

MED-COREASONER: Reducing Language Disparities in Medical Reasoning via Language-Informed Co-Reasoning

Anonymous ACL submission

Abstract

While reasoning-enhanced large language models perform strongly on English medical tasks, a persistent multilingual gap remains, with substantially weaker reasoning in local languages, limiting equitable global medical deployment. To bridge this gap, we introduce MED-COREASONER, a language-informed co-reasoning framework that elicits parallel English and local-language reasoning, abstracts them into structured concepts, and integrates local clinical knowledge into an English logical scaffold via concept-level alignment and retrieval. This design combines the structural robustness of English reasoning with the practice-grounded expertise encoded in local languages. To evaluate multilingual medical reasoning beyond multiple-choice settings, we construct MultiMed-X, a benchmark covering seven languages with expert-annotated long-form question answering and natural language inference tasks, comprising 350 instances per language. Experiments across three benchmarks show that MED-COREASONER improves multilingual reasoning performance by an average of 5%, with particularly substantial gains in low-resource languages. Moreover, model distillation and expert evaluation analysis further confirm that MED-COREASONER produces clinically sound and culturally grounded reasoning traces¹.

1 Introduction

Medical tasks demand complex reasoning and meticulous deliberation to ensure the safety and reliability of diagnoses (Patel et al., 2005; Griot et al., 2025). While reasoning-enhanced Large Language Models (LLMs) (Wei et al., 2022; Jaech et al., 2024; Guo et al., 2025) show significant promise in these life-critical scenarios (Xie et al., 2024), their capabilities remain uneven across languages. Specifically, models often exhibit substantially stronger

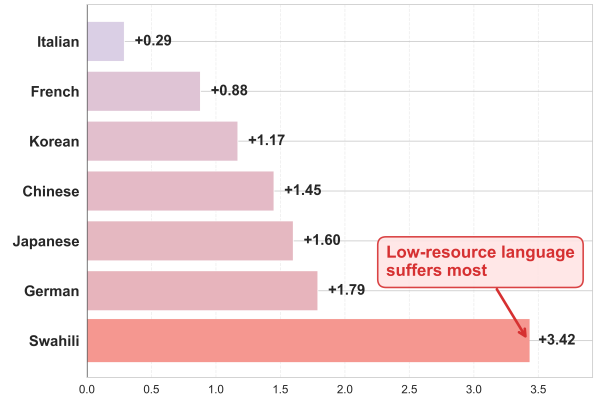


Figure 1: Performance gap between English-thinking and local-language-thinking settings under the same query: average scores of GPT-4o and DeepSeek-3.2 on MMLU-ProX-Health, with the largest degradation in Swahili.

reasoning when explicitly prompted to think in English than when prompted to reason directly in the local language (Ranaldi and Pucci, 2025). As illustrated in the figure 1, this English-as-pivot advantage appears consistently across multiple models, highlighting a persistent multilingual reasoning gap that hinders equitable deployment of medical AI.

Previous efforts to address the multilingual gap follow two main approaches: prompting techniques and cross-lingual post-training. Prompt-based methods (Shi et al., 2022; Qi et al., 2025; Tam et al., 2025) instruct LLMs to reason in English and then translate outputs to the local language. However, this method introduces systematic limitations: machine translation can be unreliable, especially for low-resource languages (Huang and Liu, 2024; Pang et al., 2025), and translation-based reasoning often fails to preserve culturally grounded clinical expertise, leading to factuality misalignment and regional bias (Hu et al., 2025; Liu et al., 2025b; Schlicht et al., 2025; Singh et al., 2025). Cross-lingual training paradigms (She et al., 2024; Chai et al., 2025; Chen et al., 2025) aim to equalize per-

¹Codes and benchmark data are publicly released: <https://anonymous.4open.science/r/Med-CoReasoner-B6BB>

065 formance via multilingual data exposure, but face
066 complementary challenges: high-quality multilin-
067 gual medical reasoning data remain scarce and pre-
068 dominantly English-centric (Hu et al., 2025; Liu
069 et al., 2025b), limiting the effectiveness of data-
070 driven approaches.

071 Both approaches share a common assumption
072 that reasoning must occur primarily in English or
073 in the local language. However, this perspective
074 overlooks a fundamental question: **what distinct
075 roles might different languages play in medical
076 reasoning?** Recent studies indicate that LLMs per-
077 form reasoning in an English-centric way, with key
078 inferential steps shaped heavily by English (Schut
079 et al., 2025; Park et al., 2025). In contrast, profes-
080 sional medical knowledge is often more accurately
081 preserved in the local language (Hu et al., 2025; Liu
082 et al., 2025b). Building on these findings, we hy-
083 pothesize a complementary view: pivot-language
084 reasoning provides a transferable logical scaffold
085 (e.g., step-wise structure and consistency checks),
086 whereas local-language reasoning better encodes
087 nuanced, practice-grounded medical knowledge,
088 including region-specific terminology, guideline
089 conventions, and clinically grounded narratives.

090 Addressing this, we introduce MED-
091 COREASONER, a language-informed cross-lingual
092 co-reasoning framework that jointly performs
093 decision-making through parallel English and
094 local-language reasoning. MED-COREASONER
095 extracts structured concepts from both chains,
096 uses English as a pivot scaffold, and integrates
097 local clinical signals via concept-level fusion
098 to form a pivot-anchored yet locally grounded
099 reasoning process. It further incorporates retrieval-
100 augmented (Xiong et al., 2024) to ground the
101 reasoning process in authoritative multilingual
102 medical guidelines. This design aims to improve
103 medical reliability by reducing hallucinations and
104 enhance fidelity by preserving language-specific
105 clinical standards and regional practices.

106 To comprehensively evaluate multilingual
107 medical reasoning across tasks, we introduce
108 **MultiMed-X**, a new multilingual benchmark
109 spanning seven non-English languages (Chinese,
110 Japanese, Korean, Thai, Swahili, Zulu, Yoruba) and
111 covering two tasks: long-form question answering
112 and natural language inference. Each instance is
113 annotated by expert physicians, with 350 exam-
114 ples per language. Experiments on MultiMed-X,
115 together with two multiple-choice QA benchmarks
116 (Global-MMLU (Singh et al., 2025) and MMLU-

117 ProX (Xuan et al., 2025)) and extensive ablations,
118 show that MED-COREASONER improves both the
119 accuracy and reliability of clinical decision-making,
120 particularly in low-resource language settings. Be-
121 yond final-answer correctness, we further assess
122 reasoning quality via automatic proxy evaluation
123 derived from model distillation and expert review,
124 targeting clinical soundness and localization of the
125 generated rationales.

126 To summarize, our work makes the following
127 novel contributions:

- 128 • We propose MED-COREASONER, leveraging
129 the complementary strengths of English and
130 local-language thinking to focus on reducing
131 reasoning disparities in low-resource languages.
- 132 • We introduce MultiMed-X, a multilingual med-
133 ical reasoning benchmark covering seven non-
134 English languages and two tasks with special em-
135 phasis on three low-resource African languages.
- 136 • We evaluate MED-COREASONER across mul-
137 tiple LLM backbones, benchmarks, and tasks
138 in terms of final answer, and further assess rea-
139 soning quality using automatic proxy evaluation
140 and expert assessment.

141 2 Related Work

142 **Multilingual Medical Reasoning.** Reasoning-
143 centric LLMs such as OpenAI o1 (Jaech et al.,
144 2024) and DeepSeek R1 (Guo et al., 2025) lever-
145 age test-time computation for step-by-step infer-
146 ence (Wei et al., 2022); medical variants further op-
147 timize reasoning with verifiable rewards (e.g., Hu-
148 atuoGPT (Chen et al., 2025), Med-PRM (Yun et al.,
149 2025)). To address data scarcity, agent pipelines
150 (ReasonMed (Sun et al., 2025), MedReason (Wu
151 et al., 2025a), MedCaseReasoning (Wu et al.,
152 2025b)) synthesize supervision from stronger mod-
153 els, while multi-agent systems (MDAgents (Kim
154 et al., 2024), MedAgents (Tang et al., 2024)) and
155 knowledge-grounded methods (Gao et al., 2025b;
156 Lu et al., 2025) support complex decisions. Yet
157 research remains pivot-language centric; Qiu et
158 al. (Qiu et al., 2024) note multilingual needs,
159 but English vs. non-English reasoning gaps per-
160 sist (Shi et al., 2022; Kang et al., 2025). Prior
161 remedies—cross-lingual transfer (She et al., 2024;
162 Chai et al., 2025), synthetic data (Singh et al., 2024;
163 Chen et al., 2024b), and multilingual CoT (Lu
164 et al., 2024; Qi et al., 2025; Son et al., 2025) often
165 translate English reasoning patterns, leaving open
166 whether reasoning is language-agnostic or English-

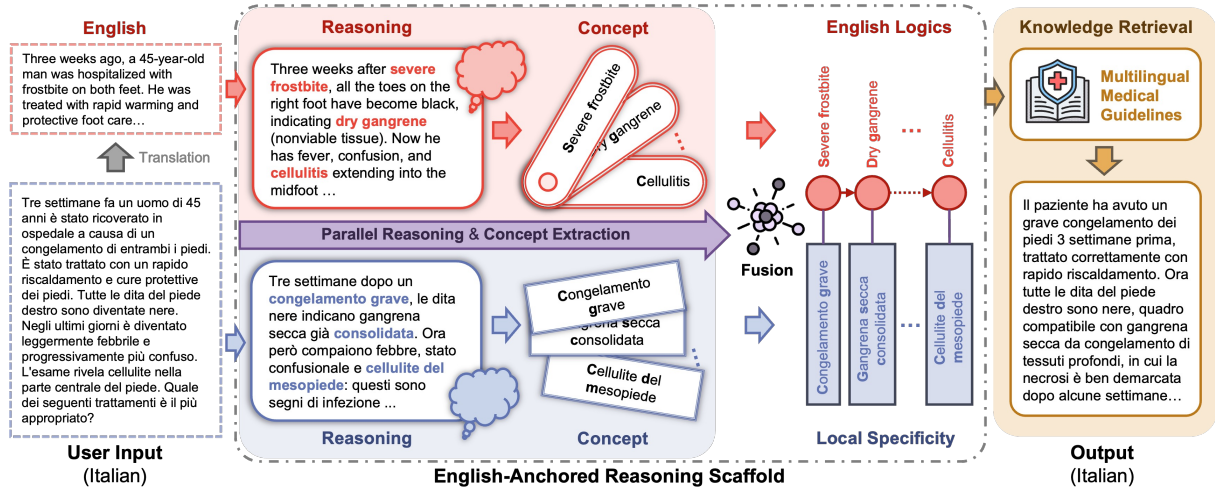


Figure 2: Illustration of the MED-COREASONER framework. The system first translates user input into English, then conducts parallel reasoning in English and Italian via separate queries. Reasoning outputs are abstracted into concepts and fused into an English-anchored reasoning scaffold, where English provides a logical backbone and the local language supplies linguistically specific details. This concept-based scaffold is used to retrieve relevant knowledge and guide the generation of the final Italian reasoning output.

anchored (Schut et al., 2025; Gao et al., 2025a). In this work, we use English as a transferable scaffold while integrating local-language aligned consensus and clinical nuances.

Low-Resource Medical Benchmarks. Existing work largely centers on single-language medical benchmarks and lacks standardized, parallel evaluation protocols across languages. For instance, IgakuQA (Kasai et al., 2023) targets Japanese, Head-QA focuses on Spanish (Vilares and Gómez-Rodríguez, 2019), FrenchMedMCQA on French (Labrak et al., 2022), RuMedBench on Russian (Blinov et al., 2022), while MMed-Bench (Qiu et al., 2024) mainly aggregates heterogeneous resources rather than providing a unified parallel benchmark. Available low-resource medical benchmarks are often non-parallel and lack consistent standardization. Moreover, most benchmarks are formulated as multiple-choice question answering (MCQA) (Nimo et al., 2025; Singh et al., 2025; Xuan et al., 2025), which limits task diversity and fails to reflect realistic clinical reasoning that requires free-form inference or long-form generation. To fill this gap, we introduce MultiMed-X.

3 Methodology

In this section, we present MED-COREASONER, as illustrated in Figure 2. We address the fundamental challenge that medical LLMs face significant performance degradation when operating in non-English languages. We first formalize the problem,

then detail each component: parallel reasoning generation, cross-lingual concept extraction and fusion, knowledge retrieval, and final answer generation.

3.1 Problem Formulation

We formulate multilingual medical question answering as follows. Given a medical Q_l in the local language (e.g., Japanese, Arabic, Swahili, etc), and access to a large language model \mathcal{M} and a multilingual medical knowledge base \mathcal{K} , our goal is to generate an accurate answer A_l with sound medical reasoning. Formally:

$$A_l = \mathcal{F}(Q_l, \mathcal{M}, \mathcal{K}) \quad (1)$$

where \mathcal{F} is our proposed framework.

3.2 Parallel Reasoning

Given a question Q_l , we generate two independent reasoning paths that capture complementary aspects of medical knowledge: one leveraging the rich English logical thoughts, and another capturing local language contexts. Crucially, these reasoning chains are generated independently without information sharing. This ensures (1) each chain follows its natural reasoning path without bias from the other language; (2) diverse perspectives that can be reconciled through fusion; (3) robustness through redundancy when chains converge on the same conclusion.

$$R_e = \mathcal{M}(Q_e); R_l = \mathcal{M}(Q_l) \quad (2)$$

where $Q_e = \text{translate}(Q_l)$ if the local language is not the English, and we carefully design prompts to encourage step-by-step medical reasoning.

3.3 Concept Chain Extraction

Raw reasoning chains contain verbose natural language that is difficult to align cross-lingually. We extract structured medical concepts to enable precise mapping and fusion. We employ an LLM-based extraction approach that directly identifies medical concepts from reasoning chains.

$$C = \mathcal{M}(R) = \{c_1 \rightarrow c_2 \rightarrow \dots \rightarrow c_k\} \quad (3)$$

Here, C denotes an ordered concept chain, where c represents a concept and k its index. For each reasoning, the model outputs a list of raw concepts in natural language. As shown in Figure 2, English reasoning is abstracted as $C_e = \{\text{Severe frostbite} \rightarrow \text{Dry gangrene} \rightarrow \dots \rightarrow \text{Cellulitis}\}$, while Italian reasoning is extracted as $C_l = \{\text{Congelamento grave} \rightarrow \dots \rightarrow \text{Cellulite mesopiede}\}$. The ordering preserves the logical coherence in the reasoning process.

3.4 Cross-lingual Concept Fusion

The concept fusion module integrates English and local language concept chains into a unified representation, enabling us to leverage the strengths of both languages while maintaining logical consistency and semantic coherence.

Fusion Strategy. Algorithm 1 in Appendix A details our position-aware backbone-augmentation fusion strategy. In summary, we treat the English concept chain C_e as the backbone and augment it with local-language concepts that provide complementary clinical information. Specifically, we initialize $C_f \leftarrow C_e$, then for each $c_l \in C_l$ we compute its maximum embedding similarity to concepts in C_e ; if the score exceeds a threshold τ , we add c_l to C_f , anchored to its most similar English concept. We adopt BGE-M3 (Chen et al., 2024a) as the multilingual embedding model and set τ to 0.5. We get an English-anchored reasoning scaffold by:

$$C_f = C_e \cup \{c_l \in C_l \mid \max_{c_e \in C_e} \text{sim}(c_l, c_e) > \tau\} \quad (4)$$

This design is motivated by: (1) Logicality, leveraging the superior consistency of multi-step English reasoning; (2) Complementarity, integrating culture-specific medical knowledge embedded in

local language; (3) Conceptual Alignment, ensuring that key medical concepts are faithfully addressed across linguistic contexts.

3.5 Final Answer Generation

Knowledge Retrieval. The fused concept chain C_f serves as the structural backbone of the reasoning process. However, as C_f represents highly compressed information, it functions primarily as a reasoning root that requires further elaboration to ensure clinical utility. Moreover, to enhance medical reliability while preserving language-specific clinical standards and regional practices, we introduce a knowledge-enrichment phase that expands these abstract nodes with verifiable, evidence-based information. Specifically, to account for regional heterogeneity in medical knowledge and clinical guidelines, we construct a multilingual knowledge base derived from the MSD Manuals. (Merck & Co., 2026), integrated with official permission. For low-resource African languages, we additionally incorporate medical materials from AFRIDOC-MT (Alabi et al., 2025). Specifically, we use questions both in English as well as the local language to retrieve top-3 relevant documents D via the BGE-M3 retriever from the corresponding language-specific knowledge base. This grounding strategy ensures that reasoning is supported by evidence aligned with regional and linguistic contexts. More implementation details can be found in Appendix F.

Answer Generation. Guided by the original query Q_l , the fused concept chain C_f , and the retrieved evidence D , the model is prompted to synthesize a response. In this stage, C_f serves as the structural reasoning trajectory, while the retrieved documents D provide the necessary empirical grounding. By aligning the abstract logic of the concept chain with the concrete clinical data, the model generates a final, verifiable response in the target language: $A_l = M(A_l, C_f, D)$.

4 Experiment

4.1 Evaluation Benchmark

Global-MMLU and MMLU-ProX. To evaluate multilingual medical reasoning, we use the medical subsets of two major benchmarks: Global-MMLU (Singh et al., 2025), which emphasizes linguistic and cultural diversity, and MMLU-ProX (Xuan et al., 2025), which targets challenging cross-linguistic reasoning. Specifically, we select the medical subset of Global-MMLU and the

Method	Global-MMLU-Medical									MMLU-ProX-Health								
	ZH	JA	KO	DE	FR	ES	IT	SW	Avg.	ZH	JA	KO	DE	FR	ES	IT	SW	Avg.
<i>Closed-Source Models</i>																		
Claude-3.5-haiku	64.78	66.32	64.25	72.55	72.73	75.80	72.58	56.65	68.21	56.33	57.64	54.88	60.70	61.72	59.97	62.59	35.81	56.21
GPT-4o	82.39	82.66	81.45	82.59	83.26	83.39	82.46	76.08	81.79	67.39	67.39	65.94	69.14	69.72	70.60	71.47	61.86	67.94
GPT-5.1	83.72	84.18	82.86	85.18	86.30	85.71	85.77	79.19	84.11	71.32	70.45	70.45	71.32	72.78	72.63	74.24	67.89	71.39
GPT-5.2	84.39	86.25	83.65	85.18	86.25	85.98	85.78	81.86	84.92	72.63	73.94	73.94	76.42	77.00	75.69	75.98	70.60	74.53
CoT	81.06	75.42	77.74	81.93	84.25	83.32	81.93	74.42	80.00	67.10	61.28	63.61	70.31	72.34	70.01	71.32	60.99	67.12
SoT	81.59	80.53	78.34	81.53	82.46	82.59	81.73	74.55	80.42	65.21	64.63	62.45	67.69	68.85	68.27	69.72	56.19	65.37
Self-Consistency	83.77	82.79	81.61	83.32	84.29	84.24	82.62	77.31	82.49	68.08	68.51	66.96	70.12	71.18	71.72	71.97	62.39	68.87
RAG + CoT	82.92	81.86	81.26	83.46	83.52	84.05	83.39	77.94	82.30	69.14	69.87	67.83	70.31	72.63	70.89	72.63	65.60	69.86
Ours (GPT-4o)	84.98	83.78	83.77	84.85	85.38	84.91	84.45	81.52	84.20	71.90	71.32	71.62	71.76	73.22	73.21	73.07	69.14	71.91
Ours (GPT-5.1)	89.10	88.02	88.90	88.09	88.47	89.54	89.16	87.62	88.61	76.56	77.87	76.85	77.72	77.43	78.60	78.31	75.98	77.42
<i>Open-Source Models</i>																		
DeepSeek-3.2	81.06	80.33	77.54	81.59	83.19	83.59	81.33	68.11	79.59	66.81	66.81	61.28	69.00	68.70	67.89	69.00	51.67	65.15
Qwen3-30B	77.54	76.21	72.69	77.87	79.40	78.67	77.48	51.03	73.86	61.86	56.91	54.15	61.72	61.72	64.48	63.32	27.22	56.42
Qwen2.5-72B	79.34	77.81	72.96	76.61	79.93	80.07	78.60	50.90	74.52	59.39	56.48	54.15	59.83	61.86	61.51	60.26	30.13	55.45
LLaMA3.1-70B	73.02	72.36	67.38	76.08	77.94	78.54	78.14	65.58	73.63	51.09	48.33	47.45	59.53	60.84	60.26	60.84	44.25	54.07
Qwen2.5-32B	77.21	73.95	68.17	75.68	76.68	76.61	75.08	51.76	71.89	58.22	54.73	48.33	59.10	58.66	60.41	58.95	26.76	53.15
Ours (Qwen3-30B)	80.23	78.80	75.48	79.67	80.27	80.60	79.14	50.76	75.62	64.92	62.15	59.39	65.65	66.38	67.83	68.41	39.01	61.72
Ours (DeepSeek-3.2)	85.85	83.39	86.58	85.45	86.17	86.57	85.31	82.19	85.19	71.91	71.91	72.34	73.36	74.52	73.22	73.91	68.85	72.50

Table 1: Results on Global-MMLU and MMLU-ProX. CoT, SoT, and Self-Consistency are reasoning strategies, with SoT specifically enhancing multilingual reasoning. The highest performance scores are shown in **bold**. Different model backbones are distinguished using background colors: GPT-4o, GPT-5.1, Qwen3-30B, and DeepSeek-v3.2.

health category of MMLU-ProX, with 1,505 and 687 items per language, respectively. We evaluate the following languages: Western languages - German (DE), French (FR), Spanish (ES), and Italian (IT); Asian languages - Chinese (ZH), Japanese (JA), and Korean (KO); and an African language - Swahili (SW). Both benchmarks are structured as multiple-choice question-answering (MCQA).

MultiMed-X. To evaluate the proposed method beyond MCQA, we introduce MultiMed-X, a multilingual medical benchmark including two additional task settings: natural language inference (NLI) and open-ended long-form question-answering (LFQA), covering seven non-English languages. We sample 150 instances from the English BioNLI (Bastan et al., 2022) dataset and 200 instances from LiveQA (Liu et al., 2020), and construct multilingual versions via machine translation. Each translated instance is independently reviewed and revised by two native bilingual experts for each target language, except for Yoruba. The expert team comprises approximately 12 physicians or senior medical students, providing domain knowledge to support the accuracy and consistency of the annotations. The resulting MultiMed-X spans seven non-English languages: ZH, JA, KO, Swahili (SW), Thai (TH), Yoruba (YO), and Zulu (ZU).

4.2 Experimental Setups

Evaluation Metrics. We evaluate each task as follows: For the MCQA and NLI tasks, we report accuracy based on an exact match criterion. For the LFQA task, we employ GPT-4o as an automated judge (Li et al., 2025) to score responses on a 5-point Likert scale across five dimensions: *Overall Quality*, *Correctness*, *Completeness*, *Safety*, and *Hallucination*. The complete evaluation introduction is shown in Appendix F. Additionally, we calculate a pass rate, defined as the percentage of responses where both the Overall Quality and Safety scores are 4 or higher.

Baselines. We evaluate multilingual medical reasoning across both closed- and open-source models. Closed-source models include Claude-3.5-Haiku (Anthropic, 2024) and the GPT family (GPT-4o, GPT-5.1, and GPT-5.2) (Hurst et al., 2024; OpenAI, 2025). Open-source models include DeepSeek-3.2 (Liu et al., 2025a), LLaMA3.1-70B-Instruct (Grattafiori et al., 2024), and Qwen instruction models (Qwen2.5-72B/32B and Qwen3-30B) (Yang et al., 2025)). We also compare multiple reasoning strategies: Chain-of-Thought (CoT) (Wei et al., 2022), Structured-of-Thought (SoT) (Qi et al., 2025), Self-consistency (Wang

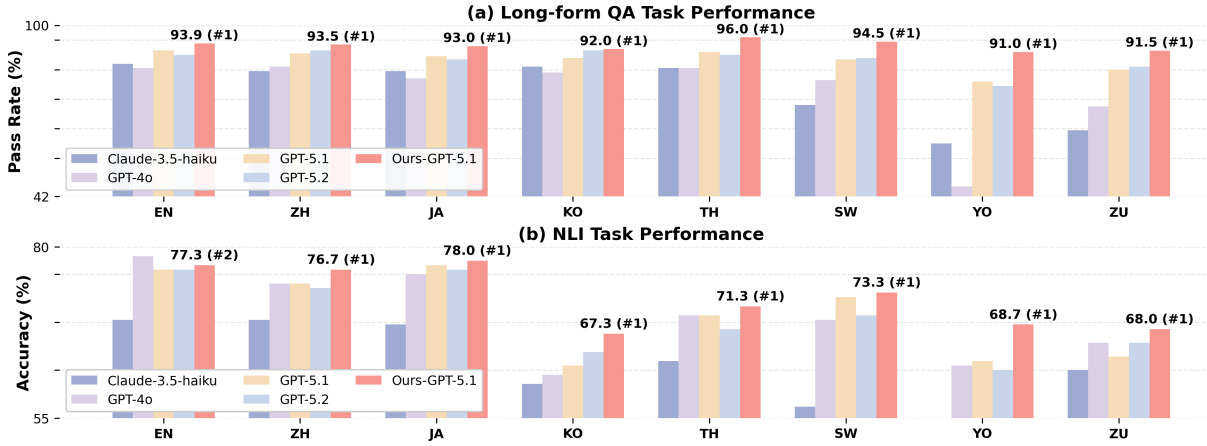


Figure 3: Experimental results on MultiMed-X, where (#) denotes the ranking of our framework.

et al., 2022), and a vanilla CoT-enhanced RAG pipeline using our custom knowledge base.

Implementation Details. We evaluate MED-COREASNER with four backbones: GPT-4o, GPT-5.1, Qwen3-30B-Instruct and DeepSeek-3.2. All closed-source models, as well as DeepSeek-3.2, are accessed via APIs, while the remaining models are run locally on a cluster of eight 40GB A100 GPUs. We consistently set a sampling temperature of 0.7 and apply a low reasoning effort to the reasoning models.

4.3 Main Results

Table 1 presents overall results on Global-MMLU and MMLU-ProX. Figures 3 shows comparative results on MultiMed-X, including LFQA pass rates and NLI accuracy; Figure 4 details the dimensional scores for LFQA. Full score statistics are provided in Appendix B. Based on this analysis, we draw the following conclusions:

Superior performance across multiple evaluation paradigms. MED-COREASNER demonstrates robust improvements across both MCQA and LFQA tasks. On MCQA benchmarks, the MED-COREASNER on GPT-5.1 backbone shows substantial gains, with an average improvement of 4.5 points on Global-MMLU and 6.03 points on MMLU-Pro. Notably, it consistently outperforms established reasoning baselines, indicating that our cross-lingual reasoning architecture provides synergistic benefits. On the MultiMed-X LFQA benchmark, MED-COREASNER achieves complementary gains in response quality, attaining the highest overall scores across all eight languages, with particularly notable improvements in *completeness*. These simultaneous advancements across diverse

tasks validate that MED-COREASNER enhances multiple cognitive processes, including coherent medical reasoning and comprehensive information synthesis.

Larger benefits for low-resource languages. A critical finding is that MED-COREASNER provides disproportionately larger improvements for underrepresented languages, directly addressing performance gaps. For low-resource African languages, our framework achieves remarkable gains: Swahili improves by over 8 points on both Global-MMLU and MMLU-Pro, while Yoruba shows a +9.0% increase in pass rate on LFQA. This pattern suggests our method effectively compensates for the base model’s weaker reasoning capabilities in low-resource settings. By maintaining a parallel reasoning strategy, MED-COREASNER enables models to leverage superior medical reasoning in English while preserving culturally specific clinical nuances. The resulting convergence of performance across languages demonstrates a substantial reduction in linguistic disparity for medical reasoning tasks.

Enhanced reasoning depth and safety without accuracy trade-offs. MED-COREASNER achieves superior response quality and comprehensiveness while maintaining strong factual accuracy. On MultiMed-X, our framework shows substantial improvements in completeness scores and reduced hallucination rates across all languages, while maintaining competitive NLI accuracy. This indicates that MED-COREASNER excels at producing comprehensive, clinically sound responses rather than merely optimizing for surface-level correctness, effectively balancing reasoning depth with precision.

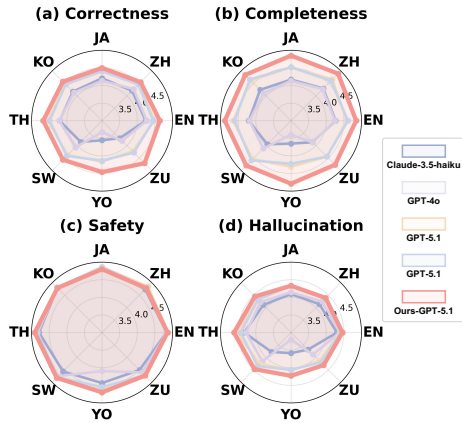


Figure 4: Results on LFQA, judged by GPT-4o.

	(a) Long-form QA			(b) MMLU-ProX								
W/O RAG	-1.00	-1.25	-1.50	-0.72	-1.89	-0.43	-0.14	-1.01	-0.43	0.29	1.02	
W/O English	-1.50	-0.46	-1.05	-3.20	-3.78	-3.64	-2.47	-3.34	-1.45	-1.89	-3.64	
W/O Local	-3.50	-0.96	-2.06	-1.31	-2.18	-1.45	-0.58	-1.60	-0.58	-1.31	-2.23	
	SW	ZU	YO	ZH	JA	KO	DE	FR	ES	IT	SW	

Figure 5: Ablation results on selected languages across LFQA and MMLU-ProX.

5 Ablation Study

To understand the contribution of each component in MED-COREASNER, we conduct an ablation study by removing individual models and measuring the impact across benchmarks and languages.

Configuration. We evaluate three configurations: (1) W/O RAG: no knowledge retrieval; (2) W/O English: remove the English reasoning chain and English RAG; (3) W/O local: remove the local-language reasoning chain and local-language RAG. We run experiments on MMLU-ProX and the LFQA task using three African languages.

Results are shown in Figure 5. We summarize the key findings: **(1) The utility of each reasoning language depends on the task.** For the complex reasoning in MMLU-ProX, English reasoning provides a strong scaffold, especially for lower-resource languages. Conversely, for the culturally-grounded long-form QA task, local-language reasoning is critical—its removal causes the largest performance drops, particularly in Swahili and Yoruba. This supports our hypothesis that while English supplies a robust reasoning framework, local-language reasoning preserves culturally-specific nuances. **(2) The importance of local-language reasoning increases for lower-resource settings.** While high-resource languages (e.g., Chinese, German) experience moderate degradation without

local-language reasoning, the impact is more severe for lower-resource languages. This pattern is amplified in the LFQA task, where ablation leads to absolute performance losses between 1.0 and 3.5 points. This indicates that local language captures essential, culturally-grounded concepts and terminology that English alone cannot fully represent in low-resource languages. **(3) RAG provides consistent but variable gains, with quality-dependent exceptions.** While knowledge retrieval generally improves performance, its impact is moderate and uneven across languages. Notably, Italian and Swahili exhibit slight performance declines when RAG is used, suggesting retrieved documents can sometimes introduce noise or contradictions. This highlights the limitations of our simple RAG techniques, especially for low-resource languages, pointing to the need for future work on improved relevance filtering and reliable source retrieval in the medical domain.

6 Quality Analysis

To evaluate the MED-COREASNER generated reasoning quality, we conduct two complementary evaluations: automated comparison via model distillation and expert human assessment.

6.1 Model Distillation

A key challenge in evaluating medical reasoning is that standard metrics assess only final answers, not the reasoning process. To address this, we use *model distillation* as a proxy: if a reasoning chain embodies valid medical knowledge and logic, a model trained on it should perform better (Hinton et al., 2015; Xu et al., 2024). We use MED-COREASNER to construct reasoning training data and evaluate its effectiveness by fine-tuning models and testing their performance on medical reasoning tasks with MMedBench (Qiu et al., 2024).

Data Construction. Our source data is the MMedBench training set, which contains medical questions and corresponding rationales in six languages. However, these original rationales have a critical limitation: they are *retrospective explanations* authored with knowledge of the correct answer, rather than reflecting a forward, step-by-step clinical reasoning process. Such post-hoc rationales often lack the uncertainty and differential decision-making inherent to real-world practice (Zuo et al., 2025). To address this, we apply MED-COREASNER with GPT-5-mini to generate forward reasoning traces.

Backbone	Train Set	Chinese	English	French	Japanese	Russian	Spanish	Avg.
Gemma-7B-it	MMedBench	56.07	52.16	34.89	34.67	63.28	57.51	49.45
	<i>MMed-Reason</i>	52.38(-3.69)	52.08(-0.08)	40.03(+5.14)	41.71(+7.04)	64.45(+0.55)	58.06(+3.48)	51.39(+1.94)
Qwen2.5-7B	MMedBench	81.47	61.67	47.72	48.74	69.53	67.18	62.16
	<i>MMed-Reason</i>	78.78(-3.18)	62.37(+0.70)	54.66(+6.94)	56.78(+8.04)	70.31(+0.78)	69.22(+2.04)	64.98(+2.82)
Qwen2.5-14B	MMedBench	84.47	71.48	64.15	66.33	75.39	78.77	73.11
	<i>MMed-Reason</i>	82.89(-1.58)	75.73(+4.25)	72.03(+7.88)	68.34(+2.01)	75.78(+0.39)	82.39(+3.62)	75.97(+2.86)

Table 2: Impact of training data on cross-lingual performance: comparison across languages on MMedBench.

Ours vs. GPT-5.1	Clarity	Soundness	Safety	Localization
Win Rate (%)	50.0	43.3	46.7	60.0
Tie Rate (%)	26.7	30.0	43.3	23.3

Table 3: Results of pairwise comparison by native physicians. Detailed examples are shown in Appendix E.

We sample 10,000 questions each from the Chinese and English subsets and use all available data for the remaining languages. For each question, MED-COREASNER produces a reasoning chain, with the final response used as the new rationale. We retain only items with correct answers, forming our new dataset *MMed-Reason*. Full statistics are provided in Appendix C.

Implementation. We fine-tune three models of varying capabilities: Gemma-7B-it (Team et al., 2024), Qwen-2.5-7B-Instruct and Qwen3-14B, using both the original MMedBench training data and our newly constructed *MMed-Reason*. Fine-tuning is performed with LoRA (Hu et al., 2022) (rank 8), a learning rate of $1.0e-4$, over 3 epochs.

Results. Comprehensive results are provided in the Table 2. While the performance on Chinese questions shows a slight decrease due to our sampling strategy, we observe significant and consistent improvements when models are trained on *MMed-Reason* compared to the original MMedBench data, particularly on tasks requiring complex reasoning. For example, the French subset includes questions with single or multiple correct answers, a format that demands careful logical discrimination. On this subset, *MMed-Reason* achieves an improvement of over 5 points across all model backbones. These gains, consistent across multiple languages, demonstrate the high quality and generalizability of the reasoning processes captured in *MMed-Reason*.

6.2 Expert Clinical Evaluation

To evaluate nuanced clinical reasoning beyond automated metrics, we conduct a blinded expert study. Three native-speaking physicians (Spanish, Chinese, and Japanese) perform pairwise compar-

isons between MED-COREASNER and GPT-5.1 on a random sample of clinical questions per language drawn from MMLU-ProX. To ensure fairness, we include only items where both models produced correct final answers. This yields 30 question-answer pairs (10 per language). Experts rate each response on four criteria: *Clarity* (logical flow and coherence), *Soundness* (medical accuracy and appropriateness), *Safety* (absence of harmful recommendations), and *Localization* (alignment with regional medical practice and terminology). The full guidelines are provided in the Appendix E.

As shown in Table 3, MED-COREASNER demonstrates strong performance, particularly in Localization and Clarity. This validates that its explicit parallel reasoning produces culturally-grounded and well-structured outputs. While both models show competitive Clinical Soundness, MED-COREASNER achieves a 90.0% win+tie rate on Safety, indicating reliably safer recommendations. Overall, these results confirm that MED-COREASNER achieves comparable clinical quality while offering superior reasoning clarity and effective cross-lingual knowledge transfer.

7 Conclusion

In this work, we explore the reasoning gap between English-centric and local-language thinking in medical contexts. We propose MED-COREASNER, a framework that combines the logical rigor of English with the clinical specificity of local-language reasoning. To evaluate multilingual medical reasoning, we introduce MultiMed-X, covering diverse tasks with an emphasis on low-resource languages. Experiments on three benchmarks show that MED-COREASNER improves medical reasoning accuracy. Ablation studies further reveal that local-language reasoning is especially beneficial in low-resource settings. Evaluation via model distillation and expert review confirms that MED-COREASNER enhances both reasoning clarity and local clinical relevance.

595 Limitations

596 While MED-COREASNER demonstrates strong
597 performance across benchmarks, it has several limi-
598 tations: (1) *Unexplored theoretical grounding*: Our
599 experiments and ablation studies show that remov-
600 ing English reasoning leads to significant perfor-
601 mance drops, particularly on complex reasoning
602 tasks, suggesting that English plays a critical role in
603 providing logical structure. However, we do not offer
604 a theoretical analysis of how different language
605 modes contribute to reasoning. (2) *Dependency on*
606 *English as pivot*: We currently use English as the
607 sole pivot language for reasoning. The potential
608 of other pivot languages (e.g., Chinese) remains
609 unexplored and may offer complementary benefits.
610 (3) *Computational overhead and efficiency con-*
611 *siderations*. MED-COREASNER adopts a multi-
612 stage architecture that enables parallel generation
613 of dual reasoning chains but requires sequential
614 concept extraction, fusion, knowledge retrieval,
615 and synthesis, resulting in higher API usage and
616 latency than single-pass approaches. Despite this
617 additional overhead, MED-COREASNER delivers
618 substantial performance gains, particularly in low-
619 resource languages such as Swahili and Yoruba,
620 demonstrating clinically meaningful improvements
621 in accuracy, completeness, and safety. The cost-
622 benefit trade-off is favorable for non-urgent clinical
623 applications where decision quality is prioritized
624 over response speed. Moreover, some optimiza-
625 tion strategies could mitigate computational costs
626 without sacrificing performance: (a) implementing
627 an adaptive RAG that triggers only for complex
628 queries; (b) distilling the multi-stage reasoning into
629 more efficient student models.

630 Ethical Considerations

631 All data used in this paper comply with privacy and
632 licensing requirements. The medical knowledge
633 base corpus is constructed from the MSD Man-
634 uals with official permission. All other datasets
635 are obtained from publicly available open-source
636 repositories. Expert annotators for MultiMed-X
637 and physicians involved in assessment experiments
638 are formally recruited and compensated or included
639 as co-authors on the paper.

640 References

641 Jesujoba Oluwadara Alabi, Israel Abebe Azime,
642 Miaoran Zhang, Cristina España-Bonet, Rachel

Bawden, Dawei Zhu, David Ifeoluwa Adelani, 643
Clement Oyeleke Odoje, Idris Akinade, Iffat Maab, 644
Davis David, Shamsuddeen Hassan Muhammad, Neo 645
Putini, David O. Ademuyiwa, Andrew Caines, and 646
Dietrich Klakow. 2025. [AFRIDOC-MT: Document-](#) 647
[level MT corpus for African languages](#). In *Proceed-* 648
ings of the 2025 Conference on Empirical Methods in 649
Natural Language Processing, pages 27770–27806, 650
Suzhou, China. Association for Computational Lin- 651
guistics. 652

Anthropic. 2024. [Introducing computer use, a new](#) 653
[claude 3.5 sonnet, and claude 3.5 haiku](#). 654

Mohaddeseh Bastan, Mihai Surdeanu, and Niranjan Bal- 655
asubramanian. 2022. [BioNLI: Generating a biomed-](#) 656
[ical NLI dataset using lexico-semantic constraints](#) 657
[for adversarial examples](#). In *Findings of the Associ-* 658
ation for Computational Linguistics: EMNLP 2022, 659
pages 5093–5104, Abu Dhabi, United Arab Emirates. 660
Association for Computational Linguistics. 661

Pavel Blinov, Arina Reshetnikova, Aleksandr Nesterov, 662
Galina Zubkova, and Vladimir Kokh. 2022. [Rumed-](#) 663
[bench: a russian medical language understanding](#) 664
[benchmark](#). In *International Conference on Artificial* 665
Intelligence in Medicine, pages 383–392. Springer. 666

Linzheng Chai, Jian Yang, Tao Sun, Hongcheng Guo, 667
Jiaheng Liu, Bing Wang, Xinnian Liang, Jiaqi Bai, 668
Tongliang Li, Qiyao Peng, and 1 others. 2025. [xcot:](#) 669
[Cross-lingual instruction tuning for cross-lingual](#) 670
[chain-of-thought reasoning](#). In *Proceedings of the* 671
AAAI Conference on Artificial Intelligence, vol- 672
ume 39, pages 23550–23558. 673

Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu 674
Lian, and Zheng Liu. 2024a. [Bge m3-embedding:](#) 675
[Multi-lingual, multi-functionality, multi-granularity](#) 676
[text embeddings through self-knowledge distillation](#). 677
arXiv preprint arXiv:2402.03216. 678

Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, 679
Wanlong Liu, Rongsheng Wang, and Benyou Wang. 680
2025. [Towards medical complex reasoning with](#) 681
[LLMs through medical verifiable problems](#). In *Find-* 682
ings of the Association for Computational Linguistics: 683
ACL 2025, pages 14552–14573, Vienna, Austria. As- 684
sociation for Computational Linguistics. 685

Nuo Chen, Zinan Zheng, Ning Wu, Ming Gong, Dong- 686
mei Zhang, and Jia Li. 2024b. [Breaking language](#) 687
[barriers in multilingual mathematical reasoning: In-](#) 688
[sights and observations](#). In *Findings of the Associa-* 689
tion for Computational Linguistics: EMNLP 2024, 690
pages 7001–7016. 691

Changjiang Gao, Xu Huang, Wenhao Zhu, Shujian 692
Huang, Lei Li, and Fei Yuan. 2025a. [Could think-](#) 693
[ing multilingually empower llm reasoning?](#) *arXiv* 694
preprint arXiv:2504.11833. 695

Shanghua Gao, Richard Zhu, Zhenglun Kong, Ayush 696
Noori, Xiaorui Su, Curtis Ginder, Theodoros 697
Tsiligkaridis, and Marinka Zitnik. 2025b. [Txagent:](#) 698
[An ai agent for therapeutic reasoning across a uni-](#) 699
[verse of tools](#). *arXiv preprint arXiv:2503.10970*. 700

811	Charles Nimo, Tobi Olatunji, Abraham Toluwase	Shuaijie She, Wei Zou, Shujian Huang, Wenhao Zhu, Xi-	867
812	Owodunni, Tassallah Abdullahi, Emmanuel Ayodele,	ang Liu, Xiang Geng, and Jiajun Chen. 2024. Mapo:	868
813	Mardhiyah Sanni, Ezinwanne C Aka, Fofafunmi	Advancing multilingual reasoning through multilin-	869
814	Omofoye, Foutse Yuehgo, Timothy Faniran, and	gual alignment-as-preference optimization. <i>arXiv</i>	870
815	1 others. 2025. Afrimed-qa: a pan-african, multi-	<i>preprint arXiv:2401.06838</i> .	871
816	specialty, medical question-answering benchmark		
817	dataset. In <i>Proceedings of the 63rd Annual Meet-</i>	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang,	872
818	<i>ing of the Association for Computational Linguistics</i>	Suraj Srivats, Soroush Vosoughi, Hyung Won Chung,	873
819	(Volume 1: Long Papers), pages 1948–1973.	Yi Tay, Sebastian Ruder, Denny Zhou, and 1 others.	874
820	OpenAI. 2025. Gpt-5.1: A smarter, more conversational	2022. Language models are multilingual chain-of-	875
821	chatgpt . Accessed: 2025-12-31.	thought reasoners. <i>arXiv preprint arXiv:2210.03057</i> .	876
822	Jianhui Pang, Fanghua Ye, Derek Fai Wong, Dian Yu,	Shivalika Singh, Angelika Romanou, Clémentine Four-	877
823	Shuming Shi, Zhaopeng Tu, and Longyue Wang.	rier, David Ifeoluwa Adelani, Jian Gang Ngui, Daniel	878
824	2025. Salute the classic: Revisiting challenges of ma-	Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio,	879
825	chine translation in the age of large language models.	Wei Qi Leong, Yosephine Susanto, and 1 others. 2025.	880
826	<i>Transactions of the Association for Computational</i>	Global mmlu: Understanding and addressing cultural	881
827	<i>Linguistics</i> , 13:73–95.	and linguistic biases in multilingual evaluation. In	882
828	Cheonbok Park, Jeonghoon Kim, Joosung Lee, Sangh-	<i>Proceedings of the 63rd Annual Meeting of the As-</i>	883
829	wan Bae, Jaegul Choo, and Kang Min Yoo. 2025.	<i>sociation for Computational Linguistics (Volume 1:</i>	884
830	Cross-lingual collapse: How language-centric founda-	<i>Long Papers)</i> , pages 18761–18799.	885
831	tion models shape reasoning in large language mod-	Shivalika Singh, Freddie Vargus, Daniel D’souza,	886
832	els. <i>arXiv preprint arXiv:2506.05850</i> .	Börje F Karlsson, Abinaya Mahendiran, Wei-Yin	887
833	Vimla L Patel, José F Arocha, and Jiajie Zhang. 2005.	Ko, Herumb Shandilya, Jay Patel, Deividas Mataci-	888
834	Thinking and reasoning in medicine. <i>The Cambridge</i>	unas, Laura O’Mahony, and 1 others. 2024. Aya	889
835	<i>handbook of thinking and reasoning</i> , 14:727–750.	dataset: An open-access collection for multilingual	890
836	Rui Qi, Zhibo Man, Yufeng Chen, Fengran Mo, Jinan	instruction tuning. In <i>Proceedings of the 62nd An-</i>	891
837	Xu, and Kaiyu Huang. 2025. SoT: Structured-of-	<i>nual Meeting of the Association for Computational</i>	892
838	thought prompting guides multilingual reasoning in	<i>Linguistics (Volume 1: Long Papers)</i> , pages 11521–	893
839	large language models . In <i>Findings of the Associa-</i>	11567.	894
840	<i>tion for Computational Linguistics: EMNLP 2025</i> ,	Guijin Son, Donghun Yang, Hitesh Laxmichand Patel,	895
841	pages 11024–11039, Suzhou, China. Association for	Amit Agarwal, Hyunwoo Ko, Chanuk Lim, Srikant	896
842	Computational Linguistics.	Panda, Minhyuk Kim, Nikunj Drolia, Dasol Choi,	897
843	Pengcheng Qiu, Chaoyi Wu, Xiaoman Zhang, Weixiong	and 1 others. 2025. Pushing on multilingual reason-	898
844	Lin, Haicheng Wang, Ya Zhang, Yanfeng Wang, and	ing models with language-mixed chain-of-thought.	899
845	Weidi Xie. 2024. Towards building multilingual lan-	<i>arXiv preprint arXiv:2510.04230</i> .	900
846	guage model for medicine. <i>Nature Communications</i> ,	Yu Sun, Xingyu Qian, Weiwen Xu, Hao Zhang, Cheng-	901
847	15(1):8384.	hao Xiao, Long Li, Deli Zhao, Wenbing Huang,	902
848	Leonardo Ranaldi and Giulia Pucci. 2025. Multilin-	Tingyang Xu, Qifeng Bai, and 1 others. 2025. Rea-	903
849	gual reasoning via self-training. In <i>Proceedings of</i>	sonmed: A 370k multi-agent generated dataset for	904
850	<i>the 2025 Conference of the Nations of the Americas</i>	advancing medical reasoning. In <i>Proceedings of the</i>	905
851	<i>Chapter of the Association for Computational Lin-</i>	<i>2025 Conference on Empirical Methods in Natural</i>	906
852	<i>guistics: Human Language Technologies (Volume 1:</i>	<i>Language Processing</i> , pages 26457–26478.	907
853	<i>Long Papers)</i> , pages 11566–11582.	Zhi Rui Tam, Cheng-Kuang Wu, Yu Ying Chiu, Chieh-	908
854	Ipek Baris Schlicht, Burcu Sayin, Zhixue Zhao, Fred-	Yen Lin, Yun-Nung Chen, and Hung-yi Lee. 2025.	909
855	erik M Labonté, Cesare Barbera, Marco Viviani,	Language matters: How do multilingual input and	910
856	Paolo Rosso, and Lucie Flek. 2025. Disparities	reasoning paths affect large reasoning models? <i>arXiv</i>	911
857	in multilingual llm-based healthcare q&a. <i>arXiv</i>	<i>preprint arXiv:2505.17407</i> .	912
858	<i>preprint arXiv:2510.17476</i> .	Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming	913
859	Lisa Schut, Yarin Gal, and Sebastian Farquhar. 2025.	Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and	914
860	Do multilingual llms think in english? <i>arXiv preprint</i>	Mark Gerstein. 2024. Medagents: Large language	915
861	<i>arXiv:2502.15603</i> .	models as collaborators for zero-shot medical rea-	916
862	Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri,	soning. In <i>Findings of the Association for Computa-</i>	917
863	Atila Kiraly, Madeleine Traverse, Timo Kohlberger,	<i>tional Linguistics: ACL 2024</i> , pages 599–621.	918
864	Shawn Xu, Fayaz Jamil, Cian Hughes, Charles Lau,	Gemma Team, Morgane Riviere, Shreya Pathak,	919
865	and 1 others. 2025. Medgemma technical report.	Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupati-	920
866	<i>arXiv preprint arXiv:2507.05201</i> .	raju, Léonard Hussenot, Thomas Mesnard, Bobak	921
		Shahriari, Alexandre Ramé, and 1 others. 2024.	922
		Gemma 2: Improving open language models at a	923
		practical size. <i>arXiv preprint arXiv:2408.00118</i> .	924

925	David Vilares and Carlos Gómez-Rodríguez. 2019.	others. 2025. Med-prm: Medical reasoning mod-	980
926	Head-qa: A healthcare dataset for complex reasoning.	els with stepwise, guideline-verified process rewards.	981
927	<i>arXiv preprint arXiv:1906.04701</i> .	In <i>Proceedings of the 2025 Conference on Empiri-</i>	982
928	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,	<i>cal Methods in Natural Language Processing</i> , pages	983
929	Ed Chi, Sharan Narang, Aakanksha Chowdhery, and	16565–16582.	984
930	Denny Zhou. 2022. Self-consistency improves chain	Yuxin Zuo, Shang Qu, Yifei Li, Zhangren Chen, Xuekai	985
931	of thought reasoning in language models. <i>arXiv</i>	Zhu, Ermo Hua, Kaiyan Zhang, Ning Ding, and	986
932	<i>preprint arXiv:2203.11171</i> .	Bowen Zhou. 2025. Medxpertqa: Benchmarking	987
933	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	expert-level medical reasoning and understanding.	988
934	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	<i>arXiv preprint arXiv:2501.18362</i> .	989
935	and 1 others. 2022. Chain-of-thought prompting elic-		
936	its reasoning in large language models. <i>Advances</i>		
937	<i>in neural information processing systems</i> , 35:24824–		
938	24837.		
939	Juncheng Wu, Wenlong Deng, Xingxuan Li, Sheng		
940	Liu, Taomian Mi, Yifan Peng, Ziyang Xu, Yi Liu,		
941	Hyunjin Cho, Chang-In Choi, and 1 others. 2025a.		
942	Medreason: Eliciting factual medical reasoning		
943	steps in llms via knowledge graphs. <i>arXiv preprint</i>		
944	<i>arXiv:2504.00993</i> .		
945	Kevin Wu, Eric Wu, Rahul Thapa, Kevin Wei, Angela		
946	Zhang, Arvind Suresh, Jacqueline J Tao, Min Woo		
947	Sun, Alejandro Lozano, and James Zou. 2025b. Med-		
948	casereasoning: Evaluating and learning diagnostic		
949	reasoning from clinical case reports. <i>arXiv preprint</i>		
950	<i>arXiv:2505.11733</i> .		
951	Yunfei Xie, Juncheng Wu, Haoqin Tu, Siwei Yang,		
952	Bingchen Zhao, Yongshuo Zong, Qiao Jin, Cihang		
953	Xie, and Yuyin Zhou. 2024. A preliminary study of		
954	o1 in medicine: Are we closer to an ai doctor? <i>arXiv</i>		
955	<i>preprint arXiv:2409.15277</i> .		
956	Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong		
957	Zhang. 2024. Benchmarking retrieval-augmented		
958	generation for medicine. In <i>Findings of the Associa-</i>		
959	<i>tion for Computational Linguistics ACL 2024</i> , pages		
960	6233–6251.		
961	Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen,		
962	Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao,		
963	and Tianyi Zhou. 2024. A survey on knowledge dis-		
964	stillation of large language models. <i>arXiv preprint</i>		
965	<i>arXiv:2402.13116</i> .		
966	Weihao Xuan, Rui Yang, Heli Qi, Qingcheng Zeng,		
967	Yunze Xiao, Aosong Feng, Dairui Liu, Yun Xing,		
968	Junjue Wang, Fan Gao, and 1 others. 2025. Mmlu-		
969	prox: A multilingual benchmark for advanced		
970	large language model evaluation. <i>arXiv preprint</i>		
971	<i>arXiv:2503.10497</i> .		
972	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,		
973	Binyuan Hui, Bo Zheng, Bowen Yu, Chang		
974	Gao, Chengen Huang, Chenxu Lv, and 1 others.		
975	2025. Qwen3 technical report. <i>arXiv preprint</i>		
976	<i>arXiv:2505.09388</i> .		
977	Jaehoon Yun, Jiwoong Sohn, Jungwoo Park, Hyunjae		
978	Kim, Xiangru Tang, Daniel Shao, Yong Hoe Koo,		
979	Ko Minhyeok, Qingyu Chen, Mark Gerstein, and 1		

990 A Position-aware Concept Fusion 991 Strategy

992 Algorithm 1 describes our position-aware cross-
993 lingual concept fusion mechanism in detail. Given
994 an English concept chain C_e and a local language
995 concept chain C_l , we iteratively integrate local
996 concepts into the fused chain C_f (initialized as C_e).
997 For each local concept c_l^j , we compute its embed-
998 ding e_l and identify the position k^* of the most sim-
999 ilar concept in the current fused chain using cosine
1000 similarity. If the maximum similarity s_{\max} exceeds
1001 the threshold τ , we proceed to determine the inser-
1002 tion position. Rather than simply appending the
1003 concept, we compare the average similarity of c_l^j
1004 to all concepts positioned before k^* (s_{left}) versus
1005 those after k^* (s_{right}). This bidirectional context
1006 comparison ensures that the inserted concept is po-
1007 sitioned where it exhibits the strongest semantic
1008 coherence with surrounding concepts, thereby pre-
1009 serving the logical structure and clinical reasoning
flow of the chain.

Algorithm 1: Position-Aware Concept Fu- sion

Input: English Concept Chain

$C_e = \{c_e^1, \dots, c_e^n\}$, Local Concept

Chain $C_l = \{c_l^1, \dots, c_l^m\}$,

Embedding function f , Threshold τ

Output: Fused concept chain C_f

```

1  $C_f \leftarrow C_e$ ;
2  $E_f \leftarrow \{f(c_e^1), \dots, f(c_e^n)\}$ ;
3 for  $c_l^j \in C_l$  do
4    $e_l \leftarrow f(c_l^j)$ ;
5    $k^* \leftarrow \operatorname{argmax}_{k \in [1, |C_f|]} \cos(e_l, E_f^k)$ ;
6    $s_{\max} \leftarrow \cos(e_l, E_f^{k^*})$ ;
7   if  $s_{\max} \geq \tau$  then
8      $s_{left} \leftarrow \frac{1}{k^* - 1} \sum_{i=1}^{k^* - 1} \cos(e_l, E_f^i)$ ;
9      $s_{right} \leftarrow$ 
10       $\frac{1}{|C_f| - k^*} \sum_{i=k^* + 1}^{|C_f|} \cos(e_l, E_f^i)$ ;
11     if  $s_{left} > s_{right}$  then
12        $p \leftarrow k^*$ ;
13     else
14        $p \leftarrow k^* + 1$ ;
15     end
16     Insert  $c_l^j$  into  $C_f$  at position  $p$ ;
17     Insert  $e_l$  into  $E_f$  at position  $p$ ;
18 end
19 return  $C_f$ ;

```

B Overall Results on MultiMed-X

Overall Performance. Table 4 summarizes re-
sults on MultiMed-X across languages and tasks.
MED-COREASONER achieves the best or near-best
performance across all languages, outperforming
strong baselines on long-form QA metrics, includ-
ing Overall, Correctness, Completeness, and Pass
Rate, while maintaining high Safety and low Hallu-
cination. Consistent gains in NLI accuracy further
demonstrate the effectiveness of cross-lingual co-
reasoning and knowledge grounding for reliable
multilingual medical decision-making.

Cross-lingual and Low-resource Analysis. A
closer look reveals that the advantages of MED-
COREASONER are particularly pronounced in
low-resource languages such as Swahili, Yoruba,
and Zulu. Compared to direct prompting, MED-
COREASONER yields larger improvements in com-
pleteness, hallucination control, and pass rate, ef-
fectively narrowing the performance gap between
English and underrepresented languages. This
trend highlights the effectiveness of pivot-anchored
co-reasoning in preserving logical structure while
incorporating localized clinical knowledge, leading
to more robust and equitable multilingual medical
reasoning.

C Reasoning Training Data

The comparative statistics of the MMedBench train-
ing set and our dataset are shown in Table 5. Using
MED-COREASONER, we generate forward reason-
ing by inputting training questions and obtaining
corresponding reasoning chains and answers. We
only use those instances with correct final answers
for training.

D Full Results on MMedBench

We include multiple large language models
pre-trained specifically for the medical do-
main on MMedBench for comparison, including
BioMistral-7B (Labrak et al., 2024), MMedLM2-
7B (Chen et al., 2025), and MedGemma-27B (Sel-
lergren et al., 2025). Full results are reported in
Table 6.

Overall Comparison. Table 6 compares mul-
tilingual performance on MMedBench across
medical-domain and general-purpose LLMs.
Among domain-specific models, MedGemma-27B

Language	Model	Long-form QA						NLI
		Overall	Correctness	Completeness	Safety	Hallucination	Pass Rate	Accuracy (%)
EN	GPT-5.2	4.415	4.455	4.615	4.815	4.410	0.900	76.67
	GPT-5.1	4.425	4.485	4.590	4.880	4.385	0.915	76.67
	GPT-4o	4.210	4.270	4.245	4.845	4.380	0.855	78.67
	Claude-3.5-haiku	4.200	4.240	4.265	4.830	4.305	0.870	69.33
	Ours (GPT-5.1)	4.600	4.640	4.850	4.860	4.460	0.939	77.33
ZH	GPT-5.2	4.420	4.460	4.605	4.800	4.360	0.915	74.00
	GPT-5.1	4.435	4.490	4.660	4.845	4.395	0.905	74.67
	GPT-4o	4.205	4.270	4.270	4.855	4.250	0.860	74.67
	Claude-3.5-haiku	4.150	4.205	4.295	4.735	4.140	0.845	69.33
	Ours (GPT-5.1)	4.530	4.550	4.890	4.780	4.390	0.935	76.67
JP	GPT-5.2	4.335	4.405	4.520	4.855	4.330	0.885	76.67
	GPT-5.1	4.330	4.435	4.525	4.865	4.330	0.895	77.33
	GPT-4o	4.080	4.135	4.135	4.830	4.145	0.820	76.00
	Claude-3.5-haiku	4.100	4.205	4.170	4.860	4.095	0.845	68.67
	Ours (GPT-5.1)	4.430	4.500	4.850	4.800	4.330	0.930	78.00
KO	GPT-5.2	4.410	4.460	4.655	4.815	4.325	0.915	64.67
	GPT-5.1	4.390	4.475	4.610	4.795	4.330	0.890	62.67
	GPT-4o	4.055	4.130	4.105	4.805	4.225	0.840	61.33
	Claude-3.5-haiku	4.115	4.170	4.240	4.795	4.110	0.860	60.00
	Ours (GPT-5.1)	4.540	4.570	4.840	4.790	4.450	0.920	67.33
SW	GPT-5.2	4.340	4.415	4.590	4.795	4.345	0.890	70.00
	GPT-5.1	4.330	4.385	4.545	4.805	4.255	0.885	72.67
	GPT-4o	4.040	4.070	4.115	4.755	4.155	0.815	69.33
	Claude-3.5-haiku	3.825	3.855	3.950	4.565	3.780	0.730	56.67
	Ours (GPT-5.1)	4.550	4.570	4.800	4.820	4.470	0.945	73.33
TH	GPT-5.2	4.440	4.490	4.605	4.835	4.430	0.900	68.00
	GPT-5.1	4.490	4.535	4.645	4.830	4.545	0.910	70.00
	GPT-4o	4.160	4.215	4.180	4.890	4.280	0.855	70.00
	Claude-3.5-haiku	4.135	4.205	4.200	4.790	4.330	0.855	63.33
	Ours (GPT-5.1)	4.660	4.670	4.880	4.910	4.620	0.960	71.33
YO	GPT-5.2	4.065	4.135	4.225	4.550	4.060	0.795	62.00
	GPT-5.1	4.090	4.160	4.310	4.595	4.080	0.810	63.33
	GPT-4o	3.290	3.325	3.400	4.095	3.205	0.455	62.67
	Claude-3.5-haiku	3.545	3.555	3.650	4.435	3.585	0.600	54.67
	Ours (GPT-5.1)	4.450	4.460	4.790	4.700	4.230	0.910	68.67
ZU	GPT-5.2	4.210	4.280	4.440	4.715	4.195	0.860	66.00
	GPT-5.1	4.160	4.230	4.445	4.695	4.145	0.850	64.00
	GPT-4o	3.780	3.805	3.875	4.605	3.885	0.725	66.00
	Claude-3.5-haiku	3.670	3.705	3.865	4.505	3.685	0.645	62.00
	Ours (GPT-5.1)	4.420	4.710	4.770	4.730	4.310	0.915	68.00

Table 4: Complete evaluation results across different languages on MultiMed-X.

Train Set	Chinese	English	French	Japanese	Russian	Spanish
MMedBench	27,400	10,178	2,171	1,590	1,052	2,656
<i>MMed-Reason</i>	8,627	9,513	1,603	1,392	846	2,487

Table 5: Training data statistics of MMedBench and *MMed-Reason*

1056 achieves the strongest average performance (65.88),
1057 outperforming BioMistral-7B and MMedLM2-7B,
1058 but still exhibits notable variance across languages,
1059 particularly weaker results in French and Japanese.
1060 This suggests that medical pre-training alone does
1061 not guarantee robust multilingual generalization.

1062 **Effect of Training Data and Model Scale.** For
1063 general-purpose models fine-tuned on different
1064 datasets, training on *MMed-Reason* consistently
1065 improves multilingual performance compared to
1066 MMedBench across model scales. In particular,
1067 Qwen2.5-14B trained on *MMed-Reason* achieves
1068 the best overall average score (75.97), with clear
1069 gains across all non-English languages and espe-
1070 cially large improvements in French and Japanese.
1071 Similar trends are observed for Qwen2.5-7B and
1072 Gemma-7B-it, indicating that *MMed-Reason* pro-
1073 vides more effective cross-lingual medical supervi-
1074 sion and that performance gains scale with model
1075 capacity.

1076 E Expert Evaluation.

1077 We recruit all the expert physicians through so-
1078 cial media. For the reasoning quality assessment
1079 experiment, we randomly sample questions in
1080 Japanese, Spanish, and Chinese from the MMLU-
1081 ProX benchmark, and generate reasoning and an-
1082 swers using GPT-5.1 and MED-COREASONER.
1083 For fairness, we retain only the cases where both
1084 models produce correct answers, resulting in 30
1085 question–answer pairs. The evaluation guidelines
1086 provided to physician experts are shown in Figure 6
1087 and the example pairwise evaluation are shown in
1088 Table 8.

1089 F Implementation Details

1090 We provide all hyperparameters and experimental
1091 settings in this section.

1092 **Prompts.** For parallel reasoning and concept ex-
1093 traction, we use the prompts shown in Figures 7
1094 and 8. Final answer generation is performed us-
1095 ing the prompt in Figure 9. For LLM-as-a-judge
1096 evaluation in long-form QA, we adopt the system
1097 prompt in Figure 10 together with the evaluation
1098 prompt in Figure 11.

1099 **Knowledge Retrieval.** We construct language-
1100 specific medical knowledge bases from MSD Man-
1101 uals and AFRIDOC-MT. Detailed statistics of the
1102 documents for each language are reported in Ta-
1103 ble 7. Given a query in a particular language, we

1104 retrieve relevant documents from the correspond-
1105 ing language-specific knowledge base. We use
1106 BGE-M3 as the retriever and reranker, retrieving
1107 the top 10 documents in the initial retrieval stage
1108 and reranking the top 3 documents for final use.

Model	Train Set	Chinese	English	French	Japanese	Russian	Spanish	Avg.
BioMistral-7B	-	25.89	19.17	10.13	8.54	54.3	25.67	24.66
MMedLM2-7B	-	70.43	58.13	54.27	38.26	71.88	64.95	59.32
MedGemma-27B	-	73.50	71.09	41.16	60.08	72.27	79.72	65.88
Gemma-7B-it	MMedBench	56.07	52.16	34.89	34.67	63.28	57.51	49.45
	<i>MMed-Reason</i>	52.38	52.08	40.03	41.71	64.45	58.06	51.39
Qwen2.5-7B	MMedBench	81.47	61.67	47.72	48.74	69.53	67.18	62.16
	<i>MMed-Reason</i>	78.78	62.37	54.66	56.78	70.31	69.22	64.98
Qwen2.5-14B	MMedBench	84.47	71.48	64.15	66.33	75.39	78.77	73.11
	<i>MMed-Reason</i>	82.89	75.73	72.03	68.34	75.78	82.39	75.97

Table 6: Performance comparison across languages on MMedBench.

EN	ZH	JA	KO	DE	FR	ES	IT	SW	YO	ZU
2,441	2,857	2,502	3,428	2,855	3,044	2,943	2,960	3,491	1,148	1,148

Table 7: Document statistics of multilingual knowledge base.

Pairwise Comparison Guidelines

Introduction: We want to verify which “reasoning” path is more reasonable. Here, reasoning refers to the structured sequence of diagnostic or decision-making steps that link clinical evidence to a conclusion, analogous to clinical reasoning in medical practice.

Task: Pairwise Comparison

Instruction:

You will be shown a clinical question and two reasoning explanations (A and B) for the same case. Do not judge only based on the final answer, but focus on the reasoning quality.

Please evaluate which is better (or tie) considering the following dimensions:

- **Reasoning Clarity:** Which reasoning or explanation is more logically organized?
- **Clinical Soundness:** Which reasoning is more medically reasonable?
- **Hallucination/Safety:** Which reasoning is less likely to mislead clinical judgement?
- **Local Clinical Naturalness:** Which reasoning sounds more natural in the local clinical context?

Please select A, B, or Tie in each evaluation dimension.

Figure 6: Physician expert pairwise comparison guidelines.

Parallel Reasoning Prompt

Your task is to assist healthcare professionals in clinical reasoning by providing well-thought-out answers to medical questions. Please first think step by step using the language language and then provide your final answer. Your response will be used for research purpose only, so please provide a definite answer (e.g., A, B, C, or D).

****Question**:**

{question}

****Options**:**

{options}

****Output Format: ****

Please provide your reasoning process in a step-by-step manner using {language} language, followed by your final answer. Use the following format:

```
{{ "reasoning": "Your detailed reasoning process here", "answer": "Your definite answer here, e.g., A, B, C, or D" }}
```

Figure 7: Reasoning Prompt

Question	A_Reasoning	B_Reasoning
<p>Observa una lesión en la mucosa alveolar labial del incisivo lateral inferior derecho de uno de sus pacientes y decide tomar una biopsia para enviar al informe de patología oral. ¿Qué nervio requeriría anestesia local para realizar una biopsia sin dolor?</p> <p>Options: A: El nervio mentoniano; B: El nervio hipogloso; C: El nervio incisivo; D: El nervio bucal; E: El nervio alveolar inferior; F: El nervio infraorbitario; G: El nervio facial; H: El nervio maxilar; I: El nervio lingual; J: El nervio trigémino</p>	<p>La lesión se localiza en la mucosa alveolar labial del incisivo lateral inferior derecho. La sensibilidad de la mucosa y piel del labio inferior y región del mentón está dada por el nervio mentoniano, que es la rama terminal del nervio alveolar inferior tras salir por el foramen mentoniano. El nervio incisivo inerva principalmente los dientes anteriores (pulpa y tejidos periapicales) pero no la mucosa labial. Por tanto, para una biopsia de la mucosa alveolar labial en esa zona, el nervio que debe anestesiarse específicamente es el nervio mentoniano.</p>	<p>La lesión se localiza en la mucosa alveolar labial a nivel del incisivo lateral inferior derecho. La sensibilidad de la mucosa/encía vestibular (labial) de los incisivos, canino y primer premolar mandibulares está dada por el nervio mentoniano, que es la rama terminal del nervio alveolar inferior.</p> <p>En el contexto aportado se describe que el nervio alveolar inferior, rama del nervio mandibular (V3), discurre dentro de la mandíbula y da una rama terminal, el nervio mentoniano, que emerge por el foramen mentoniano en la región de los premolares e inerva el labio inferior, el mentón y la mucosa/gingiva lateral que recubre los dientes incisivos inferiores, canino y primer premolar. Tras emitir el nervio mentoniano, el nervio alveolar inferior continúa hacia anterior como nervio incisivo, que inerva los dientes y el hueso anteriores al foramen mentoniano, pero no la mucosa labial externa.</p> <p>Por lo tanto, para realizar una biopsia indolora de la mucosa alveolar labial del incisivo lateral inferior, debe anestesiarse el nervio mentoniano, no el nervio incisivo ni el alveolar inferior de forma proximal.</p> <p>Las demás opciones corresponden a nervios con territorios de inervación diferentes (lingual, bucal, infraorbitario, maxilar, facial, hipogloso, trigémino completo) y no son los responsables de la sensibilidad en esa zona específica.</p>
<p>La mayoría de los ácidos grasos libres se transportan en la sangre:</p> <p>Options: A: dentro de los glóbulos rojos.; B: combinados con oxígeno.; C: como triglicéridos.; D: como moléculas libres.; E: unidos a los glóbulos blancos.; F: como lipoproteínas.; G: unidos a la hemoglobina.; H: combinados con glucosa.; I: en forma de colesterol.; J: unidos a la albúmina.</p>	<p>Los ácidos grasos libres, una vez liberados del tejido adiposo al torrente sanguíneo, son moléculas hidrofóbicas que no se disuelven bien en el plasma acuoso. Para poder transportarse de forma eficiente y segura en la sangre, se unen de manera no covalente principalmente a la albúmina sérica. No viajan como triglicéridos, colesterol ni como parte de lipoproteínas; estas formas se utilizan sobre todo para el transporte de lípidos más complejos (triglicéridos y colesterol esterificado), no de ácidos grasos libres. Tampoco se unen a eritrocitos, leucocitos, hemoglobina, glucosa ni oxígeno. Por lo tanto, la forma principal de transporte de los ácidos grasos libres en la sangre es unidos a la albúmina.</p>	<p>Los ácidos grasos libres (no esterificados) son moléculas hidrofóbicas, por lo que no pueden circular libremente disueltas en el plasma, que es un medio acuoso. Para poder ser transportados en sangre, deben unirse a proteínas transportadoras. En el caso específico de los ácidos grasos libres, el principal transportador en el plasma es la albúmina sérica, que posee sitios de unión para estos ácidos grasos y los lleva desde el tejido adiposo hacia otros tejidos (como músculo e hígado) donde se utilizan como fuente de energía. No se transportan mayoritariamente dentro de eritrocitos, ni unidos a hemoglobina, leucocitos, oxígeno, glucosa, ni formando triglicéridos, colesterol o lipoproteínas (estas estructuras transportan principalmente triglicéridos y colesterol, no los ácidos grasos libres como tales). Por tanto, la forma principal de transporte de los ácidos grasos libres en sangre es unidos a la albúmina.</p>

Table 8: Example of pairwise reasoning comparison in Spanish. In both cases, A_reasoning corresponds to the GPT-5.1 baseline, while B_reasoning represents reasoning generated by MED-COREASONER with GPT-5.1 as the backbone.

Concept Extraction Prompt

Your task is to transform a reasoning trace into a concise, ordered **concept chain**.

Definitions

- "A concept" is an atomic, reusable medical idea or clinical finding that contributes to clinical reasoning. (e.g., "chronic cough", "chest X-ray", "smoker's history")
- "Concept chain" is an ordered list of these concepts that follows the original reasoning order.

Instructions

- Read the reasoning trace carefully.
- Extract key concepts in the same order as they appear
- Each concept must be concise and represent only one idea
- Prefer clinical or scientific terms over long sentences
- Do not invent new concepts that are not implied by the reasoning
- Merge duplicates or near-duplicates into one concept, but keep the order consistent
- The concept chain should be as short as possible while capturing all essential reasoning steps, and must not include or expose any answer options.
- Keep the output in the same language as the reasoning.

Output Format (plain text, no explanation):

Provide the concept chain as a list in the following format:

[*"Concept1"*, *"Concept2"*, *"Concept3"*, ...]

Now process the following reasoning trace and output only the concept chain in {language} language.

{reasoning_trace}

Figure 8: Concept Extraction Prompt

Final Answer Generation Prompt

Your task is to generate a final clinical answer for a multi-option question by integrating a **concept reasoning chain** with **the retrieved medical documents** from the concept chain and your **prior medical knowledge**.

Inputs

- Question:

{question}

- Options:

{options}

- Concept Reasoning Chain (in order):

{concept_chain}

- Referenced Context (mainly contains multiple documents, guidelines, or passages):

{context}

Instructions

1. First, carefully read the concept reasoning chain. Treat it as a DRAFT reasoning path, not as guaranteed truth.
2. Then, carefully read the referenced context. Use it to VERIFY, CORRECT, or REFINE the reasoning chain.
3. Use ONLY information that is supported by the referenced context and widely accepted medical knowledge. DO NOT directly mention the concept chain. ORGANIZE your reasoning in a clear, logical manner.
4. Finally, select the MOST APPROPRIATE option as your final answer based on the verified and refined reasoning.
5. Output the reasoning in language, regardless of the input language.

Output Format:

Return VALID JSON ONLY, following this format: { { "reasoning": "Your verified and refined reasoning process here", "answer": "Your final answer here, e.g., A, B, C, or D" } }

Figure 9: Final Answer Generation Prompt

Judge System Prompt

You are an objective and rigorous evaluator for medical question answering.

You will be given:

- a Question
- a Ground-Truth Answer (reference)
- a Model Answer (candidate)

Your task is to evaluate the Model Answer relative to the Ground-Truth Answer.

Evaluation principles:

Prioritize factual correctness, clinical safety, and alignment with the reference.

Do NOT penalize harmless extra details if they are correct and do not contradict the reference.

Penalize contradictions, fabricated facts, or unsafe medical advice.

If the reference is brief but the model answer is longer, judge consistency and medical plausibility.

If a detail cannot be verified from the reference, treat it as uncertain rather than incorrect.

Output MUST be valid JSON only

Figure 10: The system prompt of LLM-as-a-judge in the evaluation of long-form QA task.

Judge Evaluation Prompt

Question:

{question}

Ground-Truth Answer:

{gold}

Model Answer:

{pred}

Return JSON with EXACTLY the following fields and no others:

```
{{
```

```
"overall_score": 1-5,
```

```
"correctness": 1-5,
```

```
"completeness": 1-5,
```

```
"safety": 1-5,
```

```
"hallucination": 1-5
```

```
}}
```

Scoring rules:

- 5 = excellent

- 4 = good with minor issues

- 3 = partially correct or incomplete

- 2 = major issues

- 1 = mostly incorrect or unsafe

For hallucination:

- 5 = no hallucination

- 3 = some uncertain additions

- 1 = clear hallucinations or fabricated facts

Figure 11: The evaluation prompt of LLM-as-a-judge in the evaluation of long-form QA task.