise.

Limitation

Our benchmark currently focuses exclusively on clinical medical diagnosis scenarios, with all data sourced solely from patient case records. While diagnosis represents one of the most critical and challenging tasks within the medical domain, it is important to acknowledge that other scenarios also play vital roles. For example, medical treatment planning, patient monitoring, and healthcare management involve complex decision-making processes that require integration of diverse data types such as longitudinal health records, medical imaging, and real-time sensor data. Furthermore, public health surveillance and preventive care demand models capable of handling population-level data and early risk detection. Therefore, although our benchmark serves as a crucial step towards evaluating AI capabilities in diagnosis, expanding it to encompass these additional healthcare domains is essential for broader applicability and impact.

Acknowledgements

This work was supported by Major Frontier Exploration Program (Grant No. C10120250085) from the Shenzhen Medical Academy of Research and Translation (SMART), Shenzhen Medical Research Fund (B2503005), NSFC grant 72495131, the 1+1+1 CUHK-CUHK(SZ)-GDSTC Joint Collaboration Fund, Guangdong Provincial Key Laboratory of Mathematical Foundations for Artificial Intelligence (2023B1212010001), and the International Science and Technology Cooperation Center, Ministry of Science and Technology of China (under grant 2024YFE0203000).

References

Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, and 1 others. 2024. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*.

Anthropic. 2024. Introducing claude 3.5 sonnet. <https://www.anthropic.com/news/claude-3-5-sonnet>.

Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, and 1 others. 2025. Healthbench: Evaluating large language models

towards improved human health. *arXiv preprint arXiv:2505.08775*.

Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, Jianye Hou, and Benyou Wang. 2024a. Huatuoqpt-o1, towards medical complex reasoning with llms. *arXiv preprint arXiv:2412.18925*.

Junying Chen, Chi Gui, Ruyi Ouyang, Anningzhe Gao, Shunian Chen, Guiming Hardy Chen, Xidong Wang, Ruifei Zhang, Zhenyang Cai, Ke Ji, and 1 others. 2024b. Huatuoqpt-vision, towards injecting medical visual knowledge into multimodal llms at scale. *arXiv preprint arXiv:2406.19280*.

Qifan Chen, Jin Cui, Cindy Duan, Yushuo Han, and Yifei Shi. 2025. Accurate and interpretable postmenstrual age prediction via multimodal large language model. *arXiv preprint arXiv:2508.02525*.

Jan Clusmann, Fiona R Kolbinger, Hannah Sophie Muti, Zunamys I Carrero, Jan-Niklas Eckardt, Narmin Ghaffari Laleh, Chiara Maria Lavinia Löfler, Sophie-Caroline Schwarzkopf, Michaela Unger, Gregory P Veldhuizen, and 1 others. 2023. The future landscape of large language models in medicine. *Communications medicine*, 3(1):141.

Chengfeng Dou, Chong Liu, Fan Yang, Fei Li, Jiyuan Jia, Mingyang Chen, Qiang Ju, Shuai Wang, Shunya Dang, Tianpeng Li, and 1 others. 2025. Baichuanm2: Scaling medical capability with large verifier system. *arXiv preprint arXiv:2509.02208*.

Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. 2022. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*.

Google. 2024. Gemini 2.0 flash. <https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-0-flash>.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.

Yutao Hu, Tianbin Li, Quanfeng Lu, Wenqi Shao, Junjun He, Yu Qiao, and Ping Luo. 2024. Omnimedvqa: A new large-scale comprehensive evaluation

- benchmark for medical lvlm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22170–22183.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146*.
- Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. Biomistral: A collection of open-source pretrained large language models for medical domains. *arXiv preprint arXiv:2402.10373*.
- Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang Liu. 2024. Agent hospital: A simulacrum of hospital with evolvable medical agents. *arXiv preprint arXiv:2405.02957*.
- Valentin Liévin, Christoffer Egeberg Hother, Andreas Geert Motzfeldt, and Ole Winther. 2024. Can large language models reason about medical questions? *Patterns*, 5(3).
- Tianwei Lin, Wenqiao Zhang, Sijing Li, Yuqian Yuan, Binhe Yu, Haoyuan Li, Wangui He, Hao Jiang, Mengze Li, Xiaohui Song, and 1 others. 2025. Healthgpt: A medical large vision-language model for unifying comprehension and generation via heterogeneous knowledge adaptation. *arXiv preprint arXiv:2502.09838*.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, and 1 others. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Daniel McDuff, Mike Schaekermann, Tao Tu, Anil Palepu, Amy Wang, Jake Garrison, Karan Singhal, Yash Sharma, Shekoofeh Azizi, Kavita Kulkarni, and 1 others. 2025. Towards accurate differential diagnosis with large language models. *Nature*, pages 1–7.
- Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, and 1 others. 2023. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *Medicine*, 84(88.3):77–3.
- OpenAI. 2024. Gpt-4o. <https://openai.com/index/hello-gpt-4o/>.
- OpenAI. 2024. Gpt-5. <https://openai.com/zh-Hans-CN/gpt-5/>.
- Ankit Pal and Malaikannan Sankarasubbu. 2024. Openbiollms: Advancing open-source large language models for healthcare and life sciences. <https://huggingface.co/aaditya/OpenBioLLM-Llama3-70B>.
- Ankit Pal, Logesh Kumar Umaphathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Conference on health, inference, and learning*, pages 248–260. PMLR.
- Khaled Saab, Tao Tu, Wei-Hung Weng, Ryutaro Tanno, David Stutz, Ellery Wulczyn, Fan Zhang, Tim Strother, Chunjong Park, Elahe Vedadi, and 1 others. 2024. Capabilities of gemini models in medicine. *arXiv preprint arXiv:2404.18416*.
- Arrigo Schieppati, Jan-Inge Henter, Erica Daina, and Anita Aperia. 2008. Why rare diseases are an important medical and social issue. *The Lancet*, 371(9629):2039–2041.
- Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cían Hughes, Charles Lau, and 1 others. 2025. Medgemma technical report. *arXiv preprint arXiv:2507.05201*.
- James K Stoller. 2018. The challenge of rare diseases. *Chest*, 153(6):1309–1314.
- Xiangru Tang, Daniel Shao, Jiwoong Sohn, Jiapeng Chen, Jiayi Zhang, Jinyu Xiang, Fang Wu, Yilun Zhao, Chenglin Wu, Wenqi Shi, and 1 others. 2025. Medagentsbench: Benchmarking thinking models and agent frameworks for complex medical reasoning. *arXiv preprint arXiv:2503.07459*.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, and 1 others. 2025. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*.
- Qwen Team. 2025. Qwen2.5 models. <https://huggingface.co/Qwen>.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine*, 29(8):1930–1940.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others. 2024. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*.
- Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. 2024. Benchmarking benchmark leakage in large language models. *arXiv preprint arXiv:2404.18824*.

- Weiwen Xu, Hou Pong Chan, Long Li, Mahani Aljunied, Ruifeng Yuan, Jianyu Wang, Chenghao Xiao, Guizhen Chen, Chaoqun Liu, Zhaodonghui Li, and 1 others. 2025. Lingshu: A generalist foundation model for unified multimodal medical understanding and reasoning. *arXiv preprint arXiv:2506.07044*.
- Zhonghao Yan, Muxi Diao, Yuxuan Yang, Jiayuan Xu, Kaizhou Zhang, Ruoyan Jing, Lele Yang, Yanxi Liu, Kongming Liang, and Zhanyu Ma. 2025. Medreasoner: Reinforcement learning drives reasoning grounding from clinical thought to pixel-level precision. *arXiv preprint arXiv:2508.08177*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Jin Ye, Guoan Wang, Yanjun Li, Zhongying Deng, Wei Li, Tianbin Li, Haodong Duan, Ziyang Huang, Yanzhou Su, Benyou Wang, and 1 others. 2024. Gmai-mmbench: A comprehensive multimodal evaluation benchmark towards general medical ai. *Advances in Neural Information Processing Systems*, 37:94327–94427.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.
- Hongbo Zhang, Junying Chen, Feng Jiang, Fei Yu, Zhihong Chen, Jianquan Li, Guiming Chen, Xiangbo Wu, Zhiyi Zhang, Qingying Xiao, and 1 others. 2023a. Huatuogpt, towards taming language model to be a doctor. *arXiv preprint arXiv:2305.15075*.
- Kaiyan Zhang, Sihang Zeng, Ermo Hua, Ning Ding, Zhang-Ren Chen, Zhiyuan Ma, Haoxin Li, Ganqu Cui, Bqing Qi, Xuekai Zhu, and 1 others. 2024. Ultramedical: Building specialized generalists in biomedicine. *arXiv preprint arXiv:2406.03949*.
- Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023b. Pmc-vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Yuchen Duan, Hao Tian, Weijie Su, Jie Shao, and 1 others. 2025. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*.
- Yuxin Zuo, Shang Qu, Yifei Li, Zhangren Chen, Xuekai Zhu, Ermo Hua, Kaiyan Zhang, Ning Ding, and Bowen Zhou. 2025. Medxpertqa: Benchmarking expert-level medical reasoning and understanding. *arXiv preprint arXiv:2501.18362*.

A Ethics Statement

All data used in this work were obtained exclusively from **freely and publicly accessible** sources. We have carefully curated the dataset by retaining only a small and representative subset of the original data. To ensure compliance with **U.S. fair use laws**, all questions underwent rephrasing, and answer options were shuffled to prevent any direct replication. Importantly, the dataset **does not contain any personal, sensitive, or identifiable information**, strictly avoiding any privacy violations or ethical concerns related to personal data. No content involving individual identities, medical records, or confidential information has been included. To mitigate the potential data leakage risks, we refrain from releasing the data sources and request that you **do not share any example of benchmark online**, whether in plain text, image, or any other format.

B Related Work

Medical LLMs. The success of LLMs has sparked interest in creating medical-specific models, leading to the emergence of numerous powerful medical LLMs (Nori et al., 2023; Saab et al., 2024; Li et al., 2024). For example, UltraMedical collections (Zhang et al., 2024) refine LLaMA-3 models with premium datasets, achieving top benchmark performance and advancing online preference learning. Baichuan-M2 (Dou et al., 2025) enhanced the capability of language models in medical diagnosis tasks through a large-scale and realistic dynamic validation system. BioMistral (Labrak et al., 2024), an open-source model pre-trained on PubMed Central, excels in English medical QA tasks. HuatuoGPT series (Chen et al., 2024a,b; Zhang et al., 2023a), trained on high-quality medical data, significantly enhances automated capabilities in diagnosis, triage, and medical imaging, providing valuable support for clinical decision-making and patient care. Building on this trend, recent models (Yan et al., 2025; Chen et al., 2025) such as Lingshu (Xu et al., 2025) and HealthGPT (Lin et al., 2025) further advance the field by improving medical reasoning, multimodal understanding, and have demonstrated strong empirical performance across a range of medical benchmarks.

Medical Benchmarks. With the advancement of medical LLMs, corresponding benchmarks have

also evolved. Early datasets such as MedQA (Jin et al., 2021) and MedMCQA (Pal et al., 2022) focus on multiple-choice questions from USMLE and Indian medical exams, assessing models’ factual knowledge across various medical domains. PubMedQA (Jin et al., 2019) emphasizes reasoning over biomedical research abstracts. On the other hand, the emergence of MMLU-Pro (Wang et al., 2024) and GPQA (Frantar et al., 2022) benchmarks provides a more effective evaluation of long-chain reasoning models like OpenAI-o1. Additionally, MedXpert-QA (Zuo et al., 2025) introduces expert-level questions derived from advanced medical exams, significantly increasing benchmark difficulty. However, these benchmarks remain predominantly exam-oriented and fail to capture the complexity of real-world clinical scenarios. In this paper, we focus on realistic medical scenarios by constructing a benchmark based on real-world clinical case questions. The recent emergence of HealthBench (Arora et al., 2025) has introduced rubric-based evaluation for open-ended medical tasks, where large language models are required to score responses according to predefined criteria. While such a rubric benchmark is suitable for open-ended assessments, it adds additional complexity to the evaluation process. By contrast, our ClinBench adopts a multiple-choice format, enabling more straightforward, consistent, and reliable evaluation.

C Detailed Process for ClinBench Construction

In this section, we provide a detailed description of the construction process of the **ClinBench** benchmark.

C.1 Dataset Construction

We recruited 30 medical undergraduate students to assist in the dataset construction process. These students, all majoring in medicine, possess solid medical foundations. We provided them with comprehensive annotation guidelines, instructing them to carefully construct and verify question–answer pairs. Specifically, the guidelines detailed two core tasks: (1) constructing the ClinBench_{Text} questions (see Table 12); and (2) constructing the multimodal ClinBench_{MM} questions (see Table 13). Following these guidelines rigorously, the students successfully constructed a total of 2,014 high-quality questions.

C.2 Human Verification

Our human verification process involves two critical steps. First, medical undergraduate students carefully review the question stems and candidate options to ensure clarity and medical accuracy, following detailed guidelines as illustrated in Table 14. Subsequently, experienced practicing physicians (e.g., attending doctors) perform a sampling-based inspection of the reviewed questions. Any identified quality issues are systematically summarized and communicated back to the medical students for targeted revision, guided by the criteria provided in Table 15. This iterative feedback and revision mechanism ensures the overall reliability and high quality of the benchmark.

C.3 Data Specialty Distribution

We analyze the distribution of medical specialties covered by the ClinBench_{Text} and ClinBench_{MM} datasets, as illustrated in Figure 6. As shown, each medical specialty includes at least 50 questions, ensuring sufficient coverage for comprehensive evaluation.

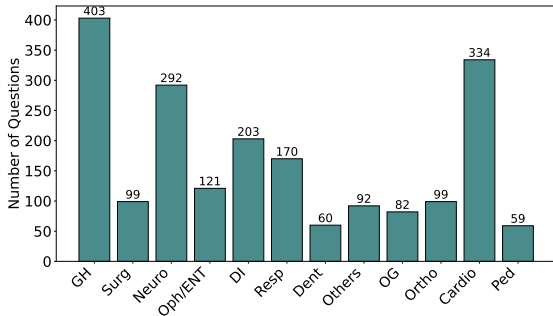


Figure 6: The pipeline for constructing ClinBench. Specialty abbreviations: GH (Gastroenterology and Hematology), Surg (Surgery), Neuro (Neurosciences), Oph/ENT (Ophthalmology and ENT), DI (Dermatology and Immunology), Resp (Respiratory and Thoracic Medicine), Dent (Dentistry), OG (Obstetrics and Gynecology), Ortho (Orthopedics), Cardio (Cardiovascular and Internal Medicine), Ped (Pediatrics).

C.4 Computational Cost for Dataset Construction and Model Evaluation

To facilitate reproducibility and transparency of our experiments, we report the detailed computational costs for both dataset construction and model evaluation. All API-based costs were incurred via official service channels, ensuring compliance and consistency. GPU hours refer to actual runtime used for local inference on NVIDIA A100 80GB GPUs unless otherwise specified.

Model	GPU Hours (h)	API Quotas (\$)
Deepseek-V3.1	-	2.090
Llama-4-maverick	-	1.098
GPT-5-0807	-	2.160
Gemini-2.5-Pro	-	4.230
o3-mini-2025-01-31	-	5.460
HuatuoGPT-o1-7B	0.42	-
Qwen3-8B	0.37	-

Table 6: Model Evaluation Costs (All Questions) on ClinBench_{Text}

Dataset Component	API Cost (\$)
ClinBench (Total)	720
ClinBench _{Text}	620
ClinBench _{MM}	100

Table 7: Dataset Construction Costs

The dataset construction phase involved substantial computational resources, primarily for data collection, preprocessing, annotation validation, and quality assurance procedures. For the multimodal dataset construction, additional costs were incurred for image processing, medical image anonymization, and cross-modal alignment verification. The annotation process required multiple rounds of expert review and inter-annotator agreement validation to ensure clinical accuracy and reliability.

Model evaluation costs encompass inference time across all benchmark tasks, including both text-only and multimodal assessments. The evaluation costs for ClinBench_{MM} are nearly identical to those for ClinBench_{Text} in terms of API calls and computational overhead, as the primary difference lies in the input modality rather than the evaluation pipeline complexity. However, multimodal models generally require slightly longer inference times due to the additional visual encoding step. We include both wall-clock time and normalized GPU hours to provide comprehensive cost metrics for future researchers planning similar evaluation campaigns. All cost estimates are provided in Table 6 and Table 7 and represent the resources required for complete evaluation across our benchmark suite.

C.5 Expert Evaluation Details

Our experts are experienced attending physicians from Longgang People’s Hospital in Shenzhen, Guangdong Province, China, specializing in Gastroenterology and Hematology, Respiratory and Thoracic Medicine, Obstetrics and Gynecology,

Cardiovascular and Internal Medicine, and Pediatrics. For the evaluation, we randomly select 50 questions from each of these medical specialties, totaling 500 questions, and assign each expert questions from two specialties: **one within their own area of expertise and another outside of their specialization**. During the answering process, experts are permitted to consult relevant medical literature and textbooks; however, the use of AI-assisted tools is strictly prohibited.

For the evaluation, we developed an online assessment platform based on the streamlit framework, allowing physicians to answer questions following the provided guidelines. A screenshot of the evaluation platform is shown in Figure 10.

D Experimental Details and Additional Results

D.1 Evaluation Details

We evaluate all models under a zero-shot prompt setting. For models such as GPT-5-0807, GPT-4o-2024-1120, Deepseek-V3.1, Gemini-2.5-Pro-0617, Deepseek-R1-0528, Grok-4, doubao-1-5-pro-32k-250115, o1-mini, o3-mini, o4-mini, and Qwen3-235B-a22b, we utilize the official APIs provided by the official. The remaining models are evaluated on a local setup consisting of 8 A800 80GB GPUs. During evaluation, we set the temperature to $t = 0.6$ and report the average results over three independent runs. The prompts used are illustrated in Figure 7 and Figure 8.

D.2 Evaluated Models

Large Language Models: We evaluate a wide range of LLMs on ClinBench_{Text}. The general LLMs include GPT-4o-2024-11-20 (OpenAI, 2024), Deepseek-V3.1 (Liu et al., 2024), Grok-4, Phi-4 (Abdin et al., 2024), LLaMA3.1-70B-Instruct, LLaMA3.1-8B-Instruct, Qwen2.5-32B-Instruct, Qwen2.5-7B-Instruct and Qwen3-8B (Yang et al., 2025). Gemma-3-12B-it (Team et al., 2025) We also include general reasoning models such as DeepSeek-R1, DeepSeek-R1-Distill-Qwen-32B (Guo et al., 2025), Qwen-3-235b-a22b (Yang et al., 2025), OpenAI-o1, OpenAI-o3-mini, OpenAI-o3-mini⁵, Claude-opus-4-20250514 (Anthropic, 2024), Grok-4, GPT-5-0807 (OpenAI, 2024),

⁵<https://openai.com/o1/>

Llama-4-maverick⁶. In the medical domain, we assess medical-domain LLMs including OpenBioLLM-8B (Pal and Sankarasubbu, 2024), and UltraMedical-8B (Zhang et al., 2024), Baichuan-M2-32B (Dou et al., 2025), Lingshu-32B (Xu et al., 2025), as well as the medical reasoning model HuatuoGPT-o1 (Chen et al., 2024a).

Large Multimodal Models (LMMs): We evaluate several large multimodal models on ClinBench_{MM}, including proprietary models such as GPT-4o-2024-11-20 (OpenAI, 2024), Gemini-2.0-Flash (Google, 2024), Gemini-2.5-Flash (Google, 2024), as well as open-source models such as Doubao-1.5-Vision-Pro-32k, InternV3-14B (Zhu et al., 2025), Qwen2.5-VL-72B (Team, 2025) and Llama-3.2-11B-Vision (Grattafiori et al., 2024). For medical llms, we evaluate the HuatuoGPT-Vision (Chen et al., 2024b), Lingshu-32B (Xu et al., 2025), MedGemma-27B-Instruct (Sellergren et al., 2025) and HealthGPT-14B (Lin et al., 2025).

D.3 Data Leakage Analysis

To evaluate the potential risk of data leakage in the ClinBench benchmark, we follow work (Xu et al., 2024) by employing perplexity (PPL) and N-gram-based metrics (ROUGE-L and edit distance similarity) as evaluation criteria. Specifically, we concatenate the original question with a prompt such as "Answer:" as input and calculate the model's perplexity on the generated output. Additionally, to assess the similarity between the model's generated rationale and the reference explanations we collected, we compute both ROUGE-L scores and edit distance similarity. We hypothesize that if "data leakage" occurred, the models would have previously learned specific reference texts, resulting in lower perplexity (PPL), higher textual overlap (Rouge-L), and higher edit distance similarity (EDS). Conversely, if no leakage occurred, we would expect the opposite pattern.

We evaluate two models, including LLaMA-3.1-70B-Instruct, and Qwen2.5-72B-Instruct, on ClinBench. As shown in Table 8:

- **PPL (Perplexity):** All models exhibit extremely high PPL values (e.g., LLaMA-3.1-

⁶<https://ai.meta.com/blog/llama-4-multimodal-intelligence/>

Models	PPL \uparrow	Rouge-L \downarrow	EDS \downarrow
Qwen2.5-72B-Instruct	1.12E+7	0.1794	0.2493
LLaMA-3.1-70B-Instruct	9.73E+6	0.1597	0.2285

Table 8: Data leakage analysis results on different models. PPL denotes Perplexity, and EDS stands for Edit Distance Similarity.

Model	Key Information Omission	Clinical Context Misunderstanding	Medical Knowledge Deficiency	Priority Judgment Error	Temporal Sequence Confusion
Deepseek-V3.1	112	337	438	168	67
GPT-5	117	264	303	195	98
Huatuogpt-o1-70B	196	327	222	308	182
Baichuan-M2-32B	203	305	210	368	198

Table 9: Error Analysis Across Different Models (Number of Cases). The highest number of errors in each category is shown in **bold**.

70B-Instruct: 9.73×10^6), indicating substantial difficulty in predicting the reference texts and lack of prior exposure.

- **Rouge-L (Textual Overlap):** Rouge-L scores are uniformly low (maximum 0.1794), reflecting minimal textual overlap with reference texts and no evidence of memorization.
- **EDS (Edit Distance Similarity):** EDS values are similarly low (maximum 0.2493), consistent with minimal content overlap and absence of memorization.

Our analysis finds no evidence of data leakage. This is primarily due to: (1) the questions in ClinBench being sourced from professional medical case reports absent from the models’ training corpora; and (2) the inherent complexity and rich clinical context making memorization or direct reproduction from prior exposure highly unlikely. These findings collectively contradict the data leakage hypothesis, confirming that ClinBench poses minimal risk of leakage and supports its validity and robustness as a medical benchmark.

D.4 Error Analysis

To better understand common failure modes of LLMs on ClinBench, we conducted an in-depth error analysis. We manually analyzed a large number of incorrect cases from different models’ responses and categorized errors into the following five types, each defined as follows:

1. **Key Information Omission:** The model neglects critical clinical information in the question (e.g., important symptoms, signs, labo-

ratory results, or medical history), leading to diagnostic errors.

2. **Clinical Context Misunderstanding:** The response reflects an incorrect interpretation of the clinical scenario or question context, leading to inappropriate reasoning paths.
3. **Medical Knowledge Deficiency:** The model correctly interprets the question but produces an incorrect answer due to insufficient or inaccurate medical knowledge.
4. **Priority Judgment Error:** The model fails to identify the most primary, urgent, or contextually appropriate diagnosis among multiple plausible options.
5. **Temporal Sequence Confusion:** The model misunderstands the timeline of disease progression, acute versus chronic course, or the sequence of clinical events.

We conducted an error analysis experiment on ClinBench_{text} for Deepseek-V3.1, GPT-5, Huatuogpt-o1-70B and Baichuan-M2-32B. Specifically, we used GPT-5 to classify error cases according to the five error types defined above. The detailed results are shown in Table 9 and the prompt is shown in Figure 9.

These results reveal distinct error distribution patterns across models. For example, Huatuogpt-o1-70B and Baichuan-M2-32B exhibits a notably high rate of Priority Judgment Error (318 cases) and Temporal Sequence Confusion (172 cases), whereas Deepseek-V3.1 shows the highest incidence of Medical Knowledge Deficiency (438 cases) and Clinical Context Misunderstanding (337

cases). These findings indicate that medical-domain LLMs demonstrate better medical knowledge injection, but their understanding capabilities and fundamental reasoning abilities exhibit deficiencies. In contrast, general-purpose LLMs may lack specialized medical expertise, but possess relatively stronger comprehensive capabilities. Such findings highlight the varying weaknesses of LLMs in complex clinical reasoning tasks.

E Addressing Clinical Metrics: Beyond Simple Accuracy

E.1 Motivation

While accuracy serves as a fundamental evaluation metric, it lacks the clinical nuance necessary for comprehensive assessment of medical AI systems. In clinical practice, the ability to correctly identify positive cases (sensitivity) and accurately exclude negative cases (specificity) are critical performance indicators that directly impact patient safety and clinical decision-making. To address this limitation and provide a more clinically meaningful evaluation framework, we introduce sensitivity and specificity metrics validated against expert-annotated standards.

E.2 Metric Definitions

E.2.1 Sensitivity

Sensitivity measures the model’s ability to correctly identify all correct answers in clinical question-answering scenarios, analogous to correctly diagnosing a patient as "having the disease" in traditional clinical settings.

Formula:

$$\text{Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

where:

- **TP (True Positive):** The model correctly selects the correct answer (e.g., correctly identifying option A as the correct diagnosis).
- **FN (False Negative):** The model fails to select the correct answer (e.g., missing the correct option A).

E.2.2 Specificity

Specificity measures the model’s ability to correctly exclude incorrect answers, analogous to correctly identifying a patient as "not having the disease" in clinical practice.

Formula:

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}} \quad (1)$$

where:

- **TN (True Negative):** The model correctly excludes incorrect answers (e.g., correctly rejecting options B, C, D as incorrect).
- **FP (False Positive):** The model incorrectly selects an incorrect answer as correct (e.g., wrongly choosing option B).

E.3 Results with Clinical Metrics

We re-evaluate all models on ClinBench using these clinically-relevant metrics to provide deeper insights into model performance characteristics. For text-based models, we evaluate Deepseek-V3.1, HuatuoGPT-o1-70B, GPT-4o, Gemini 2.5 Pro, GPT-5, and Qwen3-32B. For multimodal models, we evaluate GPT-4o, Qwen2.5-VL-72B, Qwen3-VL-235B-A22B, and Lingshu-32B.

Table 10 presents the sensitivity and specificity results for text-based language models on ClinBench_{Text}. Table 11 shows the sensitivity and specificity results for multimodal models on ClinBench_{MM}.

Model	Sensitivity (%)	Specificity (%)
Deepseek-V3.1	42.66	44.38
GPT-5	48.60	51.49
Gemini 2.5 Pro	46.46	50.92
HuatuoGPT-o1-70B	40.71	52.33
GPT-4o	41.75	37.48
Qwen3-32B	41.17	33.07

Table 10: Sensitivity and Specificity Results on ClinBench_{Text}

Model	Sensitivity (%)	Specificity (%)
Lingshu-32B	42.85	48.72
Qwen3-VL-32B	42.31	47.16
GPT-4o	40.28	39.54
Qwen-VL-72B	39.62	37.83

Table 11: Sensitivity and Specificity Results on ClinBench_{MM}

E.4 Analysis and Insights

These clinically-oriented metrics reveal important characteristics of model performance that simple accuracy metrics cannot capture:

- **Superior Balanced Performance in General LLMs:** State-of-the-art general models GPT-5 and Gemini 2.5 Pro achieve the best overall performance with well-balanced metrics (GPT-5: 48.60% sensitivity, 51.49% specificity; Gemini 2.5 Pro: 46.46% sensitivity, 50.92% specificity), demonstrating that large-scale pre-training and advanced reasoning mechanisms can achieve robust clinical reasoning capabilities across both sensitivity and specificity.
- **Specificity Advantage in Medical-Specialized Models:** HuatuoGPT-o1-70B demonstrates the highest specificity (52.33%) among all text models, though with lower sensitivity (40.71%). This indicates that medical domain-specific training particularly improves the model's ability to avoid false positives—a critical consideration for clinical deployment where minimizing incorrect positive predictions is essential.
- **Domain Specialization in Multimodal Settings:** Lingshu-32B, as a medical-specialized multimodal model, achieves the highest performance in both metrics among multimodal models (42.85% sensitivity, 48.72% specificity), outperforming larger general-purpose models like Qwen-VL-72B. This highlights the value of domain adaptation in multimodal medical AI systems.
- **Multimodal-Text Gap:** A noticeable performance gap persists between text-based and multimodal models. The best multimodal model (Lingshu-32B: 42.85% sensitivity, 48.72% specificity) lags behind top text models (GPT-5: 48.60% sensitivity, 51.49% specificity), suggesting that visual encoding and cross-modal reasoning remain challenging areas requiring further architectural improvements.

F Cases of ClinBench

We present two specific cases, one multimodal and one text-based, each including detailed questions, options, and the model's responses (from GPT-4o and Llama-3.2-11B-Vision) in Figure 11 and 12.

Evaluation Prompt for ClinBench_{TEXT}

Question: {Question}
Options: {Options}

Instruction: Given the following multiple-choice question and options, provide a concise answer based on accurate knowledge. Please think carefully and finally conclude your response with the correct option in the format: The answer is [Option].

Figure 7: Evaluation Prompt for ClinBench_{TEXT}

Evaluation Prompt for ClinBench_{MM}

Question: {Question}
Images: {Images}
Options: {Options}

Instruction: Given the following multiple-choice question, associated medical images, and options, analyze the images and question to provide a concise, accurate answer based on medical knowledge. Please think carefully and finally conclude your response with the correct option in the format: The answer is [Option].

Figure 8: Evaluation Prompt for ClinBench_{MM}

Error Classification Prompt

Role: You are an expert medical education evaluator tasked with analyzing errors made by large language models on clinical reasoning questions.

Task: Given a clinical question, the correct answer, and a model's incorrect response, classify the error into ONE of the following five categories:

Error Categories:

- 1. Key Information Omission:** The model neglects critical clinical information in the question (e.g., important symptoms, signs, laboratory results, or medical history), leading to diagnostic errors.
- 2. Clinical Context Misunderstanding:** The response reflects an incorrect interpretation of the clinical scenario or question context, leading to inappropriate reasoning paths.
- 3. Medical Knowledge Deficiency:** The model correctly interprets the question but produces an incorrect answer due to insufficient or inaccurate medical knowledge.
- 4. Priority Judgment Error:** The model fails to identify the most primary, urgent, or contextually appropriate diagnosis among multiple plausible options.
- 5. Temporal Sequence Confusion:** The model misunderstands the timeline of disease progression, acute versus chronic course, or the sequence of clinical events.

Input:

Question: {Question}

Correct Answer: {Correct Answer}

Model's Answer: {Model's Answer}

Model's Reasoning: {Model's Reasoning}

Instruction: Choose only ONE category that best describes the primary error. Provide your analysis in the following format: Error Category: [One of the five categories] and Justification: [A brief explanation (2-3 sentences) explaining why this error falls into the selected category].

Figure 9: Error Classification Prompt for Clinical Reasoning Error Analysis

ClinBench Medical Case Assessment

Case 1/50 Dermatology

Guidelines for Experts
Dear Expert, welcome to participate in the ClinBench evaluation. Please select the most appropriate option to complete the following case study.

Patient Report

General Information

Patient: Male, 8 days old

History:

- 24 weeks gestation: Detected a solid, isoechic mass on the right sole, indicating a hemangioma, which grew slowly.
- Postnatal: Right sole swelling noted, elevated from the skin, resembling a prominent mass with a large black spot on the surface.
- Additional findings: Scattered black moles on the skin of the right thigh and right knee joint. The patient's mental state remains normal, with no significant abnormalities in diet, bowel movements, or development.

Examination

Specialized Examination:

- Right Sole Mass: Elevated from the skin, presenting a prominent mass with a large black spot, no hair growth, and uniform color.
- Same-sided Limbs: Scattered small black spots on the skin surface of the right knee and thigh.

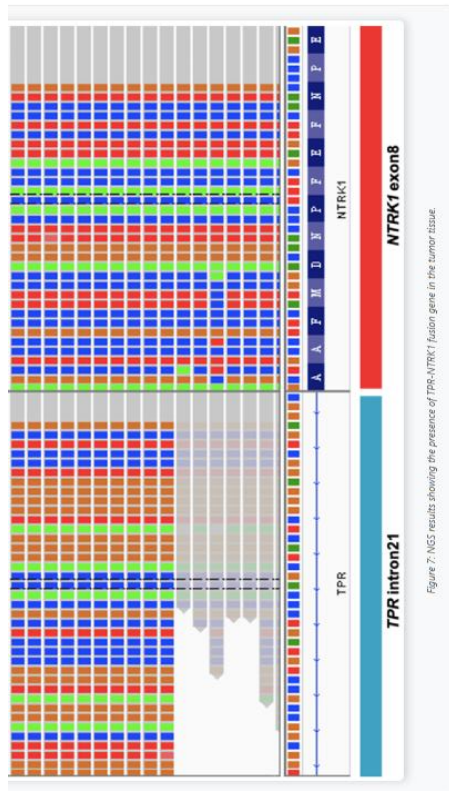
Ultrasound Examination:

- 24 weeks gestation: Solid, isoechic, mass, 21 mm x 17 mm.
- 30 weeks gestation: Increased to 42 mm x 31 mm.
- Postnatal: Mass size 40 mm x 50 mm x 30 mm, presenting mixed echo mass indicative of a hemangioma with slow growth.

MRI Examination:

- Right Sole: Large, solid soft tissue mass, homogeneous signal. T1 sequence shows slightly low signal, PD fat suppression sequence indicates slightly high signal, with many thickened blood vessels within the mass.
- Right Foot Bone: No significant abnormalities detected.





What is the most probable diagnosis for this patient?

- A. Pigmented Epithelioid Melanocytoma
- B. Synovial Sarcoma
- C. Malignant Peripheral Nerve Sheath Tumor
- D. Congenital NTRK-rearranged Melanoma with Congenital Melanocytic Nevus
- E. Congenital Hemangioma
- F. Ewing Sarcoma
- G. Neurotrophic Tyrosine Kinase (NTRK) Fusion Sarcoma

Submit Answer

Figure 10: The online evaluation of ClinBench.

Guideline for Constructing ClinBench_{Text}

Dear contributor, welcome to the ClinBench question construction task. You are now assigned to build high-quality samples for the ClinBench_{Text} dataset. Please carefully follow the guidelines below.

1. Constructing the Question Stem

Objective: Use the original content of the medical journal to construct a question stem, while minimizing modification of the source content.

1.1 For Chinese-language journals, directly extract the *Patient Information* and *Medical Examination Results* sections.

1.2 For English-language journals, locate the *Patient Information* and *Medical Examination Results* and use a professional translation tool to translate them into fluent, accurate English. Discard journals that do not contain both sections.

1.3 If the case contains images (e.g., X-ray, MRI, CT), **remove the images**, but ensure the question stem includes a concise and accurate textual description of the image. If such a description cannot be provided, discard the journal.

1.4 The question stem should be self-contained, medically accurate, and unambiguous, providing enough context for diagnosis.

2. Constructing Candidate Answer Options

Objective: Create a set of multiple-choice options, including one correct answer and several plausible distractors.

2.1 **Identifying the Correct Answer:** Extract the final diagnosis from the Diagnosis/Conclusion section. The correct answer must be medically precise and supported by evidence from the journal.

2.2 **Selecting Distractors from the Discussion:** Extract suspected or differential diagnoses from the Discussion section to serve as medically relevant distractors.

2.3 **Generating Additional Distractors Using AI:** Use GPT-5 or a similar model to generate plausible but incorrect distractors. They should reflect realistic diagnostic confusion within the relevant medical context.

2.4 **Merging Options:** Use an AI model or human review to ensure:

- (1) Exactly one correct answer is included.
- (2) No ambiguity or semantic overlap among options.
- (3) At least six options in total.

2.5 **Final Verification:** All options must be:

- (1) Medically plausible and realistic.
- (2) Mutually exclusive and clearly distinct.
- (3) Professionally formatted, with consistent language and style.

3. Saving the Output

Please save each constructed question in a single-entry JSON format, as shown below:

```
{
  "id": "ClinBench-text-00001",
  "question_stem": "",
  "options": "",
  "answer": "",
}
```

Table 12: Guideline for constructing ClinBench_{Text} questions.

Guideline for Constructing ClinBench_{MM}

Dear contributor, welcome to the ClinBench multimodal question construction task. Your task is to build multimodal (text-image) questions based on previously constructed ClinBench_{Text} questions. Please carefully follow the guidelines below.

1. Adding Medical Images

Objective: Enhance the existing ClinBench_{Text} question stems by appropriately inserting medical images from the original case journals.

1.1 Identify suitable medical images (e.g., X-ray, MRI, CT scans, histological images) within the original case journal.

1.2 Insert the identified image at a relevant location within the question stem. Ensure the image directly supports and complements the textual content.

1.3 If no suitable medical images can be found or inserted into the text, discard the question.

2. Modifying the Question Stem

Objective: Adjust the existing text-based question stem by removing textual descriptions of medical examination results depicted in the inserted images.

2.1 Remove explicit textual descriptions of medical findings or examination results that are clearly illustrated by the medical image.

2.2 Retain only concise image references within the text (e.g., *"The patient's lung condition is shown in Figure 1."*).

2.3 Ensure the modified question stem remains clear, medically accurate, and contextually complete when combined with the inserted image.

3. Saving the Output

Please save each multimodal question as a JSON file following the single-entry format below:

```
{
  "id": "ClinBench-mm-00001",
  "question_stem": "",
  "options": "",
  "answer": "",
  "image": ""
}
```

Ensure each image filename corresponds clearly with the JSON entry and store images in a designated folder. Submit JSON files along with associated images.

Table 13: Guideline for constructing ClinBench_{MM} multimodal questions.

Guideline for ClinBench Verification (Medical Students)

Dear contributor, welcome to the ClinBench question quality check task. Your role is critical to ensuring the high quality of constructed ClinBench_{Text} and ClinBench_{MM} questions. Please carefully follow the guidelines below.

1. Checking for Missing Information

Objective: Identify and correct significant information omissions in the constructed question stem compared to the original medical journal case.

1.1 Carefully review the constructed question stem against the original medical journal case to identify missing critical medical information necessary for understanding or answering the question.

1.2 Specifically, for ClinBench_{MM} questions, pay special attention to information inadvertently omitted during the manual removal of image-related descriptions.

1.3 If you identify missing crucial information, supplement the question stem by adding the necessary content clearly and concisely, ensuring the revised question stem provides adequate context for accurate diagnosis.

2. Checking for Inappropriate Candidate Options

Objective: Evaluate and refine the candidate answer options to ensure clarity, distinctiveness, and medical accuracy.

2.1 Review candidate options carefully for any overlaps with the correct answer or ambiguous phrasing that may confuse test-takers or reduce question clarity.

2.2 Remove or revise any candidate options identified as overlapping with the correct answer or ambiguously phrased, ensuring the final set of options is clear, distinct, and unambiguous.

3. Saving the Corrected Output

After performing the above checks and corrections, save each question in a single-entry JSON format as follows:

```
{
  "id": "ClinBench-check-00001",
  "question_stem": "(Corrected and complete question stem)",
  "options": "(Verified and corrected candidate options)",
  "answer": "(Confirmed correct answer)",
  "image": "(Image filename if applicable, otherwise empty)"
}
```

Ensure each corrected JSON file is named appropriately and clearly linked with any associated images.

Table 14: Guideline for checking and verifying constructed ClinBench questions.

Guideline for Expert Review of ClinBench.

Dear expert reviewer, thank you for participating in the quality assurance of the ClinBench dataset. Your expertise is crucial for ensuring the accuracy and clinical relevance of our medical questions. Please carefully follow the guidelines outlined below.

Objective: Evaluate each question comprehensively from the following four perspectives: clarity of expression, clinical rationality, factual accuracy, and appropriateness of candidate options.

Evaluation Criteria:

- **Clarity of Expression:** Ensure the question stem and options are clearly phrased, understandable, and professionally articulated.
- **Factual Accuracy:** Verify that the question stem and correct answer are medically accurate, evidence-based, and free from factual errors.
- **Appropriateness of Candidate Options:** Check whether the options provided are medically plausible, clearly distinct, non-overlapping, and appropriately challenging.

Providing Feedback:

If any issues or concerns are identified during your evaluation, clearly outline your feedback as follows:

- Describe the specific issue identified (clarity, rationality, accuracy, or appropriateness).
 - Provide detailed recommendations or corrections for resolving the identified issue.
 - Highlight any relevant medical literature or sources that support your feedback if necessary.
-

Table 15: Guideline for expert review of ClinBench questions.

F.1 Example of GPT-4o

Text-Modal task

Question:

Patient Information

** General Information:

** Patient: Female, 38 years old

** Chief Complaint: Abnormal urine test for 2 months

** Symptoms: None reported (rash, purpura, joint pain, fever, Raynaud's phenomenon, nausea, vomiting)

** Past Medical History:

- Hepatitis B carrier, no antiviral treatment,
- No history of hepatitis C
- Hypertension detected in March 2022, no antihypertensive therapy started

** Surgical History:** Cholecystectomy in March 2022

** Lifestyle:** Non-smoker, non-drinker.

** Family History: No familial hereditary disease or cancer.

** Physical Examination:

- Anemic appearance, no rashes, no enlarged lymph nodes
- Cardiovascular, respiratory, and abdominal examinations: Negative

** Electron Microscopy: - Features consistent with immune complex-mediated glomerulonephritis

** Genetic Test:

** MYD88 L265P Mutation: Negative (result may be influenced by low abnormal cell count)

Answer Choices: (A) Chronic Lymphocytic Leukemia (B) Essential Mixed Cryoglobulinemia (C) Heavy Chain Disease (D) Membranoproliferative Glomerulonephritis type I (E) Monoclonal Gammopathy of Renal Significance (F) Monoclonal Gammopathy of Undetermined Significance (G) Type I Cryoglobulinemic Glomerulonephritis (H) Type II Cryoglobulinemic Glomerulonephritis (I) Waldenstrom's Macroglobulinemia

Correct Answer: (I)

Model Answer: (E)

Model Analysis:

To determine the correct answer, the patient presents with abnormal urine findings (proteinuria, hematuria, elevated ACR, and UTP) and has laboratory and biopsy results indicative of immune complex-mediated glomerulonephritis. Key findings include:

** Renal Biopsy: Immune complexes (IgM, C3, Kappa), suggest monoclonal paraprotein involvement.

** Blood and Urine Studies: Presence of monoclonal protein (IgM and Kappa light chains), increased free Kappa, positive SPE (Serum Protein Electrophoresis), and unchanged renal structure.

** Thus, the findings fulfill the criteria for **Monoclonal Gammopathy of Renal Significance (MGRS)***, wherein monoclonal proteins cause renal disease without meeting criteria for malignancy.

** The answer is E.

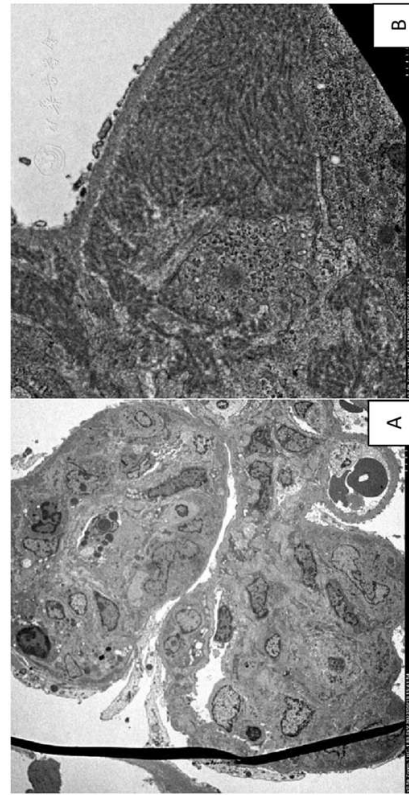


Figure 11: Electron Microscopy Image

F.2 Example of Llama-vision

Multi-Modal task

Question:

** Clinical Case Report

** Patient Information

** Occupation: Raw material crusher in an industrial ceramics factory

** Case Details

The patient was admitted on October 16, 2018, due to recurrent cough and sputum production for over 10 years, which had worsened with chest tightness for 5 days.

** Initial Examination

** Symptoms: Lip cyanosis, reduced breath sounds in the left lung, bilateral basal lung moist rales; no pathological murmurs in valve auscultation areas.

** Supplementary Check

** Examination Findings: Palpable crepitus in the neck, left supraclavicular fossa, and left chest wall. Auscultation: Diffuse wheezing, numerous moist rales heard at lung bases.

** Chest CT Findings:

- Multiple gas shadows in mediastinum, supraclavicular fossa, and subcutaneous left chest wall.
- Small area of pneumothorax on the left side.
- Nutritional support and appropriate respiration exercises.

Answer Choices: (A) Chronic Obstructive Pulmonary Disease (B) Community-Acquired Pneumonia (C) Pneumomediastinum (D) Pneumothorax (E) Primary Spontaneous Pneumothorax (F) Pulmonary Fibrosis (G) Secondary Mediastinal Emphysema (H) Type II Cryoglobulinemic Glomerulonephritis (I) Waldenstrom's Macroglobulinemia

Correct Answer: (G)

Model Answer: (G)

Model Analysis:

** Step 1: Analyze the patient's symptoms and medical history.** The patient is a 56-year-old male with a history of recurrent cough and sputum production for over 10 years, which had worsened with chest tightness for 5 days. He works as a raw material crusher in an industrial ceramics factory.

** Step 2: Analyze the exacerbation symptoms and findings.** Aggravated frequent coughing with yellow purulent sputum, severe chest tightness, orthopnea, left-sided neck swelling and pain - Palpable crepitus in the neck, left supraclavicular fossa, and left chest wall - Diffuse wheezing, numerous moist rales heard at lung bases.

** Step 3: Determine the most likely diagnosis based on the information provided.** Given the patient's occupation in an industrial ceramics factory, the presence of pneumothorax, pneumomediastinum, and subcutaneous emphysema, along with the history of recurrent respiratory issues, the most likely diagnosis is related to occupational exposure leading to lung damage.

** Answer: G. Secondary Mediastinal Emphysema.

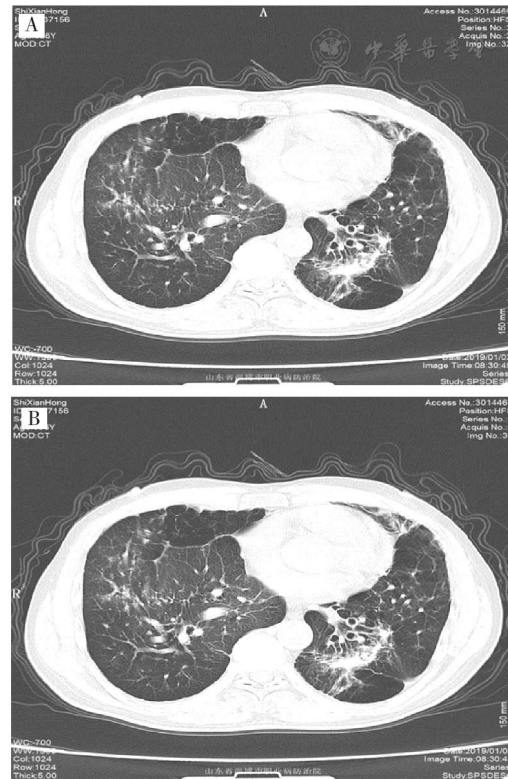


Figure 12: Chest CT Image