

HiMed: Incentivizing Hindi Reasoning in Medical LLMs

Anonymous ACL submission

Abstract

Medical large language models hold promise for reducing healthcare disparities, yet Hindi remains severely underrepresented. While medical LLMs excel in high-resource languages, their performance degrades sharply in Hindi, particularly on Indian systems of medicine. We argue that robust cross-lingual medical transfer requires Hindi reasoning. To this end, we introduce **HiMed**, a comprehensive Hindi reasoning medical corpus and benchmark suite covering both Western and Indian medicine. We further propose **HiMed-8B**, a Hindi medical reasoning LLM based on LLaMA-3.1-8B-Instruct, through the design of decaying scaffolding reward. Extensive experiments have been conducted, demonstrating consistent improvements in Hindi medical reasoning performance and substantial reduction in the English–Hindi accuracy gap. Ablation studies further validate the contribution of each training stage and reward component. All data and code are available at an anonymous GitHub repository: <https://anonymous.4open.science/r/anon-repo-54EC/README.md>.

1 Introduction

Large language models (LLMs) have revolutionized medical applications, demonstrating strong performance in high-resource languages (Wang et al., 2025a; Lawrence et al., 2024). Medical decision-making is inherently reasoning-centric: life-critical judgments require causal analysis and factual consistency (Xu et al., 2024). Recent studies further show that strengthening reasoning capabilities of medical LLMs can substantially improve decision accuracy (Chen et al., 2025; Wu et al., 2025a), underscoring the importance of reasoning.

Beyond accuracy, medical reasoning LLMs also have the potential to reduce healthcare disparities by providing scalable clinical expertise to underserved populations across languages (Chinta et al.,

2024; d’Elia et al., 2022). However, crossing linguistic boundaries introduces an additional challenge: models must not only transfer reasoning competence, but also align with the linguistic structures and cultural paradigms of the target language (Veselovsky et al., 2025; Sakai et al., 2025).

This challenge is especially pronounced in Hindi-speaking regions, where translation-based deployment exposes hallucination risks. Meanwhile, India’s medical ecosystem requires both Western and Indian systems of medicine (Chandra and Patwardhan, 2018). Many core concepts in Indian systems of medicine lack English equivalents. Together, these factors indicate that Hindi medical reasoning benefits from operating directly within Hindi linguistic representations, motivating **Hindi reasoning** as a necessary objective.

In this work, we address Hindi medical reasoning at both the method and data levels. We propose **Decaying Scaffolding Reward Reinforcement Learning**, which progressively shifts optimization from reasoning behavior guidance to task-optimal objectives. To support this design and enable systematic evaluation, we introduce **HiMed**, a comprehensive Hindi medical corpus and benchmark suite spanning both Western and Indian systems of medicine. Experiments based on LLaMA-3.1-8B-Instruct (Team, 2024) produce the **HiMed-8B** model, which shows consistent improvements across medical benchmarks, with ablations validating the contribution of each training stage and reward component. An overview of HiMed and the training framework is shown in Figure 1.

Overall, this work identifies Hindi reasoning as a first-class objective and provides both the data foundation and training methodology. By jointly addressing a series of challenges, our study offers a systematic and extensible framework for adapting general-purpose LLMs to linguistically and culturally specialized medical domains.

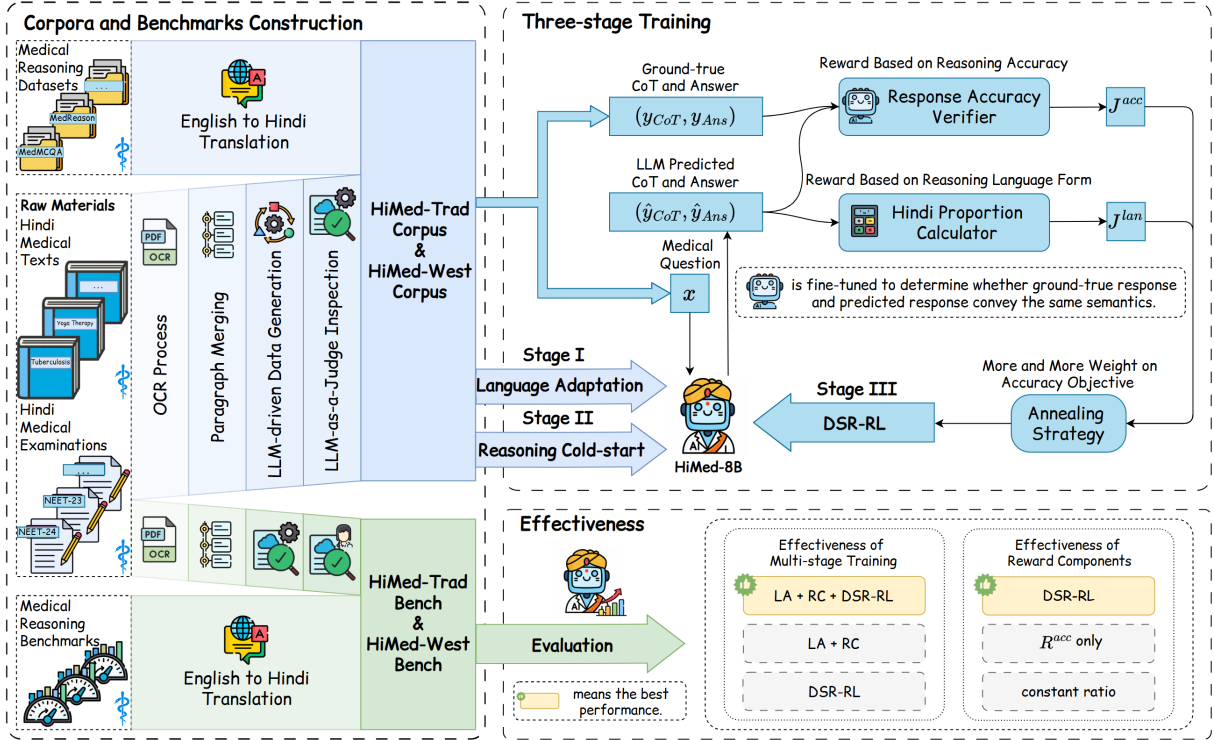


Figure 1: Overview of HiMed and the training framework.

2 Background and Motivation

2.1 Hindi Medical Reasoning Gaps

Evaluation Setup. Prior work has reported substantial cross-lingual degradation and poor alignment with Indian medical knowledge in LLMs (Jin et al., 2024; Qiu et al., 2024; Zhou et al., 2025; Wang et al., 2025b; Devane et al., 2025). To empirically characterize these gaps, we conduct a pilot evaluation on two representative medical benchmarks reflecting distinct medical paradigms: (i) the medical subsets of MMLU-Pro representing Western medicine, and (ii) RET, an Indian medical examination focusing on Indian systems of medicine. We evaluate GPT-4o (Hurst et al., 2024), HuatuoGPT-o1-8B (Chen et al., 2025), and LLaMA3-8B-UltraMedical (Zhang et al., 2024) on both benchmarks in English and Hindi. Δ denotes the accuracy drop from English to Hindi. Full evaluation details are provided in Appendix A.

Benchmark	GPT-4o	HUATUO-O1	ULTRAMED
MMLU-Pro-En	60.7%	64.7%	61.3%
MMLU-Pro-Hi	50.7%	41.3%	34.0%
MMLU-Pro- Δ	10.0%	23.4%	27.3%
RET-En	88.0%	38.0%	52.0%
RET-Hi	70.0%	26.0%	28.0%
RET- Δ	18.0%	12.0%	24.0%

Table 1: Pilot evaluation of cross-lingual performance.

Results. As shown in Table 1, all evaluated models exhibit substantial English–Hindi accuracy drops, exposing a clear cross-lingual reasoning gap. The gap remains severe for Indian systems of medicine: on RET, open-source models fall to near-random performance in Hindi despite relatively strong English results. These results suggest that medical reasoning acquired predominantly through English-centric training does not consistently generalize to Hindi, especially in settings where effective reasoning depends on culturally grounded medical concepts and practices.

Findings 1. *LLMs exhibit a clear English–Hindi performance gap in both Western medicine and Indian systems of medicine.*

Implication. Taken together, this Hindi performance gap points to a deeper limitation in cross-lingual reasoning transfer: models may produce fluent Hindi outputs without actually reasoning in Hindi, revealing a disconnect between language form and reasoning process. Prior work further shows that cross-lingual reasoning degrades without native-language reasoning supervision (Park et al., 2025; Barua et al., 2025). These observations collectively motivate treating **Hindi reasoning** as a first-class objective.

Resource	Reference	Type	Language	Size	Native Hindi	Western Med.	Indian Med.	CoT
MedMCQA	(Pal et al., 2022)	Corpus/Benchmark	English	193K	No	Yes	No	Yes
BioASQ-QA	(Krithara et al., 2023)	Corpus/Benchmark	English	5K	No	Yes	No	No
PubMedQA	(Jin et al., 2019)	Corpus/Benchmark	English	273K	No	Yes	No	No
MedNLI	(Romanov and Shivade, 2018)	Corpus/Benchmark	English	14K	No	Yes	No	No
ReasonMed	(Sun et al., 2025)	Corpus	English	1.11M	No	Yes	No	Yes
MedReason	(Wu et al., 2025a)	Corpus	English	32K	No	Yes	No	Yes
IndicInstruct	(Gala et al., 2024)	Corpus	Hindi + English	404K	No	No	No	No
UltraMedical	(Zhang et al., 2024)	Corpus	English	410K	No	Yes	No	No
AI4Bharat-IndicQA	(Kunchukuttan et al., 2020)	Corpus	Hindi + Indic	18K	Yes	No	No	No
XLingHealth	(Lab, 2024)	Corpus	Hindi + English	15K	No	Yes	No	No
MILU-cleaned	(Murthyudra, 2025)	Corpus	Hindi + Indic	74K	Yes	Yes	No	No
AyurGenixAI	(kagglekirti123, 2025)	Corpus	Primarily English	15K	Yes	No	Yes	No
Multilingual Healthcare	(Bagga, 2025)	Corpus	Hindi + English	4K	No	Yes	No	No
MedQA-USMLE	(Jin et al., 2020)	Benchmark	English	1K	No	Yes	No	-
MMLU-Pro-Med	(Wang et al., 2024b)	Benchmark	English	1K	No	Yes	No	-
HLE	(Phan et al., 2025)	Benchmark	English	2K	No	Yes	No	-
HealthBench	(Arora et al., 2025)	Benchmark	English	5K	No	Yes	No	-
GPQA-med	(Rein et al., 2023)	Benchmark	English	0.2K	No	Yes	No	-
MedXpertQA	(Zuo et al., 2025)	Benchmark	English	2K	No	Yes	No	-
BhashaBench-Ayur	(Devane et al., 2025)	Benchmark	Hindi + English	14K	Yes	No	Yes	-
HiMed-West Corpus	-	Corpus	Hindi	116K	No	Yes	No	Yes
HiMed-Trad Corpus	-	Corpus	Hindi	286K	Yes	No	Yes	Yes
HiMed-West Bench	-	Benchmark	Hindi	2K	No	Yes	No	-
HiMed-Trad Bench	-	Benchmark	Hindi	6K	Yes	No	Yes	-

Table 2: Comparison of existing medical corpora and benchmarks.

2.2 Why Translation Is Insufficient for Hindi Medical Reasoning

2.2.1 Reason I: Translation Hallucination

Evaluation Setup. A possible route to Hindi medical LLMs is translation-based deployment, which translates Hindi inputs into English for reasoning and back into Hindi for delivery. However, such methods are unreliable for terminology-dense Hindi medical content. To quantify this risk, we ask GPT-4o to translate 50 English medical sentences into Hindi and have expert translators conduct professional post-editing.

Case I: Translation Hallucination

English Source.

A 15-year-old **African-American male** with a BMI of 22 is brought to his physician by his mother to address concerns about a change in his dietary habits.

GPT Translation.

एक 15 वर्षीय **आफ्रिकान द्वीपानोसोमा संक्रमणज रोग पुरुष** जिसका BMI 22 है, उसकी माता उसे उसके पथ्यसम्बन्धी आदतों में परिवर्तन के बारे में चिंता व्यक्त करने के लिए उसके काया चिकित्सक के पास लाती हैं।

Corresponding English.

(A 15-year-old **African trypanosomal infectious disease male** with a BMI of 22 is brought by his mother to his body physician in order to express concern about a change in his dietary habits.)

Findings 2. *LLM-based translation can introduce semantic hallucinations in terminology-dense Hindi medical content.*

Results and implication. We find that 19 out of 50 translated sentences contain translation-induced semantic hallucinations, including mis-

translated medical terms and hallucinated medical information absent from the source. The above example shows how translation-induced hallucinations can silently corrupt the original semantics. Consequently, Hindi medical LLMs must support **Hindi reasoning**, allowing valid semantic structures to be formed directly in the target language.

2.2.2 Reason II: Cultural Grounding

Domain Context. Beyond translation limitations, Hindi reasoning reflects the reality of medical practice in India, where Western medicine and Indian systems of medicine coexist as complementary paradigms (World Health Organization, 2010b; Woodyard, 2011). Medical reasoning is shaped by language and culture, and prior work shows that language-concordant care improves safety and diagnostic effectiveness (Molina and Kasper, 2019; Reaume et al., 2025). Moreover, many core concepts in Indian systems of medicine are linguistically grounded and lack faithful English equivalents.

Findings 3. *Medical reasoning in Indian systems relies on linguistic and cultural grounding that is not preserved under English-centric reasoning.*

Limitation. Existing medical LLMs are predominantly shaped by Western biomedical assumptions and English-centric training data (Pfohl et al., 2024). This limits their ability to reason over Indian systems of medicine in a conceptually faithful manner. Hindi reasoning is not a linguistic preference, but a prerequisite for faithful medical reasoning in the Indian context.

2.3 Why Is Hindi Reasoning Challenging?

Developing effective Hindi medical reasoning models involves two fundamental bottlenecks: a data bottleneck and an optimization bottleneck.

Data bottleneck. As summarized in Table 2, existing medical corpora and benchmarks provide limited support for Hindi medical reasoning. Most resources are either English-only or lack coverage of Indian systems of medicine. Consequently, models lack sufficient training signals to learn Hindi medical reasoning behaviors.

Optimization bottleneck. Beyond data scarcity, Hindi medical reasoning faces a fundamental optimization challenge. Under Hindi prompts, task-level accuracy provides no supervision on how reasoning should be expressed, allowing multiple reasoning styles to be equally optimal. In the following example, both responses select the same correct option, yet differ substantially in how medical reasoning is organized and expressed.

Case II: Reasoning Styles

Question.

एक रोगी में तेज़ बुखार, गर्दन में अकड़न और प्रकाश से डर पाया गया है। सबसे संभावित निदान क्या है?

(A patient presents with high fever, neck stiffness, and fear of light. What is the most probable diagnosis?)

(A) मस्तिष्क ट्यूमर (Brain tumor) (B) मेनिन्जाइटिस (Meningitis) (C) मिर्गी (Epilepsy) (D) माइग्रेन (Migraine)

Reasoning A: English answer-driven.

The correct answer is (B), as these symptoms are consistent with meningitis.

Reasoning B: Hindi symptom-to-mechanism.

तेज़ बुखार और गर्दन में अकड़न मस्तिष्कावरण की सूजन का संकेत देते हैं, जबकि प्रकाश से डर इस सूजन के कारण तंत्रिकीय संवेदनशीलता बढ़ने से होता है। इन लक्षणों का संयोजन विकल्प (B) मेनिन्जाइटिस से मेल खाता है।

(High fever and stiffness in the neck indicate inflammation of the meninges, while fear of light occurs due to increased neurological sensitivity caused by this inflammation. The combination of these symptoms matches option (B), meningitis.)

English translations are provided solely for reader convenience and are not part of the inputs or outputs, except for the option labels (A–D).

Findings 4. Accuracy-only objectives fail to elicit Hindi reasoning form.

Implication. As both reasoning A and reasoning B are equally correct, accuracy-based objectives do not provide a gradient signal to discourage answer-driven or English rationale. During continuous optimization, models may therefore converge to objective-equivalent but linguistically mis-

aligned reasoning patterns. Collating findings 2, 3, and 4 together, Hindi medical reasoning requires both additional data and training objectives that explicitly shape reasoning behavior. This motivates the method introduced in the following section.

3 Method

3.1 Hindi Reasoning Objective

We define **Hindi reasoning** as a behavioral objective: under Hindi prompts, the model should organize medical rationales primarily in Hindi, independent of the final answer correctness. Rather than claiming to directly measure reasoning quality, we introduce a lightweight language-form scaffolding that discourages non-target language outputs by measuring the proportion of Hindi tokens in the generated sequence:

$$R^{\text{lan}} = \frac{1}{|y|} \sum_{t=1}^{|y|} \mathbb{I}[y_t \text{ is a Hindi token}], \quad (1)$$

where y denotes the generated sequence. This signal serves only as auxiliary guidance: maximizing Hindi token ratio alone does not ensure correct or well-structured reasoning. We therefore further assess Hindi reasoning through downstream accuracy and controlled human evaluations.

3.2 Decaying Scaffolding Reward Reinforcement Learning (DSR-RL)

Overview. We propose Decaying Scaffolding Reward Reinforcement Learning (DSR-RL), which optimizes reasoning tasks under two complementary rewards: a task-optimal accuracy reward and an auxiliary language-form scaffolding reward that encourages Hindi reasoning. A time-dependent schedule controls the relative influence of the two rewards, enabling early behavioral guidance without constraining the final task-optimal policy.

Notations. We treat the language model as an auto-regressive policy and optimize it accordingly. Let \mathcal{X} denote the space of prompts and \mathcal{Y} the space of output sequences. The policy $\pi_\theta(y | x)$ is parameterized by θ , where $x \in \mathcal{X}$ and $y = (y_1, \dots, y_T) \in \mathcal{Y}$. For each prompt x , we sample a group of K candidate responses $y_i \sim \pi_\theta(\cdot | x)$, indexed by $i \in \mathcal{G}(x)$. The overall reward consists of a task-optimal accuracy reward R^{acc} and an auxiliary scaffolding language-form reward R^{lan} . Their relative influence is controlled by a time-dependent coefficient $\lambda(\tau)$.

Task-optimal reward: reasoning accuracy. Each sampled response y_i receives a discrete correctness reward

$$R_i^{\text{acc}} \in \{0, 0.1, 1\}, \quad (2)$$

where 0 denotes responses without reasoning, 0.1 denotes responses with reasoning but incorrect answers, and 1 denotes responses with both reasoning and correct answers, following prior work (Trung et al., 2024; Chen et al., 2025).

Auxiliary scaffolding reward: reasoning language form. To encourage Hindi reasoning, we define a language-form reward to the whole sequence as

$$R_i^{\text{lan}} = \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \mathbb{I}[y_{i,t} \text{ is a Hindi token}] \in (0, 1]. \quad (3)$$

Group-relative advantages. For each reward type $k \in \{\text{acc}, \text{lan}\}$, we compute a group-relative normalized advantage within $\mathcal{G}(x)$:

$$\tilde{R}_i^k = R_i^k - \text{mean}_{j \in \mathcal{G}(x)}(R_j^k). \quad (4)$$

$$\hat{R}_i^k = \frac{\tilde{R}_i^k}{\text{std}_{j \in \mathcal{G}(x)}(\tilde{R}_j^k) + \varepsilon_{\text{norm}}}. \quad (5)$$

The advantage is broadcast to the token level:

$$\hat{A}_{i,t}^k = \hat{R}_i^k, \quad \forall t, k \in \{\text{acc}, \text{lan}\}. \quad (6)$$

Optimization objectives. At each token position (i, t) , we compute the likelihood ratio

$$r_{i,t}(\theta) = \frac{\pi_{\theta}(y_{i,t} | x, y_{i,<t})}{\pi_{\theta_{\text{old}}}(y_{i,t} | x, y_{i,<t})}. \quad (7)$$

Using the clipped operator

$$\mathcal{C}(r, A) = \min(rA, \text{clip}(r, 1 - \varepsilon_{\text{clip}}, 1 + \varepsilon_{\text{clip}})A), \quad (8)$$

we define the objective for each reward type as

$$\mathcal{J}^{\text{acc}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|y_i|} \sum_t \mathcal{C}(r_{i,t}(\theta), \hat{A}_{i,t}^{\text{acc}}). \quad (9)$$

$$\mathcal{J}^{\text{lan}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|y_i|} \sum_t m_{i,t} \mathcal{C}(r_{i,t}(\theta), \hat{A}_{i,t}^{\text{lan}}), \quad (10)$$

where $m_{i,t} \in \{0, 1\}$ is a token-level mask that equals 0 for English medical terms and 1 otherwise, leveraging principled code-mixing (Li et al., 2025). To prevent excessive policy drift, we include a KL regularization term:

$$\text{KL} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|y_i|} \sum_t D_{\text{KL}}(\pi_{\theta}(\cdot | h_{i,t}) \| \pi_{\text{ref}}(\cdot | h_{i,t})). \quad (11)$$

Final objective. The overall reinforcement learning objective is defined as

$$\mathcal{J}(\tau) = (1 - \lambda(\tau)) \mathcal{J}^{\text{acc}} + \lambda(\tau) \mathcal{J}^{\text{lan}} - \beta \text{KL}, \quad (12)$$

with the corresponding loss

$$\mathcal{L}_{\text{DSR}}(\tau) = -\mathcal{J}(\tau). \quad (13)$$

Policy optimization is performed using Group Relative Policy Optimization (Shao et al., 2024). Here, $\lambda(\tau) \in [0, 1]$ follows a cosine decreasing function over training iterations, so that auxiliary scaffolding reward dominates early optimization, while task-optimal reward governs the final policy.

3.3 Implementations

We implement the proposed DSR-RL objective via a three-stage training procedure, which progressively prepares the model for reinforcement learning and stabilizes Hindi medical reasoning.

Stage I: Language Adaptation via SFT (LA).

The first stage focuses on adapting the model to generate medically grounded Hindi text. This stage stabilizes basic vocabulary and medical expressions, providing a language-aligned initialization for subsequent reasoning-oriented training. Each training example consists of a short factual response without explicit reasoning.

Stage II: Reasoning Cold-start via SFT (RC).

With language generation stabilized, the second stage initializes structured medical reasoning. This stage provides a cold-start that shapes the model’s reasoning behavior, serving as a stable initialization for subsequent reinforcement learning. Each training example consists of an explicit chain-of-thought followed by a final answer, concatenated as a single output.

Optimization objective. We optimize the model using standard supervised fine-tuning with a token-averaged negative log-likelihood for stage $k \in \{\text{LA}, \text{RC}\}$. We apply teacher forcing and compute losses only on model output tokens. Sequences longer than a fixed maximum length are truncated, and padding tokens are ignored via attention masks.

$$\mathcal{L}_k = \mathbb{E}_{(x,y) \sim \mathcal{D}_k} \left[-\frac{1}{|y|} \sum_{t=1}^{|y|} \log \pi_{\theta}(y_t | x, y_{<t}) \right] \quad (14)$$

Stage III: DSR-RL. We perform reinforcement learning according to the design described in Section 3.2, with the reference model initialized from the checkpoint obtained at the end of Stage II.

3.4 Data Engineering

Overview. We design the HiMed data engineering pipeline to construct reasoning data under two settings: culture-grounded Hindi reasoning for traditional Indian medicine, and terminology-faithful reasoning for Western medicine. Accordingly, HiMed comprises two parts: HiMed-Trad and HiMed-West. All our data construction processes comply with ethical standards.

3.4.1 HiMed-Trad Corpus and Bench

Scope and coverage. HiMed-Trad is constructed through a unified OCR-based pipeline designed for noisy, real-world Hindi medical scans. A defining feature of HiMed-Trad is its comprehensive coverage of all seven officially recognized Indian systems of medicine.

Data construction pipeline. Authoritative textbook materials are first digitized using DeepSeek-OCR (Wei et al., 2025). The raw OCR outputs are then automatically repaired and page-grounded calibration is applied to recover coherent, self-contained Hindi medical passages. Each passage is subsequently converted into instruction-style instances, with source-grounded medical rationales expressed natively in Hindi.

Evaluation split and leakage prevention. To ensure reliable evaluation, all instances derived from the same source passage are treated as an indivisible unit and assigned exclusively to either the training corpus or the benchmark. This design enforces a strict passage-level split with zero overlap.

Statistics and references. Dataset statistics are summarized in Table 3, with the full taxonomy and source references provided in Appendix H.

System	Corpus Size	Bench Size
Ayurveda	105,397	1000
Yoga	8,499	1000
Naturopathy	20,735	1000
Sowa-Rigpa	119	10
Homoeopathy	130,751	1000
Unani	18,536	1000
Siddha	2,620	1000
Total	286,657	6,010

Table 3: Coverage of HiMed-Trad.

3.4.2 HiMed-West Corpus and Bench

Task focus. HiMed-West targets Western medical reasoning under Hindi prompts, where accurate handling of standardized medical terminology is

critical. We explicitly address this challenge by employing a carefully curated terminology lexicon.

Translation pipeline. We adopt a lexicon-guided translation pipeline that preserves standardized English medical terms when necessary. A high-coverage English–Hindi medical lexicon is constructed, and selected English medical corpora and benchmarks are translated using NLLB-200-3.3B (Team et al., 2022). For translated corpora, medical terms are accompanied by their original English forms in parentheses. In addition, we derive a specialized evaluation subset, HiMed-West Exam, from Indian national examinations on Western medicine: NEET UG.

Mixed-language supervision. To support effective cross-lingual transfer, only a category-stratified subset of data is translated into Hindi. This setting allows the model to be trained on both languages throughout the training process.

Evaluation safety and data sources. All translated training corpora are strictly disjoint from translated benchmarks to prevent leakage. Data sources and statistics are summarized in Table 4. GPQA-med refers to the Organic Chemistry, Molecular Biology, and Genetics subsets.

Source	Released Size	Translated Size
MedMCQA	182,822	91,411
Huatuo-o1 Corpus	19,704	9,852
Medreason	32,682	15,596
MMLU-Pro-Biology	717	717
MMLU-Pro-Health	818	818
GPQA-med	249	249
HiMed-Trad Corpus	286,657	–
HiMed-Trad Bench	6,010	–
HiMed-West Corpus	116,859	–
HiMed-West Bench	2254	–
HiMed-West Exam	470	–

Table 4: Sample composition of HiMed.

3.4.3 Data Quality

We apply multi-level quality control to ensure the reliability of the constructed data. Domain experts manually inspect randomly sampled outputs from key pipeline stages, yielding acceptance rates of 98.0% for OCR passages, 99.5% for HiMed-Trad instances, and 97.0% for translated outputs. Each inspection is conducted according to a predefined checklist. Detailed book and exam lists are provided in Appendix C and D. Complete pipeline details are provided in Appendix E, G, and I.

Model	MMLU-Pro-Biology			MMLU-Pro-Health			GPQA-med			HiMed-West Exam			HiMed-Trad	Avg. En	Avg. Hi
	En	Hi	Δ	En	Hi	Δ	En	Hi	Δ	En	Hi	Δ	Bench		
<i>Baseline Models and Our Model</i>															
GPT-4o	61.9	55.6	6.3	57.9	47.0	10.9	41.8	<u>30.5</u>	11.3	78.1	78.5	-0.4	<u>62.0</u>	59.9	54.7
UltraMedical-8B	66.8	41.7	25.1	56.1	26.3	29.8	33.3	20.9	12.4	54.7	40.2	14.5	45.1	52.7	34.8
MedReason-8B	65.8	43.7	22.1	<u>58.7</u>	27.6	31.1	27.3	29.7	-2.4	53.2	41.3	<u>11.9</u>	47.2	51.3	37.9
HuatuogPT-o1-8B	71.0	49.1	21.9	59.3	27.6	31.7	37.8	32.9	4.9	59.1	42.8	16.3	56.1	<u>56.8</u>	41.7
Qwen2.5-7B-Instruct	73.1	41.8	31.3	54.4	20.3	34.1	28.9	19.7	9.2	<u>62.0</u>	41.5	20.5	53.3	54.6	35.3
LLaMA-3.1-8B-Instruct	66.7	43.5	23.2	52.6	19.7	32.9	31.3	22.5	8.8	52.6	34.3	18.3	51.0	50.8	34.2
HiMed-8B (ours)	<u>69.9</u>	<u>53.8</u>	<u>16.1</u>	57.8	<u>32.5</u>	<u>25.3</u>	<u>38.6</u>	36.9	<u>1.7</u>	58.7	<u>46.6</u>	12.1	76.0	56.3	<u>49.2</u>
<i>Training Stage Ablation</i>															
Language Adaptation only	68.1	47.0	21.1	54.4	23.6	30.8	33.7	27.7	6.0	53.4	38.3	15.1	70.0	52.4	41.3
Reasoning Cold-start only	68.3	46.2	22.1	53.1	24.7	28.4	34.9	26.9	8.0	52.8	39.6	13.2	70.3	52.3	41.5
DSR-RL only	66.7	42.0	24.7	50.5	19.8	30.7	28.9	19.3	9.6	48.4	35.0	13.4	51.3	48.6	33.5
LA + RC	69.5	49.7	19.8	56.4	27.6	28.8	35.3	30.9	4.4	56.2	42.6	13.6	74.9	54.4	45.1
LA + DSR-RL	69.3	46.7	22.6	51.7	19.3	32.4	36.1	33.7	2.4	51.6	34.1	17.5	69.3	52.2	40.6
RC + DSR-RL	69.0	41.4	27.6	51.2	29.1	22.1	36.1	34.5	1.6	57.9	42.3	15.6	74.1	53.6	44.3
<i>Reward Ablation (DSR-RL)</i>															
Accuracy reward only	70.4	51.3	19.1	58.6	30.2	28.4	38.6	32.1	6.5	59.0	44.7	14.3	75.1	56.7	46.7
Constant rewards ratio	67.5	49.0	18.5	53.8	25.4	28.4	36.1	32.5	3.6	54.0	39.6	14.4	72.0	52.9	43.7

Table 5: Unified evaluation across medical reasoning benchmarks. Accuracy (%) is reported. Within *Baseline Models and Our Model* segment, **bold** highlights the best scores, and underlines indicate the second-best.

4 Experiments

4.1 Experimental Setup

Model. HiMed-8B is initialized from LLaMA-3.1-8B-Instruct. Training is conducted on a single node equipped with 8 NVIDIA H200 GPUs. Across all stages, the model is trained for approximately 56 hours, including 2 hours for LA, 15 hours for RC, and 39 hours for DSR-RL. Details of training data usage and hyperparameter settings are provided in Appendix K.

Benchmarks. We evaluate models on multiple medical reasoning benchmarks. Western medicine is assessed using MMLU-Pro-Biology, MMLU-Pro-Health, GPQA-med, and HiMed-West Exam, while Indian systems of medicine are evaluated using HiMed-Trad Bench. All benchmarks are evaluated under the zero-shot setting.

Task format and evaluation metrics. All evaluations use a multiple-choice format. In the LA stage, we apply greedy decoding with a four-class classification head. For all baselines and the RC and DSR-RL stages, responses are decoded using a robust answer-extraction procedure with progressively relaxed matching to ensure fair and consistent evaluation across models with diverse generation styles. The procedure ranges from strict answer-pattern detection to option-letter and option-text matching, and finally similarity-based option selection. We report accuracy and the English–Hindi accuracy gap $\Delta = \text{En} - \text{Hi}$ as evaluation metrics. All reported results are averaged over three runs.

4.2 Main Results

Table 5 summarizes performance across Western and Indian medical reasoning benchmarks. HiMed-8B consistently outperforms all open-source baselines in Hindi settings and substantially narrows the English–Hindi accuracy gap identified in findings 1, relative to the LLaMA-3.1-8B-Instruct base model, indicating effective cross-lingual transfer of medical reasoning. On HiMed-Trad Bench, HiMed-8B achieves state-of-the-art results, outperforming both open-source baselines and GPT-4o. Overall, these results show that the proposed training framework not only improves cross-lingual medical reasoning accuracy, but also enhances alignment with Hindi medical expression and domain knowledge.

4.3 Ablation Study

We conduct ablation studies to assess the contribution of each training stage and reward component in DSR-RL. As shown in Table 5, removing any stage leads to a clear drop in Hindi performance. Models trained with only the task-optimal accuracy reward or with fixed reward ratios exhibit inferior Hindi performance compared to the full decaying-reward formulation, indicating that the two rewards play complementary roles during the weight-shifting process. To verify that these gains are not driven by superficial Hindi token substitution, we further conduct a human preference study focusing on native Hindi medical reasoning. Two medical experts perform comparisons between HiMed-8B and an otherwise identi-

cal RL variant trained without the language-form reward. Across both Western medicine and Indian systems of medicine prompts, experts consistently prefer HiMed-8B. Preferred responses exhibit symptom-to-mechanism reasoning articulated directly in Hindi and avoid answer-driven explanation patterns. These results indicate that the language-form scaffolding reward induces qualitatively different reasoning behaviors. Full evaluation details are provided in Appendix M.

4.4 Effectiveness of the Language-form Reward

We analyze the training dynamics of the two rewards used in DSR-RL. As shown in Figure 2, both rewards increase throughout training. In particular, the consistent rise of the R^{lan} indicates that the model progressively produces reasoning chains with higher proportions of Hindi tokens. These results confirm that the R^{lan} provides an effective and meaningful optimization signal.

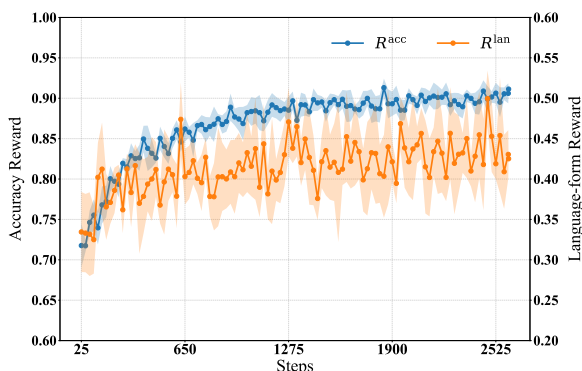


Figure 2: Training dynamics of rewards. The x-axis shows global steps, and the y-axis shows the reward signal. R^{acc} denotes the task-optimal accuracy reward, and R^{lan} denotes the auxiliary language-form reward.

4.5 Reliability of the Accuracy Reward Model

We fine-tune LLaMA-3.2-3B-Instruct model on a bilingual medical reasoning dataset as our accuracy verifier. At inference time, we apply softmax to the binary logits and use the probability of the True class as the confidence score. On the test set, the verifier achieves 0.969 Precision and 0.936 Recall. To assess reliability, we use GPT-5 to construct 300 challenging cases from two perspectives: (i) reversing key logical relations and (ii) slightly modifying numerical values in the reasoning. The verifier reaches 82.3% accuracy on this set. We further compare verifier judgments with annotations

from two medical experts on 300 randomly sampled model responses and observe a 96.7% exact agreement rate. Remaining disagreements mainly arise from inherently ambiguous cases. These results support using the verifier as a reliable reward signal in DSR-RL.

5 Related Work

After presenting our methods and empirical results, we situate our work within the broader literature. We briefly review prior work most relevant to our study and defer a comprehensive discussion to Appendix B. Existing resources for medical language modeling remain predominantly English-centric, with Hindi coverage largely limited to small-scale or translated subsets (Narayanan et al., 2024). Prior research has investigated cultural alignment (Li et al., 2024; Pham et al., 2025), reasoning enhancement (Creswell et al., 2023; Liang et al., 2024), reward modeling (Lightman et al., 2024; Ma et al., 2023), and large-scale English medical reasoning LLMs (Chen et al., 2025; Singhal et al., 2023; Zhang et al., 2024). However, these lines of work are typically studied in isolation and rarely intersect with Hindi medical reasoning. In particular, existing corpora and benchmarks fail to jointly support Hindi language supervision, Indian systems of medicine, and high-quality chain-of-thought annotations. As a result, the challenges of culturally grounded, linguistically faithful medical reasoning in Hindi remain largely unaddressed in prior work.

6 Conclusion

In summary, this work identifies Hindi medical reasoning as a critical yet underexplored requirement for medical large language models and introduces **HiMed**, the first large-scale Hindi medical corpus and benchmark suite with native rationales. Building on this foundation, we propose a three-stage adaptation framework centered on **Decaying Scaffolding Reward Reinforcement Learning**, which yields a concrete instantiation in the form of **HiMed-8B**. Extensive experiments show that HiMed-8B consistently improves Hindi medical reasoning performance across benchmarks, validating the effectiveness of the proposed data and training design. More broadly, our study offers practical insights into adapting general-purpose LLMs to linguistically and culturally specialized medical domains.

558 Limitations

559 This work has a few limitations.

560 First, while HiMed substantially expands cover- 608
561 age of Hindi medical reasoning, it does not exhaus- 609
562 tively represent all clinical scenarios or regional 610
563 variations within India. Certain subdomains, rare 611
564 conditions, and informal clinical expressions may 612
565 remain underrepresented. 613

566 Second, our evaluations focus on multiple- 614
567 choice medical reasoning benchmarks. Although 615
568 this setting enables controlled comparison across 616
569 models and languages, it does not fully cap- 617
570 ture open-ended clinical reasoning, longitudinal 618
571 decision-making, or real-world physician–patient 619
572 interactions. Extending evaluation to free-form 620
573 and interactive settings remains an important direc- 621
574 tion for future work. 622

575 Third, the proposed training framework is in- 623
576 stantiated on a single base model scale (LLaMA- 624
577 3.1-8B-Instruct). While our results demonstrate 625
578 consistent gains at this scale, further investigation 626
579 is needed to assess scalability to larger models and 627
580 generalization to other low-resource languages. 628

581 Finally, although we emphasize Hindi reasoning, 629
582 our approach does not replace the need for clinical 630
583 oversight. As with all medical language models, 631
584 HiMed-8B is intended for research and decision 632
585 support rather than autonomous clinical use. 633

586 Ethical Considerations

587 While our proposed model is a medical language 634
588 model with advanced reasoning capabilities, it may 635
589 still generate hallucinated, incomplete, or inaccur- 636
590 ate outputs. As such, it is not suitable for direct de- 637
591 ployment in real-world clinical settings. We explic- 638
592 itly restrict the use of HiMed-8B to research and 639
593 evaluation purposes, and prohibit its application 640
594 in clinical decision-making, diagnosis, treatment 641
595 planning, or other safety-critical scenarios where 642
596 errors could result in harm. Users bear ethical re- 643
597 sponsibility for respecting these limitations and for 644
598 avoiding inappropriate downstream use. 645

599 This work emphasizes Hindi medical reasoning 646
600 to support linguistic and cultural alignment in med- 647
601 ical AI research. Importantly, accuracy on bench- 648
602 marks related to Indian systems of medicine re- 649
603 flects consistency with officially examined curric- 650
604 ular, authoritative textbooks, and standardized as- 651
605 sessment materials. Such results should not be 652
606 interpreted as clinical efficacy, biomedical valida- 653
607 tion, or endorsement of any specific medical sys-

608 tem. Our benchmarks are designed to evaluate rea- 609
610 soning fidelity within established educational and 610
611 examination contexts rather than real-world thera- 611
612 peutic effectiveness. 612

613 We further acknowledge that medical knowl- 613
614 edge is culturally situated and historically contin- 614
615 gent. By constructing datasets that cover both 615
616 Western medicine and Indian systems of medicine, 616
617 our goal is not to adjudicate between medical 617
618 paradigms, but to enable transparent and faithful 618
619 modeling of how medical reasoning is expressed 619
620 within different linguistic and institutional frame- 620
621 works. The inclusion of traditional medical con- 621
622 tent does not imply clinical recommendation and 622
623 should not be used to guide health decisions with- 623
624 out qualified professional oversight. 624

625 Finally, all datasets in HiMed are derived from 625
626 publicly available or licensed sources, national ex- 626
627 amination materials, or scanned textbooks used 627
628 for educational purposes. We do not include pa- 628
629 tient records or personal health information, and 629
630 no identifiable individual data is involved. To fur- 630
631 ther ensure this, we explicitly ask GPT to ignore 631
632 these personal information during data generation. 632
633 In this study, all data, code, and models we use 633
634 and we release are under the CC-BY 4.0 license or 634
635 Apache license 2.0. We have verified that their us- 635
636 age complies with the original license agreements 636
637 and access conditions. Furthermore, participants 637
638 involved in the human evaluation phase were ex- 638
639 plicitly informed that their evaluation results would 639
640 be used solely for academic purposes. A total of 640
641 eight annotators from two countries participated 641
642 in this study, and each annotator was proficient in 642
643 both English and Hindi. 643

644 References

- 644 Asma Ben Abacha and Dina Demner-Fushman. 2019. 644
645 *A question-entailment approach to question answer-* 645
646 *ing*. *BMC Bioinformatics*, 20(1):511:1–511:23. 646
- 647 AI4Bharat. 2024. Indicqa: Question answering 647
648 dataset for indic languages. [https://huggingface.](https://huggingface.co/datasets/ai4bharat/IndicQA) 648
649 [co/datasets/ai4bharat/IndicQA](https://huggingface.co/datasets/ai4bharat/IndicQA). License: CC- 649
650 BY-4.0; includes Assamese, Bengali, Hindi, Kan- 650
651 nada, Marathi, Malayalam, Punjabi, Oriya, Tamil, 651
652 Telugu. 652
- 653 Rahul K. Arora, Jason Wei, Rebecca Soskin Hicks, Pre- 653
654 ston Bowman, Joaquin Quiñero Candela, Foivos 654
655 Tsimpourlas, Michael Sharman, Meghan Shah, An- 655
656 drea Vallone, Alex Beutel, Johannes Heidecke, and 656
657 Karan Singhal. 2025. *Healthbench: Evaluating large* 657

658	language models towards improved human health.	In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	714
659	<i>arXiv preprint arXiv:2505.08775</i> .		715
660	Kajol Bagga. 2025. Multilingual health-care text dataset (hi, en, pu). https://www.kaggle.com/datasets/kajolagga/multilingual-healthcare-text-dataset-hi-en-pu . Accessed: 2026-01-04.	CSTT. 2009. औषधि प्रतिकूल प्रतिक्रिया शब्दावली (<i>Drug Adverse Reaction Glossary</i>). वैज्ञानिक तथा तकनीकी शब्दावली आयोग (Commission for Scientific and Technical Terminology), नई दिल्ली (New Delhi). Accessed on 2025-08-30.	716
661			717
662			718
663			719
664			720
665	Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askill, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, and 12 others. 2022a. Training a helpful and harmless assistant with reinforcement learning from human feedback.		721
666			722
667			723
668			724
669			725
670			726
671			727
672			728
673			729
674	Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askill, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, and 32 others. 2022b. Constitutional AI: harmlessness from AI feedback.	CSTT. 2010. रोग-निदान एवं विकृति-विज्ञान शब्दसंग्रह (<i>Clinical Diagnostics and Pathology Glossary</i>). वैज्ञानिक तथा तकनीकी शब्दावली आयोग (Commission for Scientific and Technical Terminology), नई दिल्ली (New Delhi). Accessed on 2025-08-30.	730
675			731
676			732
677			733
678			734
679			735
680			736
681			737
682	Josh Barua, Seun Eisape, Kayo Yin, and Alane Suhr. 2025. Long chain-of-thought reasoning across languages. <i>arXiv preprint arXiv:2508.14828</i> .	CSTT. 2018a. कैंसर विज्ञान शब्दावली (<i>Oncology Glossary</i>). टाटा स्मारक अस्पताल (Tata Memorial Hospital), मुंबई (Mumbai). Accessed on 2025-08-30.	738
683			739
684			740
685	Shailaja Chandra and Kishor Patwardhan. 2018. Allopathic, ayush and informal medical practitioners in rural india—a prescription for change. <i>Journal of Ayurveda and integrative medicine</i> , 9(2):143–150.	CSTT. 2018b. शरीर रचना विज्ञान शब्द-संग्रह (<i>Glossary of Anatomy</i>). वैज्ञानिक तथा तकनीकी शब्दावली आयोग (Commission for Scientific and Technical Terminology), नई दिल्ली (New Delhi). Accessed on 2025-08-30.	741
686			742
687			743
688			744
689	Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. 2024. Humans or LLMs as the judge? a study on judgement bias.	Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca. <i>arXiv preprint arXiv:2304.08177</i> .	745
690			746
691			747
692	Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, and Benyou Wang. 2025. Towards medical complex reasoning with LLMs through medical verifiable problems. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 14552–14573, Vienna, Austria. Association for Computational Linguistics.	Alexander d’Elia, Mark Gabbay, Sarah Rodgers, Ciara Kierans, Elisa Jones, Irum Durrani, Adele Thomas, and Lucy Frith. 2022. Artificial intelligence and health inequities in primary care: a systematic scoping review and framework. <i>Family Medicine and Community Health</i> , 10(Suppl 1):e001670.	748
693			749
694			750
695			751
696			752
697			753
698			754
699	Sribala Vidyadhari Chinta, Zichong Wang, Xingyu Zhang, Thang Doan Viet, Ayesha Kashif, Monique Antoinette Smith, and Wenbin Zhang. 2024. Ai-driven healthcare: A survey on ensuring fairness and mitigating bias. <i>arXiv preprint arXiv:2407.19655</i> .	Dina Demner-Fushman and Jimmy Lin. 2007. Answering clinical questions with knowledge-based and statistical techniques. <i>Computational Linguistics</i> , 33(1):63–103.	755
700			756
701			757
702			758
703			759
704			760
705	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems.	Government of NCT of Delhi Directorate of AYUSH. Ashtāṅga ayurveda overview. https://ayush.delhi.gov.in/ayush/ayurveda . Accessed 2025-09-20.	761
706			762
707			763
708			764
709			765
710			766
711	Antonia Creswell, Murray Shanahan, and Irina Higgins. 2023. Selection-inference: Exploiting large language models for interpretable logical reasoning.	John W Ely, Jerome A Osheroff, Mark H Ebell, George R Bergus, Barcey T Levy, M Lee Chambless, and Eric R Evans. 1999. Analysis of questions asked by family doctors regarding patient care. <i>Bmj</i> , 319(7206):358–361.	767
712			768
713			769
		Pharmacopoeia Commission for Indian Medicine & Homoeopathy. Official website. https://pcimh.gov.in/ .	770

879	Richard Dufour. 2024. Biomistral: A collection of open-source pretrained large language models for medical domains . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 5848–5864. Association for Computational Linguistics.	935
880		936
881		937
882		938
883		
884		939
885		940
886	Hannah R. Lawrence, Renee A. Schneider, Susan B. Rubin, Maja J. Mataric, Daniel J. McDuff, and Megan Jones Bell. 2024. The opportunities and risks of large language models in mental health . <i>arXiv preprint arXiv:2403.14814</i> .	941
887		942
888		943
889		944
890		945
891	Jungseob Lee, Seongtae Hong, Hyeonseok Moon, and Heuseok Lim. 2025. Cross-lingual optimization for language transfer in large language models . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15100–15119, Vienna, Austria. Association for Computational Linguistics.	946
892		947
893		948
894		949
895		950
896		951
897		
898	Cheng Li, Mengzhuo Chen, Jindong Wang, Sunayana Sitaram, and Xing Xie. 2024. Culturellm: Incorporating cultural differences into large language models .	952
899		953
900		954
901		955
902	Yihao Li, Jiayi Xin, Miranda Muqing Miao, Qi Long, and Lyle Ungar. 2025. The impact of language mixing on bilingual LLM reasoning . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 32519–32536, Suzhou, China. Association for Computational Linguistics.	956
903		957
904		958
905		959
906		960
907		961
908		
909	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2024. Encouraging divergent thinking in large language models through multi-agent debate .	962
910		963
911		964
912		965
913	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let’s verify step by step .	966
914		967
915		968
916		969
917	Danni Liu and Jan Niehues. 2025. Middle-layer representation alignment for cross-lingual transfer in fine-tuned LLMs . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15979–15996, Vienna, Austria. Association for Computational Linguistics.	970
918		971
919		972
920		973
921		974
922		975
923		
924	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	976
925		977
926		978
927		
928		979
929		980
930		981
931	Renqian Luo, Luyang Sun, Yingce Xia, and et al. 2022. Biogpt: Generative pre-trained transformer for biomedical text generation and mining . <i>Briefings in Bioinformatics</i> .	982
932		983
933		984
934		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

988	Lowe. 2022. Training language models to follow instructions with human feedback .		1043
989			1044
990	Ankit Pal, Logesh Kumar Umapathi, and Malaikanan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering . In <i>Conference on Health, Inference, and Learning, CHIL 2022, 7-8 April 2022, Virtual Event</i> , volume 174 of <i>Proceedings of Machine Learning Research</i> , pages 248–260. PMLR.		1045
991			1046
992			1047
993			1048
994			1049
995			1050
996			1051
997	Jisu Park, Lucas Bandarkar, Artem Vazhentsev, Minjoon Choi, Swaroop Mishra, Chitta Baral, and Daniel Khashabi. 2025. Xtom: Exploring the multilingual theory of mind for large language models . <i>arXiv preprint arXiv:2506.02461</i> .		1052
998			1053
999			1054
1000			1055
1001			1056
1002	PCIM&H. a. Downloads: Unani pharmacopoeia of india (upi) / national formulary of unani medicine (nfum). https://www.portal.pcimh.gov.in/content/downloads .		1057
1003			1058
1004			1059
1005			1060
1006	PCIM&H. b. Homoeopathic pharmacopoeia of india (hpi) / homoeopathic pharmacopoeia laboratory (hpl) – publications portal. https://pcimh.gov.in/show_content.php?lang=1&level=1&lid=57&ls_id=59 .		1061
1007			1062
1008			1063
1009			1064
1010			1065
1011	PCIM&H. c. Siddha pharmacopoeial publications and siddha formulary of india. https://pcimh.gov.in/show_content.php?lang=1&level=1&lid=55&ls_id=57 .		1066
1012			1067
1013			1068
1014			1069
1015	PCIM&H. d. The unani pharmacopoeia of india – publication list. https://pcimh.gov.in/show_content.php?lang=1&level=1&lid=56&ls_id=58 .		1070
1016			1071
1017			1072
1018	Stephen R. Pfohl, Aleksandra Foryciarz, Harsha Nori, Michael Chen, and Tatsunori Hashimoto. 2024. A toolbox for surfacing health equity harms and biases in LLM-generated medical answers . <i>Nature Medicine</i> , 30:1731–1740.		1073
1019			1074
1020			1075
1021			1076
1022			1077
1023	Viet Thanh Pham, Zhuang Li, Lizhen Qu, and Gholamreza Haffari. 2025. CultureInstruct: Curating multicultural instructions at scale . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 9207–9228, Albuquerque, New Mexico. Association for Computational Linguistics.		1078
1024			1079
1025			1080
1026			1081
1027			1082
1028			1083
1029			1084
1030			1085
1031	Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Sean Shi, Michael Choi, Anish Agrawal, Arnav Chopra, Adam Khoja, Ryan Kim, Jason Hausenloy, Oliver Zhang, Mantas Mazeika, Daron Anderson, Tung Nguyen, Mobeem Mahmood, Fiona Feng, and 81 others. 2025. Humanity’s last exam . <i>arXiv preprint arXiv:2501.14249</i> .		1086
1032			1087
1033			1088
1034			1089
1035			1090
1036			1091
1037			1092
1038	Pengcheng Qiu, Chaoyi Wu, Xiaoman Zhang, Weixiong Lin, Haicheng Wang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2024. Towards building multilingual language model for medicine . <i>arXiv preprint arXiv:2402.13963</i> .		1093
1039			1094
1040			1095
1041			1096
1042			1097
	Hamed Rahimian and Sanjay Mehrotra. 2019. Distributionally robust optimization: A review . <i>arXiv preprint arXiv:1908.05659</i> .		1098
			1099
	Katrina T. Reaume, Maria Diaz, Sameer R. Patel, and Shoshana Cohen. 2025. Patient-physician language concordance and cardiovascular outcomes . <i>JAMA Network Open</i> , 8(1):e243803.		1100
			1101
	Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.		1102
			1103
	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. GPQA: A graduate-level google-proof q&a benchmark . <i>arXiv preprint arXiv:2311.12022</i> .		1104
			1105
	Alexey Romanov and Chaitanya Shivade. 2018. Lessons from natural language inference in the clinical domain . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1586–1596, Brussels, Belgium. Association for Computational Linguistics.		1106
			1107
	Shintaro Sakai, Jisun An, Migyeong Kang, and Hae-woon Kwak. 2025. Somatic in the east, psychological in the west?: Investigating clinically-grounded cross-cultural depression symptom expression in llms .		1108
			1109
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models . <i>arXiv preprint arXiv:2402.03300</i> .		1110
			1111
	Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Kumar Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Schärli, Aakanksha Chowdhery, Philip Andrew Mansfield, Blaise Agüera y Arcas, Dale R. Webster, and 11 others. 2022. Large language models encode clinical knowledge . <i>arXiv preprint arXiv:2212.13138</i> .		1112
			1113
	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaeckermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Andrew Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Agüera y Arcas, and 12 others. 2023. Towards expert-level medical question answering with large language models . <i>CoRR</i> , abs/2305.09617.		1114
			1115
	SNOMED International. 2024. SNOMED CT . Accessed: 2024-06-20.		1116
			1117

1099	Yu Sun, Xingyu Qian, Weiwen Xu, Hao Zhang, Chenghao Xiao, Long Li, Deli Zhao, Wenbing Huang, Tingyang Xu, Qifeng Bai, and Yu Rong. 2025. ReasonMed: A 370K multi-agent generated dataset for advancing medical reasoning . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 26457–26478, Suzhou, China. Association for Computational Linguistics.	1155
1100		1156
1101		1157
1102		1158
1103		1159
1104		1160
1105		
1106		
1107	Llama Team. 2024. The llama 3 herd of models .	
1108	NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, and 20 others. 2022. No language left behind: Scaling human-centered machine translation .	1161
1109		1162
1110		1163
1111		1164
1112		1165
1113		1166
1114		1167
1115		
1116	The State Education Department and The University of the State of New York. 2018. High school level chemistry glossary (english–hindi): Translation of chemistry terms based on the coursework for chemistry grades 9 to 12. Updated October 2018.	
1117		
1118		
1119		
1120		
1121	Luong Quoc Trung, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. 2024. Reft: Reasoning with reinforced fine-tuning . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 7601–7614. Association for Computational Linguistics.	1168
1122		1169
1123		1170
1124		1171
1125		1172
1126		1173
1127		1174
1128		1175
1129	Sshubam Verma, Mohammed Safi Ur Rahman Khan, Vishwajeet Kumar, Rudra Murthy, and Jaydeep Sen. 2024. Milu: A multi-task indic language understanding benchmark . <i>arXiv preprint arXiv:2411.02538</i> .	1176
1130		1177
1131		1178
1132		1179
1133	Veniamin Veselovsky, Berke Argin, Benedikt Stroebel, Chris Wendler, Robert West, James Evans, Thomas L. Griffiths, and Arvind Narayanan. 2025. Localized cultural knowledge is conserved and controllable in large language models .	1180
1134		1181
1135		1182
1136		1183
1137		1184
1138	Guoxin Wang, Minyu Gao, Shuai Yang, Ya Zhang, Lizhi He, Liang Huang, Hanlin Xiao, Yexuan Zhang, Wanyue Li, Lu Chen, Jintao Fei, and Xin Li. 2025a. Citrus: Leveraging expert cognitive pathways in a medical language model for advanced medical decision support . <i>arXiv preprint arXiv:2502.18274</i> .	1185
1139		1179
1140		1180
1141		1181
1142		1182
1143		1183
1144	Mingyang Wang, Lukas Lange, Heike Adel, Yunpu Ma, Jannik Strötgen, and Hinrich Schuetze. 2025b. Language mixing in reasoning language models: Patterns, impact, and internal causes . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 2637–2665, Suzhou, China. Association for Computational Linguistics.	1184
1145		1185
1146		1186
1147		1187
1148		1188
1149		1189
1150		1190
1151	Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers . In <i>NeurIPS</i> .	1191
1152		1192
1153		1193
1154		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209

1210	Juncheng Wu, Wenlong Deng, Xingxuan Li, Sheng Liu, Taomian Mi, Yifan Peng, Ziyang Xu, Yi Liu, Hyunjin Cho, Changin Choi, Yihan Cao, Hui Ren, Xiang Li, Xiaoxiao Li, and Yuyin Zhou. 2025a. Medreason: Eliciting factual medical reasoning steps in llms via knowledge graphs . <i>arXiv preprint arXiv:2504.00993</i> .	
1211		
1212		
1213		
1214		
1215		
1216		
1217	Linjuan Wu, Hao-Ran Wei, Baosong Yang, and Weiming Lu. 2025b. From english to second language mastery: Enhancing LLMs with cross-lingual continued instruction tuning . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 23006–23023.	
1218		
1219		
1220		
1221		
1222		
1223	Shaochen Xu, Yifan Zhou, Zhengliang Liu, Zihao Wu, Tianyang Zhong, Huaqin Zhao, Yiwei Li, Hanqi Jiang, Yi Pan, Junhao Chen, Jin Lu, Wei Zhang, Tuo Zhang, Lu Zhang, Dajiang Zhu, Xiang Li, Wei Liu, Quanzheng Li, Andrea Sikora, and 3 others. 2024. Towards next-generation medical agent: How o1 is reshaping decision-making in medical scenarios . <i>arXiv preprint arXiv:2411.14461</i> .	
1224		
1225		
1226		
1227		
1228		
1229		
1230		
1231	Wen Yang, Junhong Wu, Chen Wang, Chengqing Zong, and Jiajun Zhang. 2025. Implicit cross-lingual rewarding for efficient multilingual preference alignment . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 21125–21147, Vienna, Austria. Association for Computational Linguistics.	
1232		
1233		
1234		
1235		
1236		
1237		
1238	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	
1239		
1240		
1241		
1242		
1243		
1244		
1245		
1246	Kaiyan Zhang, Sihang Zeng, Ermo Hua, Ning Ding, Zhang-Ren Chen, Zhiyuan Ma, Haoxin Li, Ganqu Cui, Biqing Qi, Xuekai Zhu, Xingtai Lv, Jinfang Hu, Zhiyuan Liu, and Bowen Zhou. 2024. Ultramedical: Building specialized generalists in biomedicine . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	
1247		
1248		
1249		
1250		
1251		
1252		
1253		
1254		
1255	Huichi Zhou, Zehao Xu, Munan Zhao, Kaihong Li, Yiqiang Li, and Hongtao Wang. 2025. Moral reasoning across languages: The critical role of low-resource languages in llms . <i>arXiv preprint arXiv:2504.19759</i> .	
1256		
1257		
1258		
1259		
1260	Lianghui Zhu, Xinggang Wang, and Xinlong Wang. 2025. Judgelm: Fine-tuned large language models are scalable judges . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	
1261		
1262		
1263		
1264		
1265		
1266	Yuxin Zuo, Shang Qu, Yifei Li, Zhang-Ren Chen, Xuekai Zhu, Ermo Hua, Kaiyan Zhang, Ning Ding,	
1267		
	and Bowen Zhou. 2025. Medxpertqa: Benchmarking expert-level medical reasoning and understanding . In <i>Forty-second International Conference on Machine Learning, ICML 2025, Vancouver, BC, Canada, July 13-19, 2025</i> .	1268 1269 1270 1271 1272
	A Pilot Evaluation	1273
	For both MMLU-Pro and RET, we randomly sampled 150 test cases using a fixed random seed and translated them into Hindi to facilitate a controlled English–Hindi comparison under identical question content, following the translation pipeline described in Appendix I. For MMLU-Pro, we only included biology and health subsets. Specifically, RET comprises 100 questions from the 2023 AIAPGET exam, evenly drawn from four disciplines, and the remaining 50 questions from BNYS First Year, Part II examination papers. All RET questions were formatted as multiple-choice items with standardized options before translation and evaluation. GPT-4o was asked to output simply one option as the answer. HuatuoGPT-o1-8B and LLaMA3-8B-UltraMedical were deployed locally. For these two models, we adopted the option with the highest predicted probability as the model’s answer. This follows each model’s standard inference interface and avoids additional prompt engineering that could confound cross-model comparison. We reported the number of correct predictions out of the 150 test cases as the accuracy. The Δ value denotes the performance gap, calculated as the difference between Hindi and English accuracy. All evaluations are conducted in a zero-shot setting. The reported results are obtained through single run.	1274 1275 1276 1277 1278 1279 1280 1281 1282 1283 1284 1285 1286 1287 1288 1289 1290 1291 1292 1293 1294 1295 1296 1297 1298 1299 1300 1301
	B Related Work	1302
	Datasets. The dual nature of the Hindi medical ecosystem (Chandra and Patwardhan, 2018) imposes three key requirements on corpora for training Hindi medical reasoning LLMs: Hindi language coverage, inclusion of Indian systems of medicine, and support for chain-of-thought training. No existing resource satisfies all three requirements. MedMCQA (Pal et al., 2022), although derived from the NEET PG and AIIMS examinations, is publicly released almost entirely in English and focuses on Western medicine. UltraMedical (Zhang et al., 2024), BioASQ-QA (Krithara et al., 2023), PubMedQA (Jin et al., 2019), MedNLI (Romanov and Shivade, 2018), Reason-Med (Sun et al., 2025), and MedReason (Wu et al.,	1303 1304 1305 1306 1307 1308 1309 1310 1311 1312 1313 1314 1315 1316 1317

2025a) provide large-scale biomedical MCQ, QA, or NLI data, but are exclusively English. General multilingual corpora such as IndicInstruct (Gala et al., 2024), AI4Bharat-IndicNLP (Kunchukuttan et al., 2020), AI4Bharat-IndicQA (AI4Bharat, 2024), IndicLLM-Suite (Khan et al., 2024), XLingHealth (Jin et al., 2024; Lab, 2024), and MILU-cleaned (Verma et al., 2024; Murthyrudra, 2025) contain Hindi text, but either lack medical specialization or rely primarily on translated content. AyurGenixAI (kagglekirti123, 2025) covers Indian medicine but remains predominantly English, while the Multilingual Healthcare Text Dataset (Bagga, 2025) contains diverse medical content without providing a ready-to-use QA or MCQ format. Moreover, only MedMCQA, ReasonMed, and MedReason explicitly include CoT annotations. Overall, no existing corpus adequately covers both Western and Indian medicine in Hindi with high-quality CoT supervision, leaving a substantial gap for downstream model training.

Benchmarks. Benchmark imposes similar requirements of native Hindi coverage and inclusion of Indian medicine. Native Hindi evaluation is crucial to avoid translation bias and cultural information loss (Jin et al., 2024; Qiu et al., 2024; Wang et al., 2024a). However, few existing benchmarks meet these criteria. MedQA-USMLE (Jin et al., 2020), MMLU-Pro-Med (Hendrycks et al., 2021; Wang et al., 2024b), HLE-med (Phan et al., 2025), and HealthBench (Arora et al., 2025) are entirely English-based. This limitation extends to advanced biomedical reasoning benchmarks such as GPQA-med (Rein et al., 2023) and MedXpertQA (Zuo et al., 2025). BhashaBench-Ayur (Devane et al., 2025) consists primarily of OCR outputs from examinations, but merely focuses on Ayurveda, leaving a huge gap on other streams such as Yoga. Taken together, current benchmarks remain English-centric and fail to support comprehensive reasoning evaluation across both Western and Indian systems of medicine.

Cultural alignment. Recent work explores cultural grounding through region-specific data curation and instruction tuning (Li et al., 2024; Pham et al., 2025), or through principle-based approaches that promote normative adherence (Bai et al., 2022b). Culture-aware benchmarks have also been proposed (Wang et al., 2024c; Karinshak et al., 2024). However, most of these efforts focus

on reducing cultural bias or assessing cultural values, and remain largely English-centric. Culturally aligned methods for medical reasoning in Hindi, especially those spanning both Western and Indian systems of medicine, are still largely unexplored.

Reasoning enhancement. A broad range of techniques have been proposed to improve reasoning, including chain-of-thought prompting (Wei et al., 2022), self-consistency (Wang et al., 2023), verifier-guided decoding (Creswell et al., 2023), Tree-of-Thought (Yao et al., 2023), debate frameworks (Liang et al., 2024), and process supervision (Lightman et al., 2024). These methods improve robustness and accuracy on complex reasoning tasks, but are rarely adapted to multilingual or medical settings. None explicitly addresses the linguistic and cultural fidelity required for Hindi medical reasoning.

Reward modeling. Reward modeling is central to LLM alignment and can be broadly categorized into three paradigms. (1) Outcome reward models score final answers or pairwise preferences and form the basis of standard RLHF pipelines (Ouyang et al., 2022). (2) Process reward models (PRMs) provide step-level supervision for reasoning chains (Lightman et al., 2024), with recent work integrating PRM signals directly into reinforcement learning (Ma et al., 2023). (3) LLM-as-a-Judge (Zhu et al., 2025; Liu et al., 2023) replaces human annotation with strong LLM evaluators, and has been applied both for evaluation and as a reward source in RLHF or RLAIIF pipelines (Bai et al., 2022a; Cobbe et al., 2021; Chen et al., 2024). Multi-objective training has been explored in mathematical reasoning (Shao et al., 2024) and biomedical reasoning (Wang et al., 2025a), typically by combining reward signals or enforcing structured trade-offs. Distributionally robust optimization (DRO) (Rahimian and Mehrotra, 2019) further mitigates reward imbalance by optimizing worst-case performance. Despite these advances, reward modeling in medical domains remains largely accuracy-centric, with limited attention to language nativeness or cultural fidelity.

Medical Reasoning LLMs. Recent research has increasingly focused on enhancing the medical reasoning capabilities of large language models (LLMs). Representative systems include Med-PaLM2, BioGPT, PMC-LLaMA, BioMistral, HuatuoGPT-o1 and other domain-adapted LLMs

(Luo et al., 2022; Singhal et al., 2023; Wu et al., 2023; Chen et al., 2025; Labrak et al., 2024), which demonstrate strong performance on exam-style and reasoning-oriented benchmarks such as USMLE, MedQA, and MMLU (Jin et al., 2020; Hendrycks et al., 2021; Wang et al., 2024b). Existing approaches typically improve medical reasoning through three major methodological directions: (1) instruction tuning, which leverages curated or synthetic medical question–answer and diagnostic datasets to strengthen domain-specific reasoning (Jin et al., 2020); (2) reasoning modeling, such as Chain-of-Thought prompting or stepwise diagnostic reasoning templates, to encourage multi-step inference (Wei et al., 2022); and (3) preference tuning, including preference-based learning and reinforcement learning, to improve consistency and medical plausibility of generated reasoning paths (Chen et al., 2025).

Despite these advances, the vast majority of existing medical reasoning LLMs remain English-centric, with underlying medical knowledge and reasoning paradigms deeply rooted in Western clinical practice. Cross-lingual medical reasoning scenarios have received little systematic attention.

Cross-Lingual Transfer in Large Language Models. Cross-lingual transfer in LLMs typically combines language-adaptive pre-training, instruction tuning or alignment, and representation analysis. Language-Adaptive Pre-training (LAPT) continues pre-training on target-language corpora to build core linguistic competence, often complemented by vocabulary-oriented methods such as explicit vocabulary expansion (Cui et al., 2023) or adapter-based embedding learning with fixed model weights (Han et al., 2025). To transfer instruction-following and reasoning, prior work constructs multilingual instruction data and leverages cross-lingual in-context learning, translation-based strategies, and distillation (Wu et al., 2025b). Preference alignment has also been extended to multilingual settings by transferring preferences from English-aligned models via implicit rewards or modified DPO objectives on translation-augmented data, encouraging same-language responses (Yang et al., 2025; Lee et al., 2025). Finally, representation diagnostics suggest that middle layers are relatively language-agnostic, motivating explicit middle-layer alignment during task fine-tuning (Liu and Niehues, 2025).

C Books

This section lists representative books and official documents used in constructing the HiMedTrad corpus. The resources span multiple Indian systems of medicine, including Ayurveda, Yoga, Naturopathy, Unani, Siddha, Sowa-Rigpa, and Homoeopathy. The listed materials include classical treatises, official guidelines, and modern textbooks, reflecting both historical foundations and contemporary clinical practice.

C.1 Ayurveda

- **नाड़ी-दर्शन.** A text describing the Ayurvedic practice of pulse diagnosis.
- **Aahar Chikitsa.** A book on dietary therapy within the Ayurvedic system of medicine.
- **Anubhut Yog.** A work describing traditional Ayurvedic formulations and practices based on experiential knowledge.
- **Journal of Ayurveda Case Reports.** A peer-reviewed journal publishing case reports in Ayurveda, issued by the All India Institute of Ayurveda.
- **Essential Drugs List – Ayurveda.** An official document listing essential Ayurvedic medicines, including indications and usage contexts.
- **Charak Samhita.** A classical Ayurvedic text describing core principles and clinical practices.
- **चरक-संहिता: द्वितीय खंड.** The second volume of the Charak Samhita, a foundational Ayurvedic text.
- **Ayurvedic Standard Treatment Guidelines.** Standardized treatment guidelines prepared by the Ministry of AYUSH, Government of India.
- **रसरज महोदधि.** A comprehensive text covering a wide range of Ayurvedic medical practices.
- **संपूर्ण चिकित्सा (अष्टांग हृदयम पर आधारित).** A compendium of Ayurvedic practices based on the Ashtanga Hridayam tradition.
- **Scientific Basis for Ayurvedic Therapies.** A book discussing scientific perspectives on Ayurvedic therapeutic methods.

1514	• Sidhparikshapadati Part I. A text focusing on Ayurvedic diagnostic principles and treatment approaches.	1556
1515		1557
1516		
1517	• Sushruta Samhita. A classical Ayurvedic text with a focus on surgery and medical procedures.	1558
1518		1559
1519		
1520	• The Complete Book of Ayurvedic Home Remedies. A guide to household-level Ayurvedic remedies.	1560
1521		1561
1522		1562
1523	• अगद तंत्र (विष चिकित्सा) प्रथम भाग. A text on Agad Tantra, the Ayurvedic discipline concerned with toxicology.	1563
1524		1564
1525		1565
1526	• सुश्रुतसंहिता: शारीरस्थान हिंदी अनुवाद सहित. A Hindi translation of the Sharirasthana section of the Sushruta Samhita, covering anatomy and physiology.	1566
1527		1567
1528		1568
1529		1569
1530	• स्वास्थ्य अमृतावली. A guide to health maintenance based on Ayurvedic principles.	1570
1531		1571
1532	• प्रायोगिक वनऔषधि विज्ञान. A book on the practical application of forest-based herbal medicine in Ayurveda.	1572
1533		1573
1534		1574
1535		1575
	C.2 Yoga	
1536	• योगसन चिकित्सा. A book on the therapeutic applications of yoga postures.	1576
1537		1577
1538	• श्री योगसूत्र. A Hindi exposition of the Yoga Sutras, a foundational text of yoga philosophy.	1578
1539		1579
1540	• योगचिकित्सा: अनुपान के साथ. A work discussing yoga therapy in conjunction with medicinal practices.	1580
1541		1581
1542		1582
1543	• योग के रहस्य. A book exploring yogic practices and their philosophical foundations.	1583
1544		1584
1545		1585
	C.3 Naturopathy	
1546	• यौन स्वास्थ्योपयोगी. A text addressing sexual health through naturopathic practices.	1586
1547		1587
1548	• प्राकृतिक चिकित्सा क्यों?. A book discussing the principles and benefits of naturopathy.	1588
1549		1589
1550	• रोगों की नई चिकित्सा: New Science of Healing. A work presenting naturopathic approaches to disease treatment.	1590
1551		1591
1552		1592
1553	• Ankh-Ka: The Art of Naturopathy. A book on holistic health and natural healing methods.	1593
1554		1594
1555		1595
	• Upwash Chikitsa. A text on fasting-based therapeutic practices in naturopathy.	1556
	• स्वास्थ्य और जलचिकित्सा. A book on naturopathy and hydrotherapy.	1557
	C.4 Unani	
	• यूनानी चिकित्सा के आधारभूत सिद्धांत. A work outlining foundational principles of Unani medicine.	1558
		1559
	• Unanani Chikitsasar. A book describing Unani medical theory and practice.	1560
		1561
	• संयुक्त प्रान्त आयुर्वेदिक तथा यूनानी तिब्बी चिकित्सा पद्धति अधिनियम, 1939. A legislative document concerning Ayurvedic and Unani medical systems.	1562
		1563
	C.5 Siddha	
	• सिद्ध वनौषधि. A book on medicinal plants used in the Siddha system.	1564
		1565
	• सिद्ध प्रयोग – Part I. A Hindi commentary on Siddha medical practices.	1566
		1567
	C.6 Sowa-Rigpa	
	• नेपाली बौद्ध परम्परामा सोवारिग्पा पद्धतिको अवस्था. A study of the Sowa-Rigpa tradition within Nepali Buddhist contexts.	1568
		1569
	C.7 Homoeopathy	
	• होम्योपैथिक चिकित्सा. An introductory text on homeopathic medicine.	1570
		1571
	• होम्योपैथिक मटेरिया मेडिका (रिपर्टरी सहित). A comprehensive reference on homeopathic remedies with repertory.	1572
		1573
	• होमियोपैथिक पारिवारिक चिकित्सा (Paribarik Chikitsa). A guide to family-oriented homeopathic treatments.	1574
		1575
	• होम्योपैथिक औषधियों का सचित्र विवरण. A pictorial reference of homeopathic medicines.	1576
		1577
	• होम्योपैथी चिकित्सा. A general overview of homeopathic therapeutic practices.	1578
		1579
	• होम्योपैथिक दर्शन. A book discussing philosophical foundations of homeopathy.	1580
		1581
	• सरल होम्योपैथिक चिकित्सा. A practical guide to homeopathic treatment.	1582
		1583

- पशु रोगों में होम्योपैथी: एक वैकल्पिक चिकित्सा. A work on veterinary applications of homeopathy.
- सरल होम्योपैथी इलाज. A handbook of simplified homeopathic treatments.

D Exams

This section lists national- and university-level examination papers used as primary sources for constructing the HiMed benchmarks and corpora. All materials are publicly available examination papers and do not contain personal or private information.

National Eligibility cum Entrance Test (NEET-UG). This examination serves as the uniform entrance test for admission to undergraduate medical programs across India. Papers from the years 2018, 2020, 2021, 2023, and 2024 were used.

All India AYUSH Post Graduate Entrance Test (AIAPGET). AIAPGET is an entrance examination for postgraduate programs in Ayurveda, Unani, Siddha, and Homoeopathy. We incorporate papers from the Ayurveda, Unani, and Siddha disciplines.

Banaras Hindu University Research Entrance Test (RET) – Swasthavritta & Yoga. This Research Entrance Test (RET) paper is used for the Swasthavritta & Yoga discipline at Banaras Hindu University. The examination consists of two parts: research methodology and subject-specific knowledge. The latter includes concepts from Ayurvedic Swasthavritta, public health, Yoga, and naturopathy.

University of Patanjali – Bachelor of Naturopathy and Yogic Science (BNYS), First Year, Part II. This collection includes examination papers from the BNYS First Year, Part II examinations conducted in December 2023. The papers cover foundational subjects such as Human Anatomy and Human Physiology, each divided into two parts. Additional papers focus on the Philosophy of Nature Cure and the Fundamental Principles of Integrated Systems of Medicine, incorporating concepts from Ayurveda, Homoeopathy, and Unani.

E OCR Process

This section describes the OCR pipeline used to construct **HiMed-Trad Corpus**, **HiMed-Trad**

Bench, and **HiMed-West Exam**. Segment merging and LLM-as-a-judge merging are applied only to the HiMed-Trad, as exam questions are already provided in well-separated question-level formats and therefore do not require cross-segment reconstruction.

E.1 DeepSeek-OCR and Basic Cleaning

We first convert each raw PDF into high-resolution page images to support subsequent cross-modal verification. We then apply DeepSeek-OCR to extract page-wise text. The OCR outputs are exported as raw, unprocessed MultiMarkdown (MMD) files and retained without manual correction, in order to preserve original structure and OCR artifacts for systematic downstream handling. These outputs often contain OCR inconsistencies, including broken words, spurious line breaks, residual layout markers, and page-level noise. A basic cleaning stage normalizes punctuation and whitespace and attaches document-level metadata. We also keep page indices and source identifiers alongside each intermediate segment to maintain traceability back to the scanned pages. The goal is to produce a noise-reduced yet structurally faithful representation of the scanned material.

E.2 Segment Merging

OCR systems frequently fragment semantically continuous content across adjacent lines or pages. To restore coherence, we compare adjacent fragments using embedding-based cosine similarity with all-MiniLM-L6-v2 (Hearst, 1997; Reimers and Gurevych, 2019; Wang et al., 2020). Fragments above a similarity threshold are merged, while headings and list markers are preserved as independent units. After merging, each segment is stored as a self-contained record with its merged text and the corresponding page span, enabling later alignment and verification. This step reconstructs medically meaningful segments while maintaining the original narrative structure.

E.3 LLM-as-a-Judge Merging

For HiMed-Trad Corpus, each merged segment is classified by GPT-4o along four dimensions: (i) *Hindi validity*, whether the text is primarily written in Devanagari Hindi, (ii) *medical relevance*, whether the content expresses meaningful medical knowledge, (iii) *referential ambiguity*, whether unresolved references block medical interpretation,

Branch	Raw OCR	Cleaned	Merged	4-way Filter	Calibrated	Final "No-Problem"
Ayurveda	75,164	46,040	32,420	14,424	9,085	9,054
Yoga	11,582	5,698	3,090	1,922	844	839
Naturopathy	17,683	7,200	3,799	2,218	1,527	1,451
Sowa-Rigpa	68	57	37	21	10	10
Homoeopathy	64,752	39,953	18,400	12,673	9,316	9,226
Unani	6,447	5,148	2,900	2,477	1,487	1,486
Siddha	3,125	2,122	1,385	671	465	453

Table 6: Number of Hindi medical segments for each Indian medical branch after each stage of the OCR pipeline.

and (iv) *structural role*, whether the segment functions as a title or heading.

Segments not primarily written in Hindi are discarded. Segments flagged as ambiguous are further processed via a backward merging strategy: we iteratively concatenate the ambiguous segment with the closest preceding segment until an unambiguous segment is concatenated. Merged segments satisfying the pattern Hindi=True, Medical=True, Ambiguity=False, Heading=False are accepted into the high-confidence pool. We persist both the 4-way decision and the resulting merged unit in the intermediate files, so downstream calibration and scoring operate only on the high-confidence pool.

Prompt: LLM-as-a-Judge Merging

You are a classification tool. Analyze the following text and output exactly 4 lines:

is_hindi: True/False
is_medical: True/False
has_ambiguity: True/False
is_title_or_heading: True/False

Definitions:

- is_hindi: True if text is primarily Hindi in Devanagari.
- is_medical: True if it contains medical knowledge (Western/Indian).
- has_ambiguity: True only if the text contains a strong, clearly identifiable, critical referential ambiguity that satisfies all conditions below.

Conditions for has_ambiguity=True:

1. The ambiguity concerns a key subject/object/symptom/treatment/causal relation and blocks correct medical interpretation.
2. The ambiguous element cannot be resolved from the text itself, even with generous natural-language inference.
3. The ambiguity spans the entire text (the text never provides enough information to clarify the referent).
4. It is not a normal omission/abbreviation/vagueness/stylistic shortening/common medical phrasing.

5. It is not a minor unclear phrase that does not affect the main meaning.

If the text is understandable overall (even with minor unclear parts), then has_ambiguity must be False.

Rules:

- Output **only** the 4 lines above.
- No explanations.
- Use Python-style booleans: True/False.

Text:

{text}

E.4 LLM Calibration

Accepted segments are aligned with their corresponding page images and re-evaluated by GPT-4o, which is instructed to correct only OCR-level imperfections. This stage fixes minor transcription errors while explicitly prohibiting semantic alterations, additions, or deletions. Cross-checking against page images ensures that the calibrated text remains maximally faithful to the original scanned document.

Prompt: LLM Calibration

You are an OCR post-correction assistant. You are given OCR text that may span one or more scanned pages. You are also given all corresponding page images.

Picture:

- The corresponding page images are attached as image files.

Rules:

- Only correct characters, punctuation, spacing, and diacritics.
- Do **not** add new content or delete meaningful content.
- Do **not** translate.
- Output **only** the corrected text.

OCR text:

{text}

Page image:

{image}

E.5 LLM-as-a-Judge Scoring

GPT-4o assigns each calibrated segment one of three quality levels: `NO_PROBLEM` (readable and semantically intact), `POSSIBLE_ISSUE` (minor irregularities with preserved meaning), or `DEFINITE_ISSUE` (severe corruption obstructing interpretation). Only `NO_PROBLEM` segments are retained for downstream annotation and instruction generation; the other two categories are archived for potential manual recovery. Notably, the LLM-as-a-judge is used strictly for structural and quality control and does not introduce new content or alter medical meaning.

Prompt: LLM-as-a-Judge Scoring

You are an OCR text quality evaluator. Your main goal is to determine whether the text contains serious, meaning-breaking OCR corruption. Default to `NO_PROBLEM` unless there is clear evidence of major corruption.

Classification rules:

- 1. `NO_PROBLEM`:**
 - The text is readable, understandable, and mostly coherent.
 - Minor typos, spacing issues, missing matras, imperfect Hindi, or small artifacts are not major OCR problems.
 - If a normal Hindi reader can infer the meaning easily, choose `NO_PROBLEM`.
- 2. `POSSIBLE_ISSUE`:**
 - One or two suspicious fragments may be OCR errors, but overall meaning is still understandable.
- 3. `DEFINITE_ISSUE`:**
 - The text is seriously corrupted (broken words, incomplete fragments, nonsense sequences) and the meaning cannot be recovered.

Output format:

- Line 1: exactly one of `NO_PROBLEM` / `POSSIBLE_ISSUE` / `DEFINITE_ISSUE`
- Line 2: a short English explanation

Do not output Markdown, JSON, or extra text.

Text to evaluate:

{text}

E.6 Manual Inspection

To validate the reliability of the automated pipeline, two domain experts independently evaluated 200 randomly sampled segments. They confirmed that 98.0% of retained segments were medically cor-

rect, coherent, and contextually complete, deemed suitable for instruction generation. These results indicate that the OCR pipeline provides robust and reliable textual foundations for building Hindi medical instruction data. The experts are hired from local translation company. The experts are paid 500 currency units per hour, which is notably higher than the local translators' hourly wage. In total, 2,000 currency units are spent at this inspection stage.

E.7 Annotator Instructions for Manual Inspection

This manual inspection is conducted as a targeted quality audit of the automated OCR and data processing pipeline, rather than as full human correction or annotation. The goal is to verify whether the retained text segments are suitable for downstream instruction generation in Hindi medical reasoning tasks.

Annotator background. All annotations are performed independently by two domain experts with formal training in medicine or biomedical sciences and demonstrated familiarity with Hindi medical terminology.

Sampling protocol. Each annotator is provided with the same randomly sampled set of text segments drawn from the retained outputs of the OCR pipeline. Samples are selected across different source documents and medical topics. No information about downstream model training or benchmark usage is disclosed to the annotators.

Evaluation criteria. For each text segment, annotators assess the following criteria:

- **Medical correctness.** The medical content should be factually correct according to standard textbooks or officially examined curricula. Minor stylistic issues are acceptable as long as they do not alter medical meaning.
- **Semantic completeness.** The segment should be self-contained and convey a complete medical concept or explanation without requiring missing context from adjacent passages.
- **Coherence and readability.** The text should be logically coherent and readable in Hindi, without severe OCR artifacts that hinder understanding.

- **OCR alignment fidelity.** The textual content should faithfully reflect the original source material. Minor character-level OCR errors are acceptable if they do not affect semantic interpretation.

Decision rule. Annotators assign a binary judgment to each segment:

- **Accept:** The segment satisfies all evaluation criteria and is suitable for downstream instruction generation.
- **Reject:** The segment contains medical errors, missing context, or OCR artifacts that compromise semantic fidelity.

Disagreement handling. Annotations are conducted independently. In cases of disagreement, the segment is conservatively marked as rejected. No post-hoc adjudication or correction is performed.

Scope and limitations. This inspection focuses on verifying the usability and reliability of OCR outputs as textual foundations. It does not assess clinical safety, therapeutic efficacy, or real-world medical applicability of the content.

F Instruction Pool

F.1 Method

The instruction pool comprises high-quality medical questions from several public datasets: HealthBench (Arora et al., 2025), HealthSearchQA (Singhal et al., 2022), MedMCQA (Pal et al., 2022), MedReason (Wu et al., 2025a), and ReasonMed (Sun et al., 2025). To integrate these heterogeneous sources and improve generality, we employ a transformation pipeline in which GPT-4o abstracts the core intent shared by similar questions into a single standardized instruction.

Central to our approach is a five-type question taxonomy (Ely et al., 1999) developed from an analysis of thousands of questions posed by practicing physicians. This taxonomy aligns with structured representations common in medical question answering, such as the PICO framework, and principles of semantic unification (Demner-Fushman and Lin, 2007; Abacha and Demner-Fushman, 2019). It thus reinforces the distinctions among diagnostic, therapeutic, etiological, prognostic, and conceptual information needs.

We define five categories as follows:

- **Diagnosis.** Questions concerned with identifying a disease or condition from observed findings (e.g., “Could this patient have [disease]?”; “Is test [procedure] indicated in [context]?”).
- **Treatment.** Questions about clinical management, therapies, and medications (e.g., “What is the drug of choice for [disease]?”; “How should I treat [disease]?”; “What is the dose of [medication]?”).
- **Etiology.** Questions about causes, origins, or risk factors underlying a symptom or disease (e.g., “What causes [symptom]?”; “What factor leads to [finding]?”).
- **Prognosis.** Questions concerning expected disease course, severity, or complications.
- **Medical Knowledge.** General conceptual/definitional questions foundational to medical understanding (e.g., “What is [disease]?”).

In addition, we generalize instructions using placeholders grounded in the top-level hierarchies of SNOMED CT (SNOMED International, 2024), the world’s most comprehensive clinical terminology system. Each placeholder corresponds to a high-level clinical concept, making the abstraction scheme clear and robust. Table 7 lists the placeholders and definitions.

Prompt: Forming Instruction

You are a professional medical AI data scientist, proficient in multilingual processing. Your task is to analyze medical records and accurately extract the core medical consultation intent. You need to generate standardized outputs in three formats based on the complexity of each case and perform semantic deduplication on instructions sharing the same core intent.

Processing Steps and Requirements

1. Information Analysis and Intent Extraction

- Carefully read all information points provided by the user, such as symptoms, history, test indicators, doctor’s advice, demographics, and location.
- Identify and distill the user’s core medical consultation intent.

2. Standardized Output Generation

- Based on question complexity, generate Instruction in one of the following formats:
 - **A. Conversation Summary Format** (multiple user utterances):
Use a narrative summary: “In the conversation, the user provided [Information 1], [Information 2], ..., ultimately asking [Core Question].”

Term	Definition
Symptom	A physical or mental sign that may indicate a disease.
Disease	A specific illness or disorder with distinct symptoms and causes.
Cause	The reason why a disease happens, or the source of a health problem.
Medication	A substance (usually a drug) used to treat or prevent illness.
Procedure	A medical intervention or operation, including diagnostic and therapeutic care.
Context	Circumstances or background information relevant to the patient situation.
Age	An attribute of a person.
Gender	An attribute of a person.
Department	A specific medical department or specialty.
Treatment	A therapeutic intervention intended to cure, manage, or alleviate a disease.
Doctor	A qualified healthcare professional providing diagnosis, treatment, or guidance.

Table 7: Placeholder definitions used for instruction generalization.

- **B. QA Standard Format** (one user question):
Use a direct professional question with placeholders (e.g., [symptom], [disease], [procedure]).
- **C. MCQ Format** (requires choice):
Ask in second person and list options (A/B/C/...), using placeholders (e.g., [option]).

3. Intent Classification

- Assign one primary label: Diagnosis, Treatment, Etiology, Prognosis, or Medical Knowledge.

4. Semantic Deduplication

- Identify instructions with identical core intent and retain only the clearest, most standard one for each group.

5. Output Format

- Return a Markdown table with columns: Instruction, Category, Form.

Final Output: Return the Markdown table only.

F.2 Questionnaire

We conduct a survey to understand how people formulate health-related questions when seeking medical advice from doctors or AI assistants. The questionnaire is designed to elicit diverse symptom descriptions, diagnostic inquiries, and treatment-related questions across multiple scenarios. The survey is administered in both English and Hindi to accommodate a diverse respondent population. Below we present the full questionnaire.

Questionnaire

We are a research team developing AI to better answer health-related questions. Please share a few questions you might ask a doctor or an AI assistant.

हम एक शोध टीम हैं जो स्वास्थ्य संबंधी प्रश्नों का बेहतर जवाब देने के लिए AI (कृत्रिम बुद्धिमत्ता) विकसित कर रहे हैं। कृपया कोई ऐसे प्रश्न साझा करें जो आप एक डॉक्टर या AI (सहायक से पूछ सकते हैं।

1. Imagine you wake up feeling unwell. You have

sudden blurred vision in your left eye. You decide to message a doctor or a reliable AI health assistant for advice.

How would you describe your symptoms to a doctor or AI to get help with diagnosis?

कल्पना कीजिए कि आप बीमार महसूस करते हुए उठते हैं। आपकी बाईं आंख में अचानक धुंधला दिखाई देने लगा है। आप सलाह के लिए डॉक्टर या एक विश्वसनीय AI (स्वास्थ्य सहायक) को मैसेज करने का फैसला करते हैं।

निदान में मदद पाने के लिए आप डॉक्टर या AI (को अपने लक्षण कैसे बताएंगे?)

2. After your recent physical examination at your company/school, you received your report, which showed that several indicators (e.g., white blood cells, blood pressure) were beyond the normal range. You are worried. If you want to ask the doctor/AI why your results are abnormal and what might be causing these changes, what would you say?

आपकी कंपनी/स्कूल में हाल की शारीरिक जांच के बाद, आपको अपनी रिपोर्ट मिली, जिसमें दिखाया गया कि आपके कई मानक (जैसे सफेद रक्त कोशिकाएँ, रक्तचाप) सामान्य सीमा से बाहर थे। आप चिंतित हैं। यदि आप डॉक्टर या AI (से पूछना चाहते हैं कि “मेरे परिणाम असामान्य क्यों हैं और इन बदलावों का कारण क्या हो सकता है?” तो आप क्या पूछेंगे?)

3. If you unfortunately contract chikungunya fever (a new virus), which causes persistent fever, and you are very worried that you may not live long or have serious long-term effects, what would you ask the doctor?

यदि आप दुर्भाग्य से चिकनगुनिया बुखार (एक नया वायरस) की चपेट में आ जाते हैं, जिससे लगातार बुखार आता है, और आप बहुत चिंतित हैं कि कहीं आपकी उम्र कम न हो जाए या गंभीर दीर्घकालिक प्रभाव न हों, तो आप डॉक्टर से क्या पूछेंगे?

4. If you suspect you have a necrotizing skin infection and cannot go to the hospital, what questions would you ask an AI assistant for help?

यदि आपको संदेह है कि आपको नेक्रोटाइज़िंग स्किन इन्फेक्शन (त्वचा का गलना/मरना) है और आप अस्पताल नहीं जा सकते, तो मदद के लिए आप AI (सहायक से कौन से प्रश्न पूछेंगे?)

5. Suppose you are traveling in the Amazon region and have been experiencing watery diarrhea and cramps for two weeks. You are eager to find a way

1862

1863

1864

1865

1866

1867

1868

1869

1870

1871

1872

1873

1874

to relieve the symptoms. How would you ask a doctor for treatment methods? (Write down questions.)

कल्पना कीजिए कि आप अमेज़न क्षेत्र में यात्रा कर रहे हैं और आपको दो सप्ताह से पानी जैसा दस्त और मरोड़ हो रहे हैं। आप लक्षणों से राहत पाने का तरीका ढूँढने के लिए बेताब हैं। आप डॉक्टर से उपचार के तरीके कैसे पूछेंगे? (प्रश्न लिखें)

We ultimately collect 103 valid responses, from which 37 new instructions are extracted and categorized into five types, following generic instruction frameworks derived from prior seminal work (Abacha and Demner-Fushman, 2019; Demner-Fushman and Lin, 2007). Statistics are reported in Section F.3.

F.3 Statistics and Manual Inspection

Category	Former Work	Questionnaire	Total
Diagnosis	24	7	31
Etiology	22	7	29
Prognosis	26	6	32
Treatment	26	6	32
Med. Know.	26	7	33
Total	124	33	157

Table 8: Instruction counts by category and source.

Table 8 provides the instruction counts by category and source. The instruction pool is further reviewed manually with four focuses:

- Deduplication.** Instructions are grouped by category and semantically similar items are consolidated. For duplicates, only the clearest and most comprehensive instruction is retained. After deduplication, the remaining instructions account for 82% of the original pool.
- Complexity and clarity assessment.** Two experts independently assess each instruction for clarity, which requires the intent to be unambiguous, and executability, defined as the model’s ability to reasonably produce a coherent and accurate response. Only instructions approved by all experts are retained, yielding a pass rate of 61%.

F.4 Annotator Instructions for Instruction Pool Review

This manual inspection is conducted to assess the quality and usability of the instruction pool prior to model training. The goal is to ensure that retained instructions are non-duplicative, clear, and executable for medical reasoning tasks, rather than

to perform content rewriting or optimization. Two experts are hired from professional medical institute. The experts are paid 500 currency units per hour, which is notably higher than the local doctor’s hourly wage. In total, 1,000 currency units are spent at this inspection stage.

Annotator background. The review is performed independently by two domain experts with backgrounds in medicine or biomedical sciences and proficiency in Hindi medical language. Annotators are not informed of downstream model performance or experimental outcomes.

Review procedure. Each instruction is reviewed within its assigned category. Annotators follow the criteria below and provide binary judgments as specified.

Deduplication criterion. Annotators identify semantically duplicate or near-duplicate instructions within the same category. Two instructions are considered duplicates if they ask the same medical question or require substantially overlapping reasoning, even if surface wording differs. For duplicate groups, annotators retain a single instruction that is the clearest, most complete, and least ambiguous. All other duplicates are removed.

Clarity and executability assessment. Each remaining instruction is evaluated independently by both annotators according to the following criteria:

- Clarity.** The instruction should have a well-defined intent, with no ambiguous references, underspecified conditions, or unclear task requirements.
- Executability.** The instruction should be reasonably answerable by a language model using medical reasoning, without requiring external context, private information, or real-world clinical decision-making.

Decision rule. Annotators assign a binary decision to each instruction:

- Accept:** The instruction satisfies both clarity and executability criteria.
- Reject:** The instruction is ambiguous, underspecified, redundant, or not reasonably answerable.

Only instructions accepted by both annotators are retained. In cases of disagreement, the instruction is conservatively rejected.

Scope and limitations. This review assesses instruction quality for benchmarking and training purposes only. It does not evaluate clinical correctness, therapeutic safety, or real-world medical applicability of the instructions.

G HiMed-Trad Generation

G.1 Method

We construct HiMed-Trad via a staged pipeline that converts generation-ready merged segments into a structured, instruction-ready corpus.

Stage I: subject and type labeling. For each passage, we assign one or more intent categories from the five types (Diagnosis, Treatment, Etiology, Prognosis, Medical Knowledge), and one or more instance formats from MCQ, QA, Dialogue. Labeling is performed using GPT-4o. For multi-subject or multi-type passages, we expand them into atomic entries by taking the Cartesian product of category \times format, so that each final entry has exactly one category and one format.

Stage II: instruction generation with controlled templates. From labeled entries, we randomly select an instruction template using a predefined template pool and optional few-shot exemplars used only as style references.

Given (*text*, *subject*, *type*, *template*, *few-shot*), an LLM is prompted to generate instruction instances.

We enforce structured outputs and apply automatic validators to reject malformed generations, such as missing MCQ options, schema violations, or empty rationales when required.

Prompt: Template-Guided Instruction Generation

Role

You are an expert question writer for traditional medicine.

Your task is to generate high-quality instruction-style data (QA, MCQ, or Dialogue) grounded in traditional medicine texts, using predefined templates and optional few-shot exemplars. Ignore all the possible names or personal identifiers.

Input

For each generation instance, you are given:

- text: source paragraph
- subject: medical subject
- type: one of ["MCQ", "QA", "Dialogue"]
- template: a compatible instruction template

- few-shot: optional few-shot exemplars

Templates are NOT selected by another model. Instead, for each (subject, type) pair, a compatible template is randomly sampled to reduce selection bias and increase diversity.

Overall Generation Procedure

1. Select the given template and follow its structure strictly.
2. Generate three difficulty levels: EASY, MEDIUM, and HARD.
3. Ensure all outputs are self-contained, grounded in the source text, and consistent with traditional medicine.
4. Produce structured outputs that satisfy the required schema.

Malformed generations (e.g., missing MCQ options, schema violations, or empty rationales when required) will be automatically rejected.

Few-Shot Style References

The following few-shot blocks are provided as style references only (do NOT copy sentences verbatim):

- MCQ_FEW_SHOT {...}
- QA_FEW_SHOT {...}
- DIALOGUE_FEW_SHOT {...}

Difficulty Specification

For each instruction type, you MUST generate all three difficulty levels in the following fixed order: EASY, MEDIUM, HARD.

- **EASY**: single fact, no scenario, no reasoning required.
- **MEDIUM**: two-step reasoning, combining about two points; brief scenario optional.
- **HARD**: three to five-step reasoning; a realistic scenario is REQUIRED and must be grounded in traditional medicine.

Type-Specific Instructions

(A) MCQ

- Generate THREE multiple-choice questions (EASY, MEDIUM, HARD).
- Each question MUST have exactly five options: A., B., C., D., E.
- There must be ONE correct option.
- Provide:
 - the correct option letter
 - a clear reasoning
- MCQ is **NOT** a dialogue:

- Do NOT use "User:" or "Assistant:"
- Do NOT create multi-turn conversations

- Each item must be a single, self-contained question.
- Avoid any mention of text, paragraph, or source.

(B) QA

- Generate THREE short-answer questions (EASY, MEDIUM, HARD).
- Each item MUST consist of:
 - one direct question
 - one concise correct answer
 - reasoning
- QA is NOT a dialogue:
 - NO "User:" or "Assistant:" labels
 - NO multi-turn format
- Avoid any mention of text, paragraph, or source.

(C) Dialogue

- Generate THREE multi-turn dialogues (EASY, MEDIUM, HARD).
- Each dialogue MUST:
 - include two or more turns
 - contain explicit User: and Assistant: labels
 - end with a question from the User
- Provide:
 - the dialogue (up to the final user question)
 - the correct answer
 - reasoning
- The scenario must clearly belong to traditional medicine.

Output Format

Return ONLY the following structure (no explanations, no markdown):

```
<EASY><q>...</q><a>...</a><cot>...</cot>
<MEDIUM> ...
<HARD> ...
```

- MCQ: <q> must include options A. B. C. D. E.
- QA: must be a single direct question (no dialogue)
- Dialogue: must include multi-turn User:/Assistant: format

If the task is absolutely impossible, return: <FAIL>.

Stage III: automatic quality rating and filtering.

In the third stage of data construction, we apply an automatic quality inspection step to assess and filter generated instruction instances.

Each generated sample, consisting of a source paragraph, question, answer, and chain-of-thought rationale, is evaluated by a dedicated data quality rater prompted as a medical expert in traditional medicine.

The rater assigns continuous scores in four dimensions: contextual grounding, medical correctness, reasoning clarity, and Hindi language quality.

Based on a small pilot audit by two domain experts, we set the acceptance threshold to an average score of 0.725 (mean over the four dimensions), and retain only instances above this threshold for downstream use.

Here we provide the prompt used in this stage.

Prompt: Data Quality Rater

You are an expert data quality rater for a medical Q-A-CoT dataset in traditional medicine. You will be given one JSON object with fields:

- text: source paragraph (in Hindi)
- subject: one of ["diagnosis", "etiology", "medical knowledge", "prognosis", "treatment"]
- type: one of ["MCQ", "QA", "Dialogue"]
- question: generated question
- answer: generated answer
- cot: chain-of-thought reasoning

You must evaluate the question/answer/cot on four dimensions, each scored from 0.00 to 1.00:

1. grounded_in_context :

- Score 1.00 if all information in question/answer/cot is directly derivable from the source text
- Score lower if there is any hallucinated information or external knowledge not present in the text
- Score 0.00 if the content is completely unrelated to the source text

2. medical_correctness :

- Score 1.00 if the medical information is correct according to the source text
- Score lower if there are minor inaccuracies or misinterpretations

1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009

2010

- Score 0.00 if the medical information is incorrect or contradicts the source text

3. reasoning_clarity:

- Score 1.00 if the reasoning steps are logical, clear, and well-structured
- Score lower if reasoning is somewhat unclear or has minor logical gaps
- Score 0.00 if reasoning is illogical, confusing, or missing

4. language_quality:

- Score 1.00 if the Hindi language is fluent, natural, and grammatically correct
- Score lower if there are minor grammatical errors or awkward phrasing
- Score 0.00 if the language is severely broken or incomprehensible

Return ONLY a JSON object with 4 float scores in [0.00, 1.00]:

```
{
  "grounded_in_context": 0.00,
  "medical_correctness": 0.00,
  "reasoning_clarity": 0.00,
  "language_quality": 0.00
}
```

Do NOT add any explanations or extra keys. Return only the JSON object.

G.2 Statistics

Table 9 summarizes the corpus scale after each major stage. We report two types of statistics here: document- and page-level statistics, which reflect OCR coverage, and entry-level statistics, which reflect downstream usability.

Item	Count
Source documents	43
Total scanned pages	12,450
OCR-good MMD files	43
Segmented passages	21,932
QA instances	93,859
MCQ instances	102,494
Dialogue instances	96,314

Table 9: Summary statistics of HiMed-Trad.

G.3 Manual Inspection

Manual inspection is conducted as a targeted quality audit rather than full human correction, with the goal of verifying that the corpus is usable while preserving scalability and reproducibility. 200 samples are randomly selected and inspected. Two experts are hired from professional medical insti-

tute. The experts are paid 500 currency units per hour, which is notably higher than the local doctor’s hourly wage. In total, 3,000 currency units are spent at this inspection stage. Manual inspection is conducted on samples that have already passed automatic filtering, with the goal of validating the reliability of the quality rater rather than re-filtering the data.

Sampling protocol. We randomly sample entries across subjects, types, and instruction formats (QA, MCQ, and dialogue). To avoid bias toward easy cases, samples are stratified by passage length and script composition, ensuring coverage of both short factual passages and longer, reasoning-intensive contexts, as well as pure Hindi and code-mixed inputs.

Inspection criteria. Each sampled instance is independently reviewed by domain-aware annotators with a medical background. The inspection focuses on three aspects: (i) *semantic fidelity*, i.e., whether the question, answer, and rationale are faithful to the source passage without introducing unsupported content; (ii) *medical soundness*, i.e., whether the medical statements are correct and internally consistent; and (iii) *linguistic naturalness*, i.e., whether the Hindi expressions are fluent, coherent, and appropriate for medical communication. Importantly, annotators do not rewrite or correct instances, but only judge whether they are acceptable for inclusion.

Acceptance and outcomes. An instance is marked as *passed* if all inspected aspects meet the acceptance standard. Across all inspected samples drawn from instances that had already passed the automatic quality threshold, 99.5% are deemed acceptable without manual edits, indicating strong agreement between automated filtering and human judgment. The few rejected cases are primarily attributed to subtle linguistic infelicities or minor mismatches between rationales and source passages, rather than factual medical errors.

G.4 Annotator Instructions for Post-filter Manual Inspection

This manual inspection serves as a validation audit of the automatic quality rating and filtering stage. Annotators are instructed to assess whether instances that passed automatic filtering are acceptable for inclusion as-is. No rewriting, correction, or re-ranking of instances is permitted.

2074	Annotator background. All inspections are	H HiMed-Trad Taxonomy	2116
2075	conducted independently by domain-aware annota-		
2076	tors with a medical background and proficiency in		
2077	Hindi medical language. Annotators are blinded to		
2078	the automatic quality scores and thresholds.		
2079	Inspection setup. For each sampled instance, an-		
2080	notators are provided with: (i) the source passage,		
2081	(ii) the generated question, (iii) the generated an-		
2082	swer, and (iv) the chain-of-thought rationale. Judg-		
2083	ments must be based solely on these materials.		
2084	Evaluation criteria. Annotators evaluate each		
2085	instance along the following three aspects, which		
2086	correspond to the dimensions used by the auto-		
2087	matic quality rater:		
2088	• Semantic fidelity. The question, answer,		
2089	and rationale must be faithful to the source		
2090	passage. No hallucinated facts, unsupported		
2091	claims, or content external to the passage		
2092	should be introduced.		
2093	• Medical soundness. All medical statements		
2094	must be correct, internally consistent, and		
2095	aligned with standard medical knowledge or		
2096	officially examined curricula.		
2097	• Linguistic naturalness. The Hindi language		
2098	should be fluent, coherent, and appropriate for		
2099	medical communication. Minor stylistic im-		
2100	perfections are acceptable if they do not affect		
2101	comprehension.		
2102	Decision rule. Annotators assign a binary judg-		
2103	ment to each instance:		
2104	• Pass: All evaluation criteria are satisfied, and		
2105	the instance is acceptable for inclusion with-		
2106	out modification.		
2107	• Fail: One or more criteria are violated.		
2108	Annotators must not edit the content or provide		
2109	corrective feedback. In cases of disagreement, the		
2110	instance is conservatively marked as failed.		
2111	Scope and limitations. This inspection evalu-		
2112	ates the reliability of automated quality filtering		
2113	for dataset construction. It does not assess clini-		
2114	cal safety, real-world applicability, or downstream		
2115	model behavior.		
		To provide a balanced and authoritative founda-	2117
		tion, we draw on both peer-reviewed scholar-	2118
		ship and official documents across the major In-	2119
		dian systems of medicine. For Ayurveda, we	2120
		reference the historical and conceptual overview	2121
		(Jaiswal and Williams, 2017) together with the	2122
		Ashtāṅga Ayurveda materials from the Directorate	2123
		of AYUSH (Directorate of AYUSH). In the case	2124
		of Yoga, we rely on evidence from clinical and	2125
		quality-of-life studies (Woodyard, 2011), as well	2126
		as the official <i>Common Yoga Protocol</i> periodically	2127
		issued by the Ministry of Ayush (of Ayush and	2128
		MDNIY, 2025). Training standards are contex-	2129
		tualized with the ongoing WHO benchmarks for	2130
		Yoga (Organization, 2022). For Naturopathy, we	2131
		include parliamentary records that list its princi-	2132
		pal modalities (Ministry of Ayush, 2023), together	2133
		with the Central Council for Research in Yoga	2134
		and Naturopathy, alongside WHO’s 2010 training	2135
		benchmarks (for Research in Yoga & Naturopathy	2136
		, CCRYN; World Health Organization, 2010a).	2137
		The Sowa-Rigpa tradition is represented by the	2138
		National Institute of Sowa-Rigpa, Ministry of	2139
		Ayush portals, and pharmacopoeial initiatives by	2140
		PCIM&H (of Sowa-Rigpa , NISR; of Ayush, a; for	2141
		Research in Ayurvedic Sciences , CCRAS; for In-	2142
		dian Medicine & Homoeopathy). For Homoeopa-	2143
		thy, we cite both regulatory frameworks such as	2144
		the <i>National Commission for Homoeopathy Act</i>	2145
		(of India, 2020) and institutional sources including	2146
		CCRH, NIH Kolkata, and the Homoeopathic Phar-	2147
		macopoeia of India (for Research in Homoeopa-	2148
		thy , CCRH; National Institute of Homoeopathy ,	2149
		NIH; PCIM&H, b). Similarly, Unani medicine	2150
		is documented through AYUSH portals, CCRUM,	2151
		NIUM, and the official pharmacopoeia and formu-	2152
		lary volumes (of Ayush, b; for Research in Unani	2153
		Medicine , CCRUM; National Institute of Unani	2154
		Medicine , NIUM; PCIM&H, d,a). Finally, Sid-	2155
		dha medicine is covered via the Central Council	2156
		for Research in Siddha, the National Institute of	2157
		Siddha, and the Siddha Pharmacopoeia and Formu-	2158
		lary of India (for Research in Siddha , CCRS; Na-	2159
		tional Institute of Siddha , NIS; PCIM&H, c). Our	2160
		taxonomy-level comprehensiveness serves two pur-	2161
		poses. It avoids bias toward a single tradition and	2162
		supports fine-grained error analysis. For instance,	2163
		whether an LLM struggles more with Ayurveda,	2164
		Unani, or Siddha.	2165

2166	I Translation Process		
2167	We here introduce the detailed process of trans-	the OCR text minimally normalized to preserve lay-	2212
2168	lation, which is critical for the construction of	out cues and OCR artifacts.	2213
2169	HiMed-West Corpus and HiMed-West Bench.		
2170	I.1 Lexicon Table		
2171	To guarantee terminological precision and do-	I.1.3 Corpus-based Augmentation	2214
2172	main adaptability, we constructed a high-fidelity	(MedReason)	2215
2173	English–Hindi medical lexicon. The construc-	To complement static glossaries and match the vo-	2216
2174	tion process consists of utilizing authoritative glos-	cabulary distribution of <i>HiMed-West</i> , we augment	2217
2175	saries and augmenting them with data-driven ex-	the lexicon by extracting medical entities directly	2218
2176	tractions.	from the <i>MedReason</i> corpus, targeting long-tail	2219
2177		clinical terms that are absent from authoritative re-	2220
2178	I.1.1 Authoritative Medical Glossaries	sources.	2221
2179	The foundational stratum of our lexicon was cu-	I.2 Manual Correction and Translation.	2222
2180	rated from three distinct streams, comprising a to-	Two professional Hindi–English bilingual annota-	2223
2181	tal of six authoritative volumes.	tors perform line-by-line correction of OCR out-	2224
2182		puts and translation of extracted outputs from	2225
2183	General Medical Dictionary. We include a	corpus-based augmentation. The correction fo-	2226
2184	broad-coverage English–Hindi medical dictionary,	cules on mixed-script errors and confusions	2227
2185	<i>Dictionary of Medicine</i> (Gupta and Kapoor, 1955).	among visually similar Devanagari characters, ap-	2228
2186	The dictionary is organized alphabetically by En-	plied directly to the text to recover faithful content	2229
2187	glish headwords and provides Hindi equivalents,	and reduce downstream noise.	2230
2188	covering core topics such as anatomy, physiolo-		
2189	gy, diseases, symptoms, examinations, and treat-	Final Lexicon Consolidation. We merge three	2231
2190	ments.	sources: (i) digitized authoritative glossaries and	2232
2191		dictionary (Section I.1) and (ii) manual-translated	2233
2192	CSTT Specialized Glossaries. To align with of-	corpus-derived terms. Conflicts are resolved with	2234
2193	ficial Hindi terminology standards, we incorpo-	priority to official CSTT terminology. The fi-	2235
2194	rate four domain glossaries released by the <i>Com-</i>	nal consolidated lexicon contains 45,902 English–	2236
2195	<i>mission for Scientific and Technical Terminology</i>	Hindi term pairs. Authoritative glossaries con-	2237
2196	(<i>CSTT</i>), Government of India (CSTT, 2009, 2018a,	tribute 62% of entries, while corpus augmentation	2238
2197	2010, 2018b). These volumes provide standard-	contributes the other 38% .	2239
2198	ized, authority-sanctioned Hindi medical terms		
2199	and serve as our highest-priority reference during	I.3 Annotator Instructions for Lexicon	2240
2200	lexicon consolidation.	Construction and Correction	2241
2201		This annotation task focuses on constructing and	2242
2202	Educational Science Glossary. To complement	validating a high-fidelity English–Hindi medical	2243
2203	clinical resources with foundational biochemical	lexicon. Annotators are instructed to perform de-	2244
2204	and pharmaceutical nomenclature, we addition-	terministic correction and translation of medical	2245
2205	ally include the <i>Chemistry Glossary (High School</i>	terms, rather than creative paraphrasing or free-	2246
2206	<i>Level)</i> published by the <i>New York State Education</i>	form rewriting. The annotators are experts hired	2247
2207	<i>Department (NYSED) (The State Education De-</i>	from local translation company. The experts are	2248
2208	<i>partment and The University of the State of New</i>	paid 400 currency units per hour, which is notably	2249
2209	<i>York, 2018).</i>	higher than the local translators’ hourly wage. In	2250
2210		total, 30,000 currency units are spent at this inspec-	2251
2211		tion stage.	2252
2212		Annotator background. All annotations are	2253
2213		performed by professional Hindi–English bilin-	2254
2214		gual annotators with prior experience in medical	2255
2215		or biomedical translation. Annotators are familiar	2256
2216		with Devanagari script conventions and common	2257
2217		OCR artifacts.	2258

2259	Input materials. For each entry, annotators are provided with one of the following: (i) OCR-extracted glossary entries from authoritative medical volumes, or (ii) English medical terms extracted from the MedReason corpus for lexicon augmentation.	2302
2260		2303
2261		2304
2262		2305
2263		2306
2264		2307
2265	Correction and translation guidelines. Annotators follow the rules below:	2308
2266		2309
2267	• OCR correction. Correct mixed-script errors, broken characters, and confusions among visually similar Devanagari glyphs. Do not modernize spelling or alter stylistic variants unless required to restore correct meaning.	2310
2268		2311
2269		2312
2270		2313
2271		2314
2272		2315
2273	• Terminological translation. Translate each English medical term into its most appropriate Hindi equivalent. When an official CSTT term exists, it must be used preferentially.	2316
2274		2317
2275		2318
2276		2319
2277	• One-to-one mapping. Each lexicon entry must correspond to a single English term and a single Hindi rendering. Do not introduce explanations, definitions, or multiple alternatives.	2320
2278		2321
2279		2322
2280		2323
2281		2324
2282	• Conservativeness. If a term cannot be translated confidently without external context, annotators must flag it for exclusion rather than speculate.	2325
2283		2326
2284		2327
2285		2328
2286	Conflict resolution. In cases of disagreement between sources, annotators defer to CSTT terminology. Ambiguous or low-confidence entries are excluded from the final lexicon.	2329
2287		2330
2288		2331
2289		2332
2290	Scope and limitations. This process aims to ensure terminological precision and consistency for translation and benchmarking. It does not evaluate clinical validity or recommend medical usage.	2333
2291		2334
2292		2335
2293		2336
2294	I.4 Method	2337
2295	We propose a high-throughput, lexicon-guided translation pipeline that targets medical-domain fidelity while strictly preserving the structural constraints of benchmark datasets. The system combines rule-based terminology injection with neural machine translation (NMT), and is engineered for efficient large-scale processing.	2338
2296		2339
2297		2340
2298		2341
2299		2342
2300		2343
2301		2344
		2345
	I.4.1 Lexicon-Guided Hybrid Architecture	2346
	Our translation engine adopts a three-stage hybrid design: <i>lexicon injection</i> , <i>neural translation</i> , and <i>post-editing</i> . To enforce terminological consistency, we perform dictionary-based injection prior to neural inference. Specifically, using the lexicon in Section I.1, we compile medical terms into regex patterns and deterministically replace matched English terms (with word-boundary constraints) with their standardized Hindi equivalents before translation. This yields a code-mixed intermediate that constrains the NMT model to preserve domain-specific terminology.	2347
	Neural translation is performed with NLLB-200-3.3B (Team et al., 2022). For scalable processing, we segment text with <i>BlingFire</i> and apply length-aware dynamic batching to reduce padding overhead. Finally, we optionally append the original English term in parentheses as <i>Hindi Term (English Term)</i> to facilitate bilingual inspection and medical evaluation.	2348
	I.4.2 Adaptive Dataset Processing	2349
	To preserve evaluation formats across heterogeneous benchmarks, we adopt a format-adaptive pipeline that routes structured MCQs and unstructured inputs through different processing paths.	2350
	Structured inputs (MCQs). For MCQs, a structure-aware parser decomposes each instance into the instruction, question stem, and labeled options using regular expressions. Only the textual spans are translated, while option labels and structural markers are kept intact, preventing label corruption or reordering.	2351
	Unstructured inputs. For free-form prompts, we translate the input directly with the batch engine in Section I.4 to preserve discourse flow without unnecessary decomposition.	2352
	Fault-tolerant execution. All jobs run under a unified fault-tolerant executor with checkpointing and schema validation, ensuring resumable processing and that the translated outputs strictly conform to the original JSON format for downstream evaluation.	2353
	I.4.3 Manual Audit: NLLB vs. GPT-4o Translation Outputs	2354
	Comparison and implications. Using the same audit rubric, the NLLB-based pipeline requires	2355

fewer manual fixes than the GPT-4o-based pipeline and exhibits no translation-induced hallucinations. In practice, most remaining NLLB issues are localized and can often be mitigated deterministically via lexicon injection. By contrast, GPT-4o exhibits more frequent format violations and hallucinations under the same audit rubric. Overall, NLLB provides more stable, format-preserving translations, making it preferable for large-scale benchmark construction. This audit is conducted by two experts hired from local translation company, with 500 currency units hourly wage, higher than average level. In total, 2000 currency units are spent at this inspection stages. In another two experts' manual inspection consisting of 200 samples, our pipeline's results achieve an acceptance rate of 97.0%, indicating the usability and reliability of our method and practice. In total, 4,000 currency units are spent at this inspection stages.

Metric	GPT-4o	Ours
Need Fix	23	2
Hallucination	19	0

Table 10: Manual audit results of 50 samples under a unified rubric.

I.4.4 Annotator Instructions for Final Translation Quality Audit

This manual audit evaluates the quality and reliability of translated benchmark instances produced by the lexicon-guided translation pipeline. Annotators are instructed to assess usability and fidelity of translations as-is, rather than to edit or improve them.

Annotator background. Audits are conducted by bilingual reviewers with medical knowledge and proficiency in both English and Hindi. Annotators are blinded to the translation method used for each instance.

Audit setup. For each sampled instance, annotators are provided with: (i) the original English input, (ii) the translated Hindi output, and (iii) the original task format (e.g., MCQ structure, option labels).

Evaluation criteria. Annotators assess each instance according to the following criteria:

- **Terminological fidelity.** Medical terms should be translated consistently with the constructed lexicon. No critical terminology

should be mistranslated, omitted, or hallucinated.

- **Semantic equivalence.** The Hindi translation should preserve the meaning of the original English input without introducing unsupported information or altering intent.
- **Format preservation.** Structural elements such as MCQ options, labels, and JSON schemas must remain intact. Any format violation is grounds for rejection.
- **Linguistic adequacy.** The Hindi text should be fluent and comprehensible. Minor stylistic imperfections are acceptable if meaning and structure are preserved.

Decision rule. Annotators assign a binary judgment:

- **Accept:** The translation satisfies all criteria and is usable for benchmarking.
- **Reject:** One or more criteria are violated.

Annotators do not modify translations. No adjudication or correction is performed. In cases of disagreement, the instance is conservatively rejected.

Scope and limitations. This audit assesses translation reliability for dataset construction and evaluation only. It does not assess downstream model performance or clinical applicability.

I.5 Software Packages

Libraries (versions). torch 2.9.0+cu128; transformers 4.57.3; pandas 2.3.3; openpyxl 3.1.5; sentencepiece 0.2.0; blingfire 0.1.8.

J Decontamination

To ensure that benchmark performance reflects genuine generalization rather than data memorization, we apply decontamination at multiple levels.

First, **source-level separation** is enforced by construction: all books and official documents used for corpus generation are strictly disjoint from examination papers used for benchmarking, and no source appears in both splits.

Second, at the **passage level**, all data derived from the same OCR passage (as identified by a shared source document and page index) are assigned exclusively to either the training corpus or the benchmark, preventing any partial overlap across splits.

Third, for **translated benchmarks**, both the original questions and their translated or paraphrased variants are explicitly excluded from training data construction, ensuring that benchmark items are not indirectly exposed during training.

K Training Details

K.1 Training Data Overview Across Stages

Hyperparameter selection. Unless otherwise noted, all hyperparameters are determined using the training datasets described in this section, and we do not tune hyperparameters on any evaluation datasets.

Data mixing. Across all stages, we shuffle the training instances globally, mixing all data sources, task types, and languages into a single randomized stream for optimization. Table 11 summarizes the training data composition for each stage.

Dataset	LA	RC	DSR-RL
HiMed-West	46.5K	52K	10K
HiMed-Trad	102K	102K	10K
MedMCQA	46.5K	39K	
Huatuo-o1 Corpus	–	9K	–
GSM8K	–	8K	–
MedReason	–	–	5K
DailyDialog	11K	–	–
Persona-Chat	17K	–	–

Table 11: Training data composition.

In-domain pool and stage-wise split. For our Hindi traditional medical data, we use HiMed-Trad throughout the three-stage pipeline via uniform random sampling: 101.6k instances are used in LA, 101.6k instances are used in RC, and 10k instances are reserved for DSR-RL. This stage-wise allocation keeps the in-domain distribution consistent across stages while matching the intended training schedule.

Language Adaptation (LA). For LA, we use 102k in-domain instances from HiMed-Trad. We additionally incorporate MedMCQA with 46.5k English instances and 46.5k Hindi instances, all used with chain-of-thought rationales removed. We further include DailyDialog with 5.5k English instances and 5.5k Hindi instances, as well as Persona-Chat with 8.9k English instances and 8.9k Hindi instances.

Reasoning Cold-start (RC). RC is trained in two phases. Phase 1 runs for one epoch on 183.4k instances, including 101.6k in-domain examples

from HiMed-Trad and 81.8k MedMCQA examples with restored chain-of-thought rationales, comprising 40.9k English and 40.9k Hindi instances. Phase 2 uses 101.6k HiMed-Trad instances, together with Huatuo-o1 Corpus comprising 9.9k English instances and 9.0k Hindi instances after filtering samples whose chain-of-thought plus answer exceed 1536 tokens, and MedMCQA with 79.1k chain-of-thought annotated instances, including 39.5k English and 39.6k Hindi examples. In addition to these medical datasets, we include 7.5k GSM8K instances, split into 3.7k English and 3.7k Hindi examples.

DSR-RL. For DSR-RL, we use 10k in-domain instances sampled from HiMed-Trad. We additionally include MedReason with 10k Hindi instances and 5k English instances.

K.2 R_1 Reward Model Training Details

Base Model and Training Data We adopt Llama-3.2-3B-Instruct as the base model and construct the supervised fine-tuning corpus from two primary data sources. We extract 10k English medical reasoning instances from the Huatuo-o1 (Chen et al., 2025) dataset using a fixed random seed, and additionally sample 10,000 further instances that are translated into Hindi through our customized translation pipeline. We further incorporate the MedReason (Wu et al., 2025a) dataset, sampling 7.5k English instances and translating another 7.5k instances into Hindi in the same manner. To build the reward supervision labels, we use the unfine-tuned Llama-3.1-8B-Instruct model to generate Chain-of-Thought outputs for each problem and employ GPT-5 to compare these outputs with the corresponding reference CoT, producing ground-truth annotations on whether the two reasoning traces are semantically equivalent. The resulting bilingual dataset is used to train our reward model for CoT semantic equivalence evaluation. For the annotation of 300 randomly sampled model responses, two experts are hired from professional medical institute. The experts are paid 500 currency units per hour, which is notably higher than the local doctor’s hourly wage. In total, 3,000 currency units are spent at this inspection stage. Each expert was independently presented with 300 model-generated responses, paired with the corresponding question and ground-truth answer. The task was to assess whether the model’s final answer was factually and medically correct given the

question context. Experts were instructed to focus on the correctness of the conclusion rather than the surface fluency, verbosity, or stylistic quality of the reasoning. Intermediate reasoning steps were considered only insofar as they affected the validity of the final answer. Annotations were binary, labeled as True (the answer is correct) or False (the answer is incorrect or unsupported). In cases where the response involved minor wording variations, paraphrases, or alternative but medically equivalent formulations, experts were instructed to mark the answer as True as long as the medical conclusion was sound. If the case was inherently ambiguous, under-specified, or dependent on unstated assumptions that prevented a definitive judgment, experts were allowed to make their best professional judgment based on standard medical practice. The two experts performed annotations independently without access to each other’s labels or to the verifier’s predictions.

K.3 Language Adaptation Training Details

We conduct language adaptation via supervised causal language modeling on Stage-1 prompt-response pairs. Each sample is formatted with a Llama-3 style chat template, and we compute the training loss only on assistant tokens by masking the prompt tokens with -100. We apply dynamic padding using the eos_token_id (without introducing a new pad token) and tail-truncate sequences to a maximum length of 4096 tokens. Training is implemented with the Accelerate framework and DeepSpeed, using bf16 mixed precision and gradient checkpointing. We optimize with AdamW and a cosine learning-rate schedule with linear warmup. The best checkpoint is selected based on an EMA-smoothed training loss. Table 12 summarizes the hyperparameters used for language adaptation.

K.4 Reasoning Cold-start Training Details

We initialize the model’s reasoning capability via a cold-start supervised fine-tuning stage on chain-of-thought annotated data. Each training sample is formatted using a Llama-3 style chat template, where the user message is the original prompt and the assistant response is structured as: “## Thinking” followed by the annotated rationale (Complex_CoT), and “## Response” followed by the final answer (ground_truth). We compute the training loss only on assistant tokens by masking all prompt tokens with -100.

LA-Hyperparameter	Settings
Max sequence length	4096
Micro batch size (per GPU)	32
Gradient accumulation steps	1
Global batch size	256
Optimizer	AdamW
Learning rate	5e-6
Weight decay	0.01
Warmup ratio	0.03
LR scheduler	cosine
Epochs	3
Padding strategy	EOS padding
EMA decay (for selection)	0.9

Table 12: Language adaptation hyperparameter settings.

We employ dynamic padding with eos_token_id (without introducing a new pad token) and tail-truncate sequences to a maximum length of 1536 tokens. Training is implemented with Accelerate and DeepSpeed under bf16 mixed precision and gradient checkpointing, using AdamW with a cosine learning-rate schedule and linear warmup. The best checkpoint is selected based on an EMA-smoothed training loss. Table 13 summarizes the hyperparameters used in this cold-start stage.

RC-Hyperparameter	Settings
Max sequence length	1536
Micro batch size (per GPU)	8
Gradient accumulation steps	2
Global batch size	128
Optimizer	AdamW
Learning rate	1e-6
Weight decay	0.01
Warmup ratio	0.03
LR scheduler	cosine
Epochs	10
Padding strategy	EOS padding
EMA decay (for selection)	0.9

Table 13: Reasoning cold-start training hyperparameter settings.

K.5 Reinforcement Learning Training Details

We adapt the GRPO implementation from the trl framework for reinforcement learning training. Our setup includes (i) a Task-optimal reward verifier that has been validated through reliability testing, and (ii) a token-level auxiliary scaffolding reward to provide suitable low-resource language form. For LLaMA-3.1-8B-Instruct, we inject LoRA adapters into the projection modules q_proj, k_proj, v_proj, o_proj, gate_proj, up_proj, and down_proj in every decoder block. The details of training hyperparameters is summa-

2598 rized in Table 14, which we found the best during
2599 the experiments.

RL-Hyperparameter	Settings
Batch Size	8
Learning Rate	5e-6
LoRA_r	16
LoRA_alpha	32
LoRA_dropout	0.05
Task-optimal reward threshold	0.5
Reward Start Ratio	0.1:0.9
Reward End Ratio	0.9:0.1
Annealing Strategy	cosine
GRPO_clip	0.2
KL_beta	0.001
GRPO_num_generations	8

Table 14: Reinforcement Learning Hyperparameter settings.

2600 K.6 Software Packages

2601 **Libraries (versions).** torch 2.9.0; transformers
2602 4.57.3; accelerate 1.12.0; deepspeed 0.18.2; tqdm
2603 4.67.1; jinja2 3.1.6; swanlab 0.7.2.

2604 L Data Structures

2605 We organize the training corpus and benchmarks
2606 using a simple and consistent data format. Each
2607 training instance is associated with a stable identifier,
2608 a single medical intent category, and one instruction
2609 format. Training instances further include a question,
2610 its corresponding answer, and a reasoning rationale,
2611 all instantiated from predefined templates and explicitly
2612 grounded in the same source passage. Benchmark items
2613 are stored in a lighter format, containing only an identifier,
2614 the question text, and the gold answer.
2615

2616 M Human Evaluation of Native Hindi 2617 Medical Reasoning

2618 This appendix describes the human evaluation protocol
2619 used to assess whether the language scaffolding reward
2620 in DSR-RL improves native Hindi medical reasoning.
2621

2622 **Models compared.** We compare two reinforcement
2623 learning variants trained from the same Stage-II
2624 checkpoint: (1) **HiMed-8B**, trained with the full
2625 decaying-reward scaffolding framework including the
2626 language-form reward, and (2) **RL w/o Scaffolding**,
2627 trained using the same reinforcement learning procedure
2628 and data but without the language-form reward. All
2629 other training configurations, including model architecture,
2630 data, opti-

mization settings, and training duration, are kept
2631 identical. 2632

2633 **Evaluation data.** We randomly sample 200
2634 prompts from the evaluation pool, covering both
2635 Western medicine and Indian systems of medicine.
2636 For each prompt, both models generate a full
2637 response including an explicit medical rationale.
2638 Model outputs are anonymized and randomly ordered
2639 to avoid positional or identity bias.

2640 **Annotators.** Two annotators with formal medical
2641 training and native or near-native proficiency in
2642 Hindi independently perform the evaluation. Annotators
2643 are blinded to model identity and training configuration.
2644 Annotators are hired from professional medical institute.
2645 The experts are paid 500 currency units per hour, which
2646 is notably higher than the local doctor’s hourly wage.
2647 In total, 4,000 currency units are spent at this inspection
2648 stage.

2649 **Evaluation criterion.** Annotators perform a blind
2650 pairwise comparison and are asked the following question
2651 for each response pair:

2652 *Which response exhibits more native
2653 Hindi medical reasoning?*

2654 Judgments are guided by the following rubric,
2655 focusing on the *organization of medical reasoning*
2656 rather than surface language quality:

- 2657 • **Hindi-based reasoning organization:** 2658
2659 Whether medical reasoning steps (e.g., interpretation
2660 of symptoms, causal relations, or diagnostic logic) are
2661 explicitly structured within Hindi sentence forms rather
2662 than expressed as minimal or answer-driven statements.
2663
- 2664 • **Symptom–mechanism linkage in Hindi:** 2665
2666 Whether the response articulates connections between
2667 symptoms, underlying mechanisms, and conclusions
2668 directly in Hindi, instead of omitting reasoning or
2669 relying on implicit English-style shortcuts.
- 2670 • **Avoidance of English-centric reasoning patterns:** 2671
2672 Whether the response avoids answer-first or template-
2673 like explanations that could plausibly be translated
2674 verbatim from English (e.g., “The correct answer is B
2675 because ...”).
- 2676 • **Reasoning coherence:** 2677
2678 Whether the medical reasoning forms a logically
2679 connected

2678
2679
2680
2681
2682
2683

2684
2685
2686
2687
2688
2689
2690
2691

2692
2693
2694
2695
2696
2697
2698

2699
2700
2701
2702
2703
2704
2705
2706
2707

2708
2709
2710
2711
2712
2713
2714
2715
2716

2717
2718
2719
2720
2721
2722
2723
2724
2725

and internally consistent progression in Hindi, where each step (e.g., symptom interpretation, causal inference, and diagnostic conclusion) follows coherently from the previous one rather than appearing as isolated or fragmented statements.

- **Reasoning faithfulness:** Whether the reasoning process is faithfully grounded in the information explicitly stated in the prompt, with conclusions supported by articulated symptom–mechanism relations in Hindi, rather than introducing unstated assumptions, external knowledge, or implicit reasoning shortcuts.

Annotators are instructed *not* to judge factual correctness, verbosity, or stylistic fluency, and *not* to penalize principled code-mixing for medical terminology, numerals, or standardized symbols. A *Tie* option is explicitly provided and selected when neither response clearly exhibits more native Hindi medical reasoning according to the rubric.

Agreement and analysis. Annotators independently select one of three options for each response pair: *Response A*, *Response B*, or *Tie*, without access to any model identity information. Preferences are mapped back to their originating models only after annotation is completed. Preference rates are computed on non-tie decisions, while tie rates are reported separately to avoid forcing distinctions when no clear difference is perceived.

Across annotators, HiMed-8B is consistently preferred over the RL variant without the language scaffolding reward. Specifically, among non-tie judgments, HiMed-8B is preferred in 64.5% and 64.7% of cases by the two annotators, respectively. Using a conservative panel aggregation that requires annotator agreement and assigns *Tie* otherwise, HiMed-8B is preferred in 67.4% of non-tie cases, with a tie rate of 32.5%.

Inter-annotator agreement, measured by Cohen’s κ under the three-way decision setting, is 0.61, indicating substantial agreement. A two-sided binomial test on non-tie panel decisions shows that the preference for HiMed-8B is statistically significant. These results demonstrate that annealing the language scaffolding reward leads to more native Hindi medical reasoning behaviors as perceived by human experts.

N Use of AI Assistants

AI assistants were used in a limited and auxiliary manner during the preparation of this manuscript. Specifically, large language models were employed to support language polishing, grammatical refinement, and clarity improvements in portions of the text written in English. All technical content, including research ideas, methodological design, experimental setup, results, and interpretations, was conceived, implemented, and validated by the authors. AI assistants were not used to generate experimental data, conduct analyses, derive conclusions, or make scientific judgments. All code, model training, evaluations, and statistical analyses were performed by the authors.

O OCR Result Inspection UI

Here we provide the UI used in the OCR manual inspection. The interface is designed around a side-by-side verification workflow.

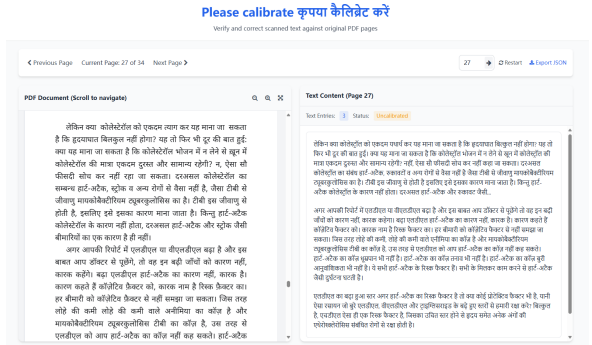


Figure 3: UI used in the OCR manual inspection

The original scanned PDF page is displayed on the left, while the corresponding OCR-extracted text for the same page is shown on the right. This layout enables annotators to directly compare source content and recognized text at the page level, quickly identify discrepancies, and perform targeted corrections. Page navigation, calibration status indicators, and export functionality are integrated to support efficient, traceable, and systematic manual OCR calibration across documents.

2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740

2741
2742
2743
2744

2745
2746
2747
2748
2749
2750
2751
2752
2753
2754