

MedicaLLM-Eval : Benchmarking AI Systems for Multilingual Reproductive and Maternal Health

Anonymous ACL submission

Abstract

Maternal health chatbots are already deployed to hundreds of millions of non-English speakers, but the benchmarks that evaluate them are almost entirely English. We built the Sakhi benchmark to close that gap. Sakhi is a parallel English, Hindi, and Marathi dataset of 380 maternal and reproductive health questions in two tracks, 149 doctor-edited pairs and 231 community-sourced pairs, validated by Accredited Social Health Activist (ASHA) workers and doctors in rural India. We evaluate 13 frontier LLMs spanning proprietary and open-weight systems with a two-layer framework: automated linguistic and semantic metrics, and a five-axis, fifteen-criterion clinical rubric (Accuracy, Completeness, Context Awareness, Communication, Terminology Accessibility) scored by an LLM judge with a three-judge calibration on a stratified 500-response subset. Four findings. No model exceeds Medical Quality Score 0.32 on the expert arm in English under the GPT-4o-mini judge, and the stricter GPT-5.1 judge drops every model a further 0.17 to 0.20 MQS on identical outputs. Terminology Accessibility is the weakest axis for every model (pass rates 0.123 to 0.206), below Context Awareness, so the dominant failure is lexical rather than broadly cultural. Apparent cross-lingual parity under a mid-tier judge is largely a judge-lenient artifact and does not survive a stricter judge. Open-weight 27B models (MedGemma 27B, Gemma 3 27B) match proprietary flagship MQS at roughly 20 times lower per-response cost, and medical fine-tuning alone does not lift rubric compliance in this cross-cultural setting. We release the dataset, reference answers, rubric, and judge-calibration records for the NeurIPS Datasets and Benchmarks track.

1 Introduction

Roughly 609 million people in India speak Hindi as a first or second language, and another 99 million speak Marathi (Express, 2025). Many of them

now get maternal-health information from conversational AI agents on messaging and social media platforms, including Meta AI and a growing set of government and NGO-run health bots (Chen et al., 2025; Longchar et al., 2025). The benchmarks that evaluate those agents are almost entirely in English (Singhal et al., 2025; Arora et al., 2025; Sellergren et al., 2025). Safety and accuracy assessments in Hindi or Marathi sit at the edge of what the field currently measures (Ghosh et al., 2025; Khullar et al., 2025), and the gap is not neutral. A response that correctly explains a pregnancy medication in English may drop the same warning in Marathi, or give the dose in a unit the reader cannot decode (Jin et al., 2024; Alonso et al., 2024; Rohera et al., 2025). Recent work documents that medical LLMs can produce unsafe recommendations and are vulnerable to generating harmful misinformation, particularly when tested outside their primary evaluation language (Moëll and Sand Aronsson, 2025; Han et al., 2024).

This paper introduces Sakhi, a community-validated benchmark for maternal and reproductive health LLMs in English, Hindi, and Marathi. We built it with Accredited Social Health Activist (ASHA) workers and doctors in rural India, anchored in a bot already deployed to that population, and we release it with two reference-answer tracks that a benchmark of this kind has not previously combined: a doctor-edited expert arm of 149 question-answer pairs and a community-sourced non-expert arm of 231 pairs, so that the MQS gap between the two is directly measurable within the same model. Nearby work touches pieces of this problem and stops short of the design we need. Khullar et al. (2025) evaluate triage classification accuracy on Indian-language maternal-health conversations, and Bhanusree et al. (2026) compare LLM responses to pregnancy queries in Telugu, but neither provides a domain-specific clinical rubric evaluation of open-ended generation across three

086 languages.

087 Our evaluation covers 13 instruction-tuned
088 LLMs drawn from a pre-declared eligibility frame
089 and stratified by access mode, training specializa-
090 tion, and capability tier. Each model produces three
091 independent responses per question per language
092 under identical prompts, for just under 45,000 total
093 generations, scored by GPT-4o-mini against a five-
094 axis, fifteen-criterion clinical rubric (Section 4) and
095 calibrated on a stratified 500-response subset with
096 two further judges (Claude Opus 4.6 and GPT-5.1).
097 The headline numbers in Section 5 support three
098 claims that bear directly on deployment. No model
099 exceeds MQS 0.32 on the doctor-edited arm in En-
100 glish, and the stricter GPT-5.1 judge drops every
101 model a further 0.17 to 0.20. Terminology Acces-
102 sibility is the weakest axis for every model, lower
103 than Context Awareness in every case, so the domi-
104 nant failure mode is lexical rather than broadly cul-
105 tural. Open-weight 27B models (MedGemma 27B,
106 Gemma 3 27B) reach the same MQS as proprietary
107 flagships at roughly one-twentieth the per-response
108 cost, and medical fine-tuning alone does not lift
109 rubric compliance in this cross-cultural setting.

110 Contributions.

- 111 1. Sakhi, a parallel English, Hindi, and Marathi
112 maternal and reproductive health benchmark
113 with two reference-answer tracks (doctor-edited
114 and community-sourced), validated through a
115 purpose-built Q&A review platform across three
116 stakeholder channels: practicing medical doc-
117 tors for clinical accuracy, Accredited Social
118 Health Activist (ASHA) frontline community
119 health workers for sociolinguistic grounding in
120 rural India, and healthcare nonprofit staff for
121 local cultural and infrastructural expertise.
- 122 2. A stratified 13-model evaluation with within-
123 family flagship/efficient pairs and a base-versus-
124 medical-finetune matched comparator, scored
125 under a five-axis clinical rubric and a three-
126 judge calibration that quantifies judge-model
127 bias.
- 128 3. Axis-level and cost-level findings that change
129 the recommended deployment choice for NGO-
130 scale maternal-health bots: prioritize lexical and
131 terminology alignment over medical fine-tuning,
132 and open-weight 27B models over proprietary
133 flagships.

2 Related Works 134

Automated Evaluation Metrics Early work on
135 automatic evaluation relied on surface overlap
136 methods such as BLEU (Papineni et al., 2002),
137 which measures n-gram precision against refer-
138 ence translations with a brevity penalty to dis-
139 courage overly short outputs, and ROUGE (Lin,
140 2004), which provides recall-oriented metrics in-
141 cluding n-gram overlap (ROUGE-N) and longest
142 common subsequence (ROUGE-L). These metrics
143 are reproducible and easy to compute, yet they
144 correlate poorly with human judgments (Callison-
145 Burch et al., 2006) as their reliance on exact lexi-
146 cal overlap means they may fail to fully capture
147 semantic adequacy or fluency in cases of valid
148 paraphrasing or synonym substitution, a limitation
149 acknowledged even in BLEU’s foundational pa-
150 per (Papineni et al., 2002). Such weaknesses are
151 particularly visible in low-resource and morpho-
152 logically rich languages (Sai et al., 2022; Datta
153 et al., 2022). To improve reproducibility, BLEU
154 score reporting guidelines and tools like Sacre-
155 BLEU were introduced (Post, 2018). The commu-
156 nity also adopted character-based metrics such as
157 chrF (Popović, 2015) and chrF++ (Popović, 2017),
158 which use character n-gram F-scores and are more
159 robust to orthographic variation, better suited for
160 morphologically rich languages, along with ME-
161 TEOR (Banerjee and Lavie, 2005), which improves
162 alignment with human preferences by combining
163 precision and recall, using flexible matching via
164 stemming and synonym databases (like WordNet),
165 and adding a fragmentation penalty to account for
166 word order. Even when combined in metric stacks
167 (BLEU/ROUGE/chrF), these methods have been
168 shown in multiple studies (Callison-Burch et al.,
169 2006; Sai et al., 2022; Kocmi et al., 2023) to miss
170 semantic and factual distortions, a limitation that
171 becomes critical in clinical evaluation, where tradi-
172 tional metrics are known to diverge from clinical
173 judgment and fail to assess semantic adequacy or
174 factuality (Tam et al., 2024; Croxford et al., 2025;
175 Abbasian et al., 2024).

To address the limitations of the aforemen-
177 tioned metrics, recent work adopts semantic em-
178 bedding-based approaches that capture contex-
179 tual meaning beyond lexical overlap. Sentence-
180 BERT (SBERT) (Reimers and Gurevych, 2019)
181 introduced a Siamese network structure to gener-
182 ate fixed-size sentence embeddings optimized
183 for semantic similarity, while SimCSE (Gao et al.,
184

2021) refined this idea through contrastive learning, producing robust representations. Similarly, contextual methods such as BERTScore (Hanna and Bojar, 2021) and MoverScore (Zhao et al., 2019) capture word-level and semantic alignment. These methods consistently outperform n-gram metrics for semantic comparisons (Sai et al., 2022; Datta et al., 2022). These approaches have been adapted for specific linguistic contexts; for instance, L3Cube-IndicSBERT (Deode et al., 2023) improves cross-lingual alignment for Hindi and Marathi, and benchmarks like IndicGenBench (Singh et al., 2024) support systematic evaluation of Indic text. Other studies assess translation quality for Indic languages using prompting and fine-tuning strategies (Nair et al., 2024), while XNLI benchmark tasks remain a standard for multilingual sentence evaluation (Conneau et al., 2018). Despite their advantages, embedding-based measures alone can overlook clinically important semantic drift (Tam et al., 2024; Chen et al., 2025) and are insufficient to ensure clinical safety, positioning them as one component of a multi-pronged evaluation anchored in explicit clinical validation.

LLM-as-an-Evaluator Approaches A third approach, which aims to scale up evaluation while mitigating the cost of human judges, involves using LLM-as-a-Judge frameworks to augment and validate other automated metrics, particularly for assessing safety and factual correctness. Benchmarks like OpenAI’s HealthBench (Arora et al., 2025) establish this foundation with physician-designed, multi-dimensional rubrics. This has spurred the development of LLM-based evaluators such as G-Eval (Liu et al., 2023) and LLM-Rubric (Hashemi et al., 2024), which calibrate judgments against human ratings to improve reliability. LLM judges, however, should not be treated as standalone assessors. They work best when cross-validated against human or rubric-based evaluations. Empirical studies provide evidence that LLM judges can exhibit biases, including self-preference where models favor their own outputs (Wataoka et al., 2025; Stureborg et al., 2024; Wang et al., 2024), and a general “agreeableness” bias, and even sycophancy (Malmqvist, 2024), which may lead to increased false-positive rates in certain evaluation contexts (Zhou et al., 2024). LLM judges also miss nuanced factual errors, which makes them less reliable for specialist correctness without oversight (Szymanski et al., 2025). In high-stakes healthcare

settings, these limitations translate into safety and equity risks, and some studies find demographic biases in evaluation outcomes themselves (Ye et al., 2024; Ji et al., 2025). Recent work now recommends anchoring LLM-judge outputs in deterministic, clinician-weighted rubrics and cross-validating them with other metrics and human review, a hybrid approach that matters most in sensitive contexts like maternal health (Chen et al., 2025; Szymanski et al., 2025).

Human Evaluation and Methodological Rigor Beyond automated metrics, evaluation research is moving past generic benchmarks toward domain-specific frameworks and tighter methodological standards. Recent examples include new proposals for conversational medical AI (Abbasian et al., 2024), structured clinical documentation assessment (Seo et al., 2024), and taxonomies for medical LLM performance (Lacerda et al., 2025). Alongside these proposals, researchers are adopting more standardized human evaluation with detailed rubrics that define clear criteria such as accuracy, coherence, and safety, rather than relying on subjective ratings (Bojic et al., 2023), and foundational guides on designing reliable user studies in NLP (Schuff et al., 2023) now inform the practice. This tighter rigor does not make human evaluation a clean gold standard. Clark et al. (Clark et al., 2021) demonstrate that human evaluation is often not fully reliable and is prone to sampling bias and low inter-annotator agreement. LLMs used as automated annotators are also vulnerable to adversarial prompt injection that can manipulate them into producing poisoned datasets and compromised findings (Baumann et al., 2025). The field now recommends using human studies to validate and contextualize automated approaches, not to replace them, with both requiring careful design to mitigate their respective limitations.

Domain-Specific Models and Evaluation Parallel to methodological advances in evaluation (Sai et al., 2022), model development has diverged into two main trajectories: general-purpose reasoning engines and domain-specialized biomedical systems. Generalist models such as the GPT-4 family, Gemini, and the multilingual Aya Expansive (Dang et al., 2024) demonstrate strong linguistic fluency and broad reasoning capabilities, yet their reliability in clinical contexts remains inconsistent, a gap that domain-specific models are explicitly designed to address (Sellergren et al., 2025; Chen

et al., 2023). By contrast, domain-adapted systems such as Meditron-70B (Chen et al., 2023) and MedGemma (Sellergren et al., 2025) are fine-tuned on biomedical corpora and clinical guidelines, often achieving superior performance on medical QA benchmarks when correctness and safety are prioritized. At the same time, comparative analyses caution that domain adaptation does not guarantee uniform improvements; robust gains depend on task-specific fine-tuning strategies (Anisuzzaman et al., 2025; Savage et al., 2025) and on careful, context-sensitive evaluation (Sai et al., 2022). This underscores the need for evaluation frameworks that are sensitive to linguistic and clinical variation, particularly in multilingual healthcare applications, a challenge highlighted by the non-English limitations of models like MedGemma and the multilingual focus of emerging generalist models (Dang et al., 2024; Sellergren et al., 2025). Multilingual medical benchmarks such as Multi-OphthaLingua (Restrepo et al., 2025) have revealed substantial cross-lingual bias in ophthalmological QA, while L2M3 (Gangavarapu, 2024) demonstrates the potential of translation-augmented LLMs for health equity in low-resource settings.

We construct a hybrid method built on methodical cross-validation with human evaluations. Our evaluation leverages automated metrics (e.g., chrF, BERTScore) (Popović, 2015; Zhao et al., 2019; Banerjee and Lavie, 2005) for reproducible, large-scale benchmarking, while acknowledging their inability to assess clinical safety by themselves (Tam et al., 2024; Croxford et al., 2025; Chen et al., 2025). Expert-grounded rubrics inspired by frameworks like HealthBench (Arora et al., 2025) are employed to encode medical correctness. LLM-judge outputs are employed but are anchored by human-expert validation (Chen et al., 2025; Szymanski et al., 2025), in line with best practices for safely deploying models in sensitive, multilingual contexts (Dang et al., 2024; Sellergren et al., 2025). While recent benchmarks such as CLINIC (Ghosh et al., 2025) evaluate multilingual trustworthiness across 15 languages and Script Gap (Khullar et al., 2025) assess triage accuracy on Indian-language maternal health conversations, neither provides domain-specific clinical rubrics for evaluating the quality of generated maternal health advice, which is the gap our framework closes.

3 Dataset

Multilingual datasets of maternal health queries from rural and semi-urban populations remain scarce (Singhal et al., 2025; Antoniak et al., 2024), and the few that exist typically mix synthetic generation with automated validation alone. Sakhi separates the two. We release the benchmark as 380 question-answer pairs across two tracks in parallel English, Hindi, and Marathi (1,140 query-language pairs total), drawn from a larger generation pool of 822 pairs through deduplication, medical-expert review, and, for the stricter track, full doctor rewriting.

Generation Pipeline. We compiled a knowledge corpus from trusted maternal health guidelines: the WHO Recommendations on Maternal Health, India’s National Antenatal Care Guideline (noa), the Auxiliary Nurse and Midwife Training Manual, and National Health Mission protocols. These documents were segmented into over 3,000 semantic chunks following retrieval-augmented generation (RAG) practices for specialized healthcare domains (Wu et al., 2025; Amugongo et al., 2025). A two-stage generator-validator produced the initial question-answer pool: Aya Expanse (Dang et al., 2024) drafted question-answer pairs from the chunks, and MedGemma (Sellergren et al., 2025) validated each pair for medical accuracy, safety, and contextual relevance.

A Purpose-Built Review System for Three Stakeholder Groups. AI-generated medical content needs human verification before it can serve as a benchmark reference, especially in high-stakes clinical domains (Amugongo et al., 2025; Fraile Navarro et al., 2025). Off-the-shelf annotation tools did not fit the workflow we needed, so we built a dedicated Q&A review platform that routes each generated pair through three distinct stakeholder groups, each contributing the kind of judgment they are actually qualified to make.

The first channel is **practicing medical doctors** who provide clinical validation. Doctors receive a pair, a review rubric (accuracy, safety, completeness of contraindications and warnings), and an editor to rewrite the reference answer if they find it insufficient. The pairs they rewrite in this channel form the **expert arm** of 149 doctor-edited question-answer pairs. The second channel is **Accredited Social Health Activist (ASHA) workers**, frontline community health workers in rural India who are

the daily point-of-care contact for most of our target users. ASHA workers provide the benchmark’s sociolinguistic grounding: they flag questions that do not sound like a patient would ask them, flag reference answers that use clinic-only vocabulary a rural user cannot parse, and flag omissions around locally-available health infrastructure (government schemes, anganwadi centres, district health posts). The third channel is **healthcare nonprofit staff** with deep local expertise in cultural norms, family dynamics, and the specific trust architecture through which health information actually reaches women in rural India. This group caught the subtle failures that clinical review alone misses, such as advice that is medically correct but culturally implausible for a daughter-in-law to act on without her mother-in-law’s approval, or for a first-time mother to raise without her husband present.

Following the gold standard for clinical NLP validation (Tam et al., 2024; Bojic et al., 2023), the three review channels feed a two-tier pipeline that produces the two released tracks. In the first tier, the medical-expert channel and the nonprofit-staff channel review the machine-generated English pairs for clinical accuracy and sociocultural appropriateness; pairs that pass both form the **non-expert arm** of 231 community-sourced pairs with doctor-reviewed reference answers. In the second tier, validated pairs are routed to practicing doctors and medical students who write their own answer versions; the 149 rewritten pairs form the **expert arm**. Both arms are then professionally translated into Hindi and Marathi by native speakers, and the ASHA worker channel reviews the translated pairs for patient-voice fidelity before release. Each review action is logged in the platform with reviewer role, timestamp, and edit diff, so that the provenance of every released reference answer is recoverable and auditable.

What the two tracks enable. The expert and non-expert arms are not redundant. They let us measure, within the same model, how much MQS drops when the reference answer gets stricter: the 0.10 to 0.13 MQS expert-minus-non-expert gap reported in Section 5.2. Benchmarks that report on a single reference-quality level cannot surface this gap. We release both arms as `sakhi_benchmark_expert.csv` and `sakhi_benchmark_non_expert.csv` along with the rubric, judge prompts, and judge-calibration records.

4 Methodology

Evaluation Framework and Problem Setting

Our work introduces a two-stage evaluation framework designed to address limitations of LLM benchmarks for medical applications. The framework emphasizes clinical validity, linguistic fidelity, and contextual appropriateness, focusing on maternal health scenarios across English, Hindi, and Marathi. Maternal health is a high-stakes domain where misinformation has direct consequences for patient safety. In multilingual settings, linguistic nuance and cultural appropriateness also shape whether a patient understands and follows the advice. Our framework evaluates model-generated responses along three complementary dimensions:

1. Semantic fidelity, which measures the preservation of clinical meaning;
2. Linguistic quality, assessing fluency and coherence;
3. Expert-based clinical validity, which evaluates the factual and ethical soundness.

The evaluation dataset (Section 5 and Section 3) comprises 380 patient-style English queries split across two reference-quality tracks, 149 doctor-edited (expert arm) and 231 community-sourced (non-expert arm), spanning ten maternal health themes. Each English query was professionally translated into Hindi and Marathi, yielding parallel evaluation sets of 380 questions per language (1,140 query-language pairs in total). The evaluation proceeds sequentially: automated benchmarking establishes reproducible linguistic and semantic baselines (Layer 1), followed by structured rubric review (Layer 2). The two layers are then synthesized into multi-dimensional indicators of model reliability.

Models Evaluated The evaluation covers 13 instruction-tuned LLMs drawn from a pre-declared eligibility frame and stratified along three axes: access mode (proprietary vs. open-weight), training specialization (general-purpose, multilingual-specialized, medical-specialized), and capability tier (flagship vs. efficient). The stratification, the full panel composition, and the matched base-versus-medical-finetune comparator (Gemma 3 27B alongside MedGemma 27B) are described in Section 5.1. Model selection prioritized public API or open-weight availability, documented

multilingual capabilities, and relevance to health-care dialogue. Models were accessed through OpenRouter, the local Claude Code CLI (Claude Opus 4.7), and HuggingFace Inference Endpoints (MedGemma variants), with identical prompts and per-family generation parameters documented in Appendix B.4. Each model generated three independent responses per (question, language) under no multi-turn context, yielding up to $13 \times 3 \times 3 \times 380 = 44,460$ total model outputs before coverage gaps are applied (see Section 5.1 and the Limitations).

Layer 1: Automated Benchmarking This layer enables scalable, language-agnostic evaluation of two core properties: semantic fidelity and linguistic quality. Automated scoring provides standardized quantitative baselines across models and languages prior to clinical review.

Semantic Fidelity Assessment Preserving clinical meaning is central to the safety of automated responses. To evaluate semantic fidelity, two complementary similarity metrics are employed. Embedding-based alignment measures overall conceptual correspondence using cosine similarity between dense sentence embeddings generated by OpenAI’s text-embedding-3-small model within a shared multilingual space. While additional embedding models (Cohere, Voyage, and alternative architectures) were evaluated during development, text-embedding-3-small was selected as the primary metric for its consistent performance across all three languages; comparative results from alternative embedding models are provided in Appendix A. BERTScore, computed with xlm-roberta-base for multilingual text and roberta-base for English, captures token-level semantic overlap and contextual entailment across Indic and Latin scripts. All configurations, model versions, and scripts are documented for reproducibility. To mitigate single-reference bias, sampled outputs were manually inspected for semantic entailment and corrected when necessary.

Linguistic Quality Assessment Linguistic quality is evaluated through complementary metrics capturing fluency, coherence, and readability. SacreBLEU assesses n-gram precision; chrF++ provides robustness to morphological variation; METEOR accounts for synonymy and paraphrasing; and ROUGE-L quantifies structural coherence through sequence overlap. Each metric provides a

normalized score, with scores reported individually to preserve diagnostic granularity across evaluation dimensions.

For English outputs, we also compute perplexity as an intrinsic measure of fluency (Jelinek et al., 2005; Meister and Cotterell, 2021). Perplexity was not computed for Hindi and Marathi due to the limited availability of validated language models for reliable perplexity estimation in these languages at the time of evaluation. Scores are normalized per response and averaged across languages. To validate the reliability of the metric, a small-scale human evaluation was conducted. Correlations between human ratings and automated scores demonstrated consistent directional alignment, supporting their interpretive adequacy.

Layer 2: Expert-Based Clinical Validation Automated metrics cannot fully capture clinical accuracy or the safety implications of a response. The second layer introduces structured rubric evaluation, grounded in a standardized criterion set validated by domain experts.

Thematic Rubric Development and Validation A panel of clinical experts (see Appendix A) designed and validated the evaluation rubrics. Through iterative content analysis, we identified ten core maternal health themes:

- Antenatal & Maternal Health Care
- Nutrition, Diet & Supplementation
- Mental, Emotional & Social Well-being
- Clinical Procedures & Guidelines
- Medication & Vaccination Safety
- Reproductive & Sexual Health (Beyond Pregnancy)
- Health Systems, Access & Provider Support
- Infection Prevention & Hygiene Practices
- Symptom Interpretation & Danger Sign Recognition
- Risk & Complication Management

Each theme was assigned a uniform 5×3 rubric structure comprising 15 binary criteria organized under five axes: Accuracy, Completeness, Context Awareness, Communication, and Terminology Accessibility.

579 Axis weights were determined through expert
580 consultation: Accuracy = 0.30, Completeness =
581 0.25, Context Awareness = 0.20, Communication
582 = 0.15, and Terminology Accessibility = 0.10.

583 Example (Antenatal & Maternal Health Care
584 theme): under the Accuracy axis, criteria include:

- 585 • Correctly emphasizing the importance and
586 recommended schedule of regular antenatal
587 check-ups.
- 588 • Accurate identification of key assessments to
589 be conducted during pregnancy.
- 590 • Recognition of the roles of healthcare
591 providers in maternal care.

592 The complete rubric framework and thematic
593 codebook appear in Appendix A.

594 The Medical Quality Score (MQS) represents
595 a single unified metric applied consistently across
596 all themes. Each theme utilizes the same five-axis
597 rubric structure with identical weighting, ensur-
598 ing comparability across maternal health topics.
599 Theme-specific rubrics differ only in the concrete
600 criteria under each axis, tailored to the clinical re-
601 quirements of each domain.

602 **Thematic Classification Process** Each ques-
603 tion was automatically assigned to a maternal
604 health theme using a dspy-based classifier imple-
605 menting few-shot structured prompting. DSPy
606 (Differentiable Structured Prompting Yield opti-
607 mizer) refines prompt templates to enhance the-
608 matic consistency across languages (Khattab et al.,
609 2023). Each theme was initialized with 10–15
610 expert-labeled seed examples to calibrate decision
611 boundaries.

612 Model classifications were manually reviewed
613 for a representative subset to ensure alignment with
614 expert expectations, particularly where thematic
615 overlap occurred (e.g., Symptom Interpretation vs.
616 Risk Management). Low-confidence predictions
617 were reviewed by clinicians before final assign-
618 ment. This hybrid process combined efficiency
619 with expert oversight, ensuring the integrity of
620 theme-level labeling.

621 **Automated Rubric Scoring** To enable eval-
622 uation at scale, rubric scoring is performed by an
623 LLM judge (GPT-4o-mini) that receives each ques-
624 tion, its reference answer, the model’s response,
625 and the 15 rubric criteria for that theme, and returns
626 binary pass/fail judgments as structured JSON. The

627 judge runs with deterministic parameters (tempera-
628 ture = 0) for reproducibility. The exact prompt is
629 given in Appendix B.3; we note in the Limitations
630 that the prompt lists the rubric as predicates over a
631 labeled Question/Reference/Response triplet rather
632 than explicitly instructing the judge to compare re-
633 sponse against reference, and we test how much
634 that structural choice matters through cross-judge
635 calibration below.

636 **Cross-Judge Calibration** Because a single
637 judge can bias the absolute MQS (Arora et al.,
638 2025), we re-score a stratified 500-response calibra-
639 tion subset (balanced across 13 models, 2 dataset
640 arms, 3 languages, and 10 themes) with two further
641 judges, Claude Opus 4.6 and GPT-5.1, under the
642 same prompt (Section 5.7). Reporting the pairwise
643 Pearson correlations and mean signed bias between
644 judges lets readers see how much of the reported
645 MQS is judge-calibrated rather than intrinsic to the
646 model. A comparable human-expert scoring pass
647 on this subset is bounded by reviewer capacity and
648 is reported as a planned extension in the Limita-
649 tions; the human validation that was feasible at this
650 stage is concentrated in the upstream creation of
651 reference answers (Section 3), not in per-response
652 scoring of the 44,460 model outputs.

653 **Score Synthesis and Analytical Approach** The
654 two-stage framework produces a multi-dimensional
655 evaluation matrix integrating automated and expert-
656 derived metrics. For each response, five axis-level
657 pass ratios (A_i) are computed. The overall Medical
658 Quality Score (MQS) is defined as:

$$659 \text{MQS} = \sum (w_i \times A_i)$$

660 where weights w_i correspond to expert-defined
661 importance factors: Accuracy = 0.30; Completeness =
662 0.25; Context Awareness = 0.20; Commu-
663 nication = 0.15; Terminology Accessibility = 0.10,
664 with $\sum w_i = 1$.

665 The framework produces separate metrics
666 for semantic fidelity (embedding similarity and
667 BERTScore), linguistic quality (SacreBLEU,
668 chrF++, METEOR, ROUGE-L, and perplexity for
669 English), and clinical validity (MQS). These di-
670 mensions are reported independently rather than
671 combined into a single composite score, as each
672 captures distinct aspects of model performance crit-
673 ical to different deployment considerations. Aggre-
674 gate interpretations consider performance across
675 all three dimensions jointly.

Per-theme and per-language analyses expose asymmetries such as high fluency combined with low medical correctness. Together, the two layers give a transparent, interpretable picture of LLM reliability in maternal healthcare and support the development of safer, context-aware clinical dialogue systems.

Reproducibility and Transparency We release the full pipeline alongside the paper. The supplementary code repository contains the prompt templates, API configurations (temperature, top_p, max_tokens), random seeds, preprocessing scripts, rubric materials, thematic classification examples, and expert evaluation protocols. Generation calls used identical API parameters within each model family (reasoning vs. standard; Appendix B.4), and every generation and judge call is logged to a timestamped JSON Lines file so that individual scores can be re-audited without re-running the API. Computational infrastructure details are in Appendix B.4.

5 Experiments and Results

This section presents the concrete experimental instantiation of the two-layer framework in Section 4 and the findings it produced. The main body carries one headline table, four figures, and the judge-calibration table; the full per-language, per-theme, per-axis, per-cost, and per-validity tables are relocated to Appendix A so that the narrative stays readable and the backing numbers remain available for review.

5.1 Experimental Setup

Stratified model panel. Rather than an ad-hoc collection of currently-popular LLMs, we fix a reproducible sampling frame and draw a stratified panel from it. The eligible frame, as of the snapshot date **2026-04-01**, is any instruction-tuned LLM that (i) offers a public inference API or open weights, (ii) documents generation in English and at least one Indic script, (iii) has been released or substantively updated within the prior 24 months, and (iv) is not a base language model. We stratify the frame along three axes shown in Section 4 to explain variance in the outcomes of interest: *access mode* (proprietary API vs. open-weight), *training specialization* (general-purpose, multilingual-specialized, medical-specialized), and *capability tier* (flagship vs. efficient). From each populated cell we include the provider’s own publicly-positioned flag-

ship, plus, where a matched-generation efficient sibling exists in the same family, the efficient variant as well. This rule yields a 13-model panel spanning five providers, with within-family flagship/efficient pairs for OpenAI (GPT-5 Mini / GPT-4o Mini), Google (Gemini 3 Pro / Gemini 3 Flash / Gemini 3.1 Flash-Lite), Anthropic (Claude Opus 4.7 / Claude Haiku 4.5), Meta (Llama 4 Maverick / Llama 3.3 70B), and Google DeepMind (MedGemma 27B / MedGemma 4B), together with Cohere’s Command A and Aya Expansive 32B. We additionally include Gemma 3 27B as a controlled base-model comparator for MedGemma to allow attribution of medical fine-tuning gains. Empty cells (no open-weight proprietary frontier; no proprietary medical-specialized) are reported rather than hidden.

Two-track dataset. All 13 models are evaluated on two disjoint question sets.

- **Expert (149):** doctor-validated pairs with professionally-translated gold answers in all three languages, released as `sakhi_benchmark_expert.csv`. This is the primary track for clinical-quality claims.
- **Non-expert (231):** community-sourced pairs with doctor-reviewed reference answers, released as `sakhi_benchmark_non_expert.csv`. This track provides broader thematic coverage and supports robustness checks.

The two tracks are reported separately throughout; no score is pooled across them. A response is generated independently for each of three runs per (model, language, question), following Section 4. Total evaluated model outputs are $13 \times 3 \times 3 \times (149 + 231) = 44,460$, minus the open-weight MedGemma variants that the OpenRouter endpoint used to dispatch the other models does not host. MedGemma is evaluated on the non-expert track only and flagged accordingly in every table.

Judge panel and calibration. Our primary rubric judge is **GPT-4o-mini**, which scores every response in both dataset arms under the zero-shot rubric prompt in Appendix B.3. To test how much the findings depend on the judge model, a methodological concern that prior work flagged with a 0.118-point MQS gap between judges on the same responses, we additionally score a **stratified 500-response calibration subset** (balanced across 13 models, 2 datasets, 3 languages, and 10 themes) with two further judges: **Claude Opus 4.6** (via local Claude Code, zero API cost) and **GPT-5.1**

(via OpenRouter). The same zero-shot prompt is used by all three judges. Judge-agreement statistics on this subset quantify how much reported MQS depends on the choice of evaluator.

Generation configuration. Proprietary models are accessed via OpenRouter with temperature=0.7 and max_tokens=400, except reasoning-capable models (GPT-5 family, Gemini 3 family, o-series) which use temperature=1.0, max_tokens=2000, and OpenRouter’s reasoning={effort=low, exclude=true} to prevent reasoning-token exhaustion while keeping reasoning content out of the scored response. Claude Opus 4.7 is invoked via the local claude -p CLI. All generation and judge calls are appended to timestamped per-(model, dataset, lang) JSON Lines files; any interruption resumes by skipping already-completed (run, q_idx) keys without re-spending API tokens. Full configuration, API parameters, and reproducibility scripts live in the supplementary repository.

5.2 Overall Cross-Model Performance

Table 1 reports MQS for each of the 13 models across the two dataset arms and three languages, computed as the mean over three runs times all questions for which the model produced a valid response. *MQS is systematically lower on the expert arm than on the non-expert arm* for every model and every language where both arms are covered. The gap is 0.110 to 0.137 on English cells, 0.079 to 0.105 on Hindi, and 0.063 to 0.109 on Marathi, with an across-cell median of 0.103. The doctor-edited expert references contain more specific clinical content, which the 15 binary rubric criteria correspondingly enforce more strictly. Collapsing the two arms into a single score would obscure this gap; we keep them separate throughout.

Under the expert arm (English), Claude Haiku 4.5 and Claude Opus 4.7 lead with 0.305 and 0.297 respectively, closely followed by GPT-5 Mini at 0.298. *No model exceeds 0.32 on expert-EN; fewer than one in three rubric criteria are reliably met.* Under the non-expert arm (English), the ordering is similar: Claude Haiku 4.5 (0.422), GPT-5 Mini (0.418), Claude Opus 4.7 (partial coverage; see below), Gemini 3 Flash (0.396), and the H2-baseline Gemma 3 27B (0.395) cluster tightly. Llama 3.3 70B and GPT-4o Mini are the weakest non-specialist proprietary / open-weight pair. Medical-specialized MedGemma 27B (0.394) is in-

distinguishable from its base Gemma 3 27B (0.395) on non-expert EN, so medical fine-tuning alone does not lift rubric compliance in this cross-cultural setting. We call this the **H2-baseline null** and the Discussion revisits it.

Figure 1 shows the same non-expert EN ranking with per-model run-to-run error bars computed by grouping each cell’s judge outputs by run index, averaging within each of the three runs, then taking the standard deviation across the three per-run means. Error bars are small (at most ~ 0.015 MQS), so rankings are stable across reruns even where the judge-vs-judge disagreement is large (Section 5.7).

Coverage gaps. Claude Opus 4.7 is the model with the largest coverage footprint in this release. The expert arm has full coverage in English and Hindi but only one complete run for Marathi (n=108 usable responses of the 447 expected), and the non-expert arm has partial coverage only on Marathi (n=113), with non-expert English and non-expert Hindi missing entirely due to sustained Claude Code CLI session limits during the dispatch window. We exclude partial and missing Claude Opus 4.7 cells from the top-3 bolding in Table 1 and show the exact coverage in Appendix Table 7. Aya Expanse 32B was not dispatched to the expert arm because its native Cohere-API route requires credentials we did not have; it is reported on the non-expert arm only. MedGemma 4B and 27B are reported on the non-expert arm only because OpenRouter does not host them and the HuggingFace Inference Endpoint route we set up was scoped to the non-expert track in this release. Gemini 3.1 Flash-Lite is retained in Appendix Table 7 to show that it failed to produce valid output in every cell; it is excluded from scored comparisons elsewhere.

5.3 Theme-Level and Cross-Lingual Structure

Figure 2 summarises the MQS surface across the 10 maternal-health themes and all 13 models on the deployment-facing non-expert arm, averaged over the three languages. Two patterns are visible without having to look at the numbers. First, themes sort left-to-right from strongest (Nutrition, Symptoms) to weakest (Clinical Procedures, Risk Management), and the ordering is consistent across models; every model in the panel is strongest on Nutrition and weakest on Risk Management. Second, models sort top-to-bottom into a tight cluster: the flagship proprietary pair (Claude Haiku 4.5,

Table 1: Medical Quality Score (MQS) by model, dataset, and language under the GPT-4o-mini judge. Cells are the mean across three runs and all questions for which the model produced a valid response. **Bold** marks the best score per column; underline marks the second best. A midrule separates the 13-model evaluated panel from the baseline Gemma 3 27B, which is reported only to anchor the medical-finetune comparison against MedGemma 27B.

Model	Expert (149)			Non-expert (231)		
	EN	HI	MR	EN	HI	MR
GPT-5 Mini	<u>0.298</u>	0.325	0.310	<u>0.418</u>	<u>0.419</u>	0.419
GPT-4o Mini	0.241	0.270	0.250	0.376	0.375	0.359
Command A	0.256	0.302	0.282	0.393	0.381	0.351
Gemini 3 Pro	0.255	0.296	0.289	0.387	0.386	0.382
Gemini 3 Flash	0.266	0.307	0.321	0.396	0.408	0.404
Claude Opus 4.7	0.297	0.354	0.353	–	–	0.363
Claude Haiku 4.5	0.305	<u>0.338</u>	<u>0.343</u>	0.422	0.434	<u>0.406</u>
Llama 3.3 70B	0.233	0.268	0.261	0.369	0.363	0.349
Llama 4 Maverick	0.254	0.280	0.279	0.364	0.376	0.369
Aya Expanse 32B	–	–	–	0.382	0.369	0.337
MedGemma 27B	–	–	–	0.394	0.392	0.383
MedGemma 4B	–	–	–	0.326	0.328	0.331
Gemma 3 27B (baseline)	0.266	0.292	0.281	0.395	0.396	0.388

GPT-5 Mini) and Gemini 3 Flash are visibly lighter than Llama 3.3 70B and MedGemma 4B, but the separation is small relative to the theme-level separation. The within-model row is flatter than the within-theme column, meaning theme-level difficulty explains more variance than model choice.

The cross-lingual story is more subtle and hinges on judge calibration. For 11 of 13 models, Hindi MQS equals or exceeds English MQS under the GPT-4o-mini judge, with ΔHI ranging from -0.013 (Aya Expanse) to $+0.058$ (Claude Opus 4.7) and a median of $+0.014$. Marathi shows a similar pattern, with ΔMR in $[-0.045, +0.061]$ and a median of $+0.005$. We read this result carefully: it does *not* mean LLMs are better at maternal-health reasoning in Hindi or Marathi. The judge-agreement analysis in Section 5.7 shows that GPT-4o-mini is systematically more lenient on non-English outputs, and GPT-5.1 as a stricter judge collapses the EN/HI/MR gap substantially. The conservative interpretation is that apparent cross-lingual parity under a mid-tier judge masks judge leniency toward Devanagari-script responses. Full per-model ΔHI and ΔMR values are in Appendix A.1 (Table 3).

5.4 Axis-Level Diagnostics

Figure 3 reports the per-axis pass rate on the non-expert EN arm. Three patterns dominate.

Communication is uniformly high. Every model scores above 0.58 on Communication (GPT-5 Mini 0.700, Gemini 3 Flash and Gemma 3 27B 0.688, Claude Haiku 4.5 0.689), confirming that

frontier LLMs reliably adopt warm, supportive phrasing and recommend professional consultation.

Terminology Accessibility is the lowest axis for every model on the non-expert EN arm

(0.123–0.206), consistently below Context Awareness in this slice. The same pattern does not hold uniformly on the expert arm, where Completeness is the lowest axis for Claude Opus 4.7 expert-EN and a handful of other (model, language) cells (Appendix A.3). The criteria under Terminology Accessibility cover defining medical abbreviations (ANC, BP, FHR), translating micronutrient and medication names into the local language, and avoiding jargon at first use. Most rubric items in this axis are failed in the non-expert EN slice that the deployed bot most directly reflects. A response can be correct, empathetic, and culturally appropriate while still using the untranslated abbreviation “FHR” that a rural user cannot decode.

Context Awareness is the second-lowest axis

(0.226–0.264), failing criteria about local cultural beliefs, resource constraints, and community-health infrastructure (ASHA workers, anganwadi centres, government health schemes). Combined with the Terminology failures, this axis pair characterises the benchmark’s central finding on the non-expert EN arm: models produce responses that are clinically correct and communicatively warm but *linguistically inaccessible* and *infrastructurally generic*, fitting Western-clinic assumptions rather than rural Indian deployment contexts.

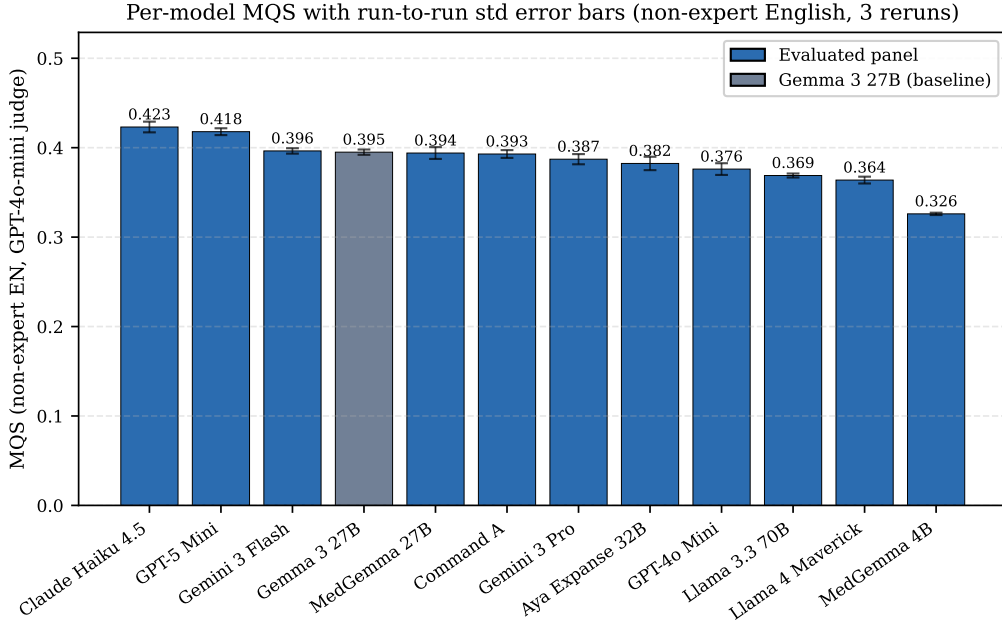


Figure 1: Per-model MQS on the non-expert EN arm with run-to-run std error bars. Means are computed across the 3 reruns \times 231 questions; error bars are the std across the three per-run means. Gemma 3 27B is shown in grey as the matched base-model baseline.

5.5 Cost-Performance Trade-Offs

Figure 4 plots per-response generation cost against non-expert EN MQS, on a log-scaled cost axis. The picture is cleaner than a table: MedGemma 27B and Gemma 3 27B sit at roughly \$0.00007 to \$0.00008 per response with MQS near 0.395, while Claude Haiku 4.5 and Gemini 3 Pro sit at \$0.0015 to \$0.0019 per response with MQS in the 0.387 to 0.422 range. The frontier flattens above \$0.0002 per response: paying 20 \times more moves MQS by at most 0.03. Claude Opus 4.7 is accessed through the local Claude Code subscription at zero marginal per-response cost; it is a research-evaluation convenience, not a fair commercial point, so we exclude it from the cost axis.

The deployment consequence is concrete: for an NGO running a maternal-health bot at million-queries-per-month scale, switching from Claude Haiku 4.5 to Gemma 3 27B saves roughly \$17,000 per month for a -0.027 MQS change. Our axis-level evidence (Figure 3) suggests that amount would produce more safety gain if spent on Hindi and Marathi medical-terminology alignment than on a frontier-tier model.

5.6 Generation Reliability

Not every response logged in a run constitutes a valid answer: some are empty, truncated, or return

Table 2: Three-judge agreement on the stratified 500-response calibration subset ($n = 383$ rows with all three judges present). Pearson ρ is computed on MQS and **bold** marks the judge pair with highest rank agreement. Mean Δ MQS is the signed bias; underline marks the smallest absolute bias.

Pair	Pearson ρ	Mean Δ MQS
GPT-4o-mini \leftrightarrow Claude Opus 4.6	0.429	<u>+0.037</u>
GPT-4o-mini \leftrightarrow GPT-5.1	0.238	+0.201
Claude Opus 4.6 \leftrightarrow GPT-5.1	0.338	+0.163

refusals. Failure rates are highest in Marathi, reflecting the familiar brittleness of LLM generation in lower-resource scripts; Llama-family models in particular show degraded Marathi coverage. Gemini 3.1 Flash-Lite fails in every cell of the panel and is excluded from scored comparisons. We report the full per-(dataset, language) valid-response rate in Appendix Table 7; deployment decisions should pair MQS with this “generation dropout” rate, since a model that scores well on valid responses is only useful to the extent that valid responses come back in the first place.

5.7 Judge-Agreement Calibration

Our primary GPT-4o-mini scores are validated on a stratified 500-response calibration subset ($n = 383$ rows with all three judges present) by two further

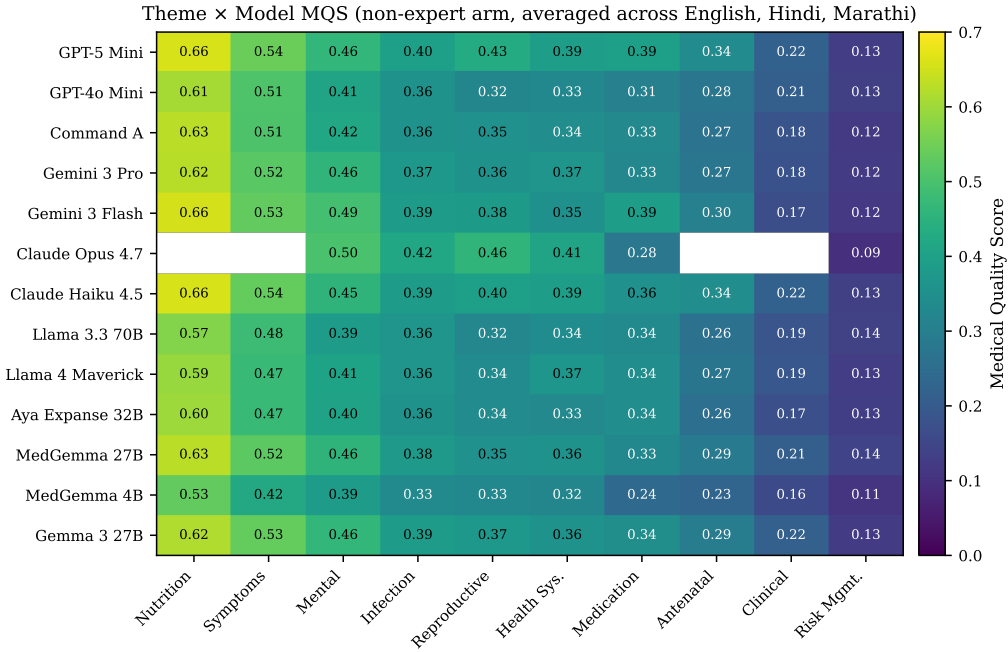


Figure 2: Theme \times Model MQS (non-expert arm, averaged across English, Hindi, and Marathi under the GPT-4o-mini judge). Themes are sorted left-to-right by cross-model mean. Blank cells for Claude Opus 4.7 on non-expert EN / HI reflect the coverage gap described in Section 5.2. Full numeric table in Appendix A.2 (Table 4).

judges, Claude Opus 4.6 and GPT-5.1. Table 2 reports pairwise Pearson correlations on MQS and the mean signed difference (bias) between each pair.

Two effects emerge. First, **GPT-5.1 is systematically stricter** than both other judges by roughly 0.17 to 0.20 MQS on identical responses. Every model’s headline score in Table 1 would shrink by that offset under the stricter evaluator. This bias is larger than the MQS gap between most pairs of models in our panel, so absolute rankings should be interpreted in light of judge choice. Second, **rank correlations are modest** ($\rho = 0.24$ GPT-4o-mini \leftrightarrow GPT-5.1; $\rho = 0.43$ GPT-4o-mini \leftrightarrow Claude Opus 4.6), meaning that judges disagree not only on absolute level but also on which responses are stronger. GPT-4o-mini and Claude Opus 4.6 agree most closely on bias ($\Delta = +0.037$), matching the pattern noted in prior work that mid- and frontier proprietary models converge on lenient rubric interpretations, while GPT-5.1 applies a more discerning standard.

We read the absolute MQS values in Table 1 as *GPT-4o-mini-calibrated* scores. Under the GPT-5.1 judge the model ordering is broadly preserved (Claude Haiku 4.5 and GPT-5 Mini remain top) but the magnitudes compress by a known, reported offset. The judge panel is itself the methodological

contribution: reporting a single judge’s scores without this triangulation understates the uncertainty that deployers should account for.

5.8 Summary of Findings

The stratified 13-model evaluation produces four findings that hold across both dataset arms:

- Rubric-compliance is low overall.** No model exceeds MQS 0.32 on expert-EN under GPT-4o-mini, and under the stricter GPT-5.1 judge the panel-wide mean compresses by a further ~ 0.17 – 0.20 MQS. Frontier LLMs meet fewer than one in three doctor-edited rubric criteria (§5.2, §5.7).
- On the deployment-facing non-expert EN arm, the dominant failure is linguistic.** Terminology Accessibility is the lowest axis for every model in this slice (pass rates 0.123–0.206), below Context Awareness (0.226–0.264). Responses are empathetic and clinically accurate but use untranslated medical shorthand and Western-clinic infrastructure assumptions. On the stricter expert arm, Completeness competes with Terminology Accessibility for the lowest axis in some cells (§5.4).
- Apparent cross-lingual parity is a judge artifact.** Under GPT-4o-mini, 11 of 13 models score equal or higher on Hindi than English.

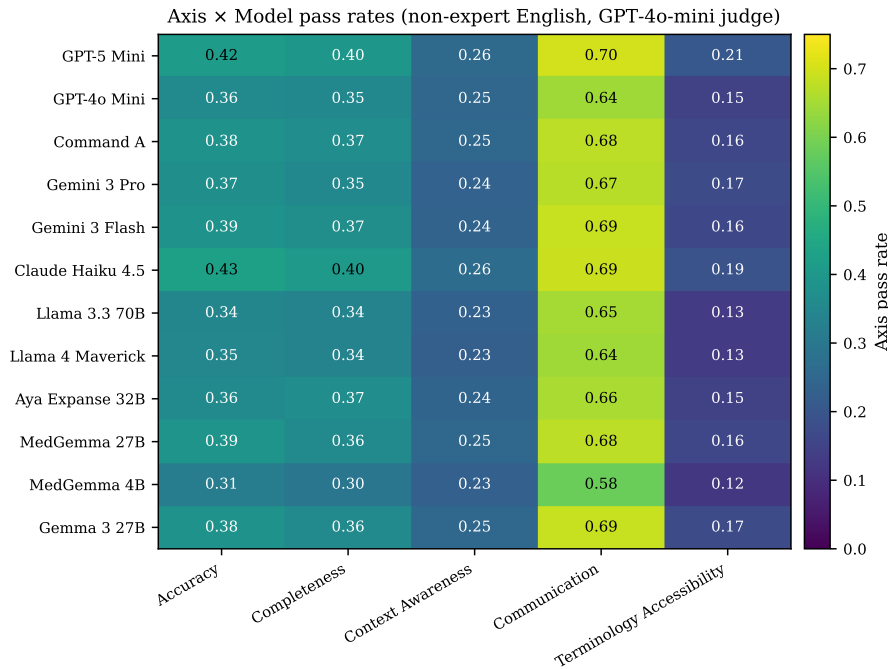


Figure 3: Axis × Model pass rates on the non-expert English arm (GPT-4o-mini judge). Each cell is the mean fraction of binary rubric criteria satisfied within that axis. Communication is uniformly high (0.58–0.70); Terminology Accessibility is uniformly low (0.12–0.21) across every model. Full numeric table in Appendix A.3 (Table 5).

Under GPT-5.1 (stricter; §5.7) the compression is not matched in the Indic direction, suggesting the parity is partly an artifact of mid-tier-judge leniency toward non-English text rather than genuine cross-lingual competence (§5.3).

4. **Open-weight medium-scale dominates the cost-performance frontier.** MedGemma 27B (\$0.00007 per response, MQS 0.394) and Gemma 3 27B (\$0.00008 per response, MQS 0.395) reach non-expert EN MQS matching proprietary flagships at over an order of magnitude lower per-response cost. The MQS/\$ curve flattens above \$0.0002 per response (§5.5).

Discussion (Section 6) revisits the implications of each finding for deployment of AI systems in rural multilingual health contexts and the methodological requirements for multi-judge clinical benchmarking.

6 Discussion

Clinically correct, linguistically inaccessible, infrastructurally generic. Across the full 13-model panel, one pattern holds in every cell: the responses are clinically accurate and warmly phrased, yet they assume a Western clinic and an English-speaking reader. No model exceeds MQS 0.32 on the expert arm in English under the GPT-4o-mini

judge, and under the stricter GPT-5.1 judge (§5.7) every model loses a further 0.17 to 0.20 MQS on identical outputs. Fewer than one in three doctor-validated rubric criteria are reliably met. The axis decomposition tells us where the shortfall lives, and where it does not. Communication is uniformly high, above 0.58 for every model (§5.4), so warm tone and the habit of recommending a doctor are not the bottleneck. Accuracy and Completeness pass rates cluster in the middle of the panel, neither the failure mode nor the strong suit. Context Awareness and, more sharply, Terminology Accessibility are the weakest axes. What this looks like on the ground: a response that correctly flags pre-eclampsia danger signs but uses the untranslated shorthand BP, never gives the Hindi or Marathi word for blood pressure, and tells the reader to visit a “healthcare provider” without mentioning the ASHA worker who is, for most of our target users, the nearest healthcare contact at all.

The dominant failure is linguistic, not just cultural. Our earlier two-model analysis, and most prior work in this space, described the shortfall as a “cultural awareness” gap. The 13-model axis decomposition locates it more precisely, and the location is actionable. Terminology Accessibility pass rates sit between 0.123 and 0.206, lower than

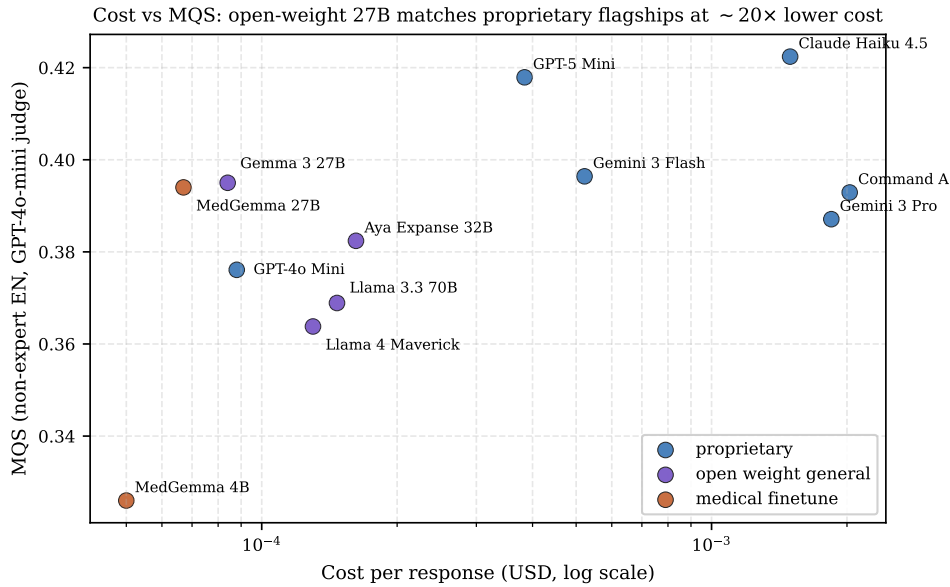


Figure 4: Cost per response (USD, log scale) versus non-expert EN MQS under the GPT-4o-mini judge. Colour encodes training family: proprietary API, open-weight general-purpose, and medical-fine-tuned. The efficient frontier is dominated by open-weight 27B models; the proprietary cluster (top right) does not separate from them on MQS despite costing $\sim 20\times$ more. Full numeric table in Appendix A.4 (Table 6).

Context Awareness (0.226–0.264) for every model in the panel. The failures are lexical: English abbreviations left untranslated at first use, Indian generic drug names rendered in Latin script instead of Devanagari, and metric shorthand that stays in English even when the question is posed in Marathi. Pointing at “culture” is too diffuse to act on. Pointing at terminology is actionable: curated parallel medical glossaries, abbreviation expansion at first use, and alignment fine-tuning on Indian Pharmacopoeia drug names are each narrow, measurable interventions. This sharpens the CLINIC finding that multilingual trustworthiness varies across scripts (Ghosh et al., 2025) by locating the weakness at the lexical level rather than at the reasoning level.

Apparent cross-lingual parity is a judge artifact. Under GPT-4o-mini, 11 of 13 models score equal or higher on Hindi than on English, and Marathi shows the same pattern for most of the panel (§5.3). Taken alone this would imply that frontier LLMs handle Indic-script maternal-health queries at least as well as English queries. Our judge-agreement analysis shows the result is not load-bearing. GPT-5.1 applied to the stratified 500-response subset scores every model 0.17 to 0.20 MQS lower on identical outputs, and the offset is not uniform across languages. Rank correlations across judges are modest: Pearson 0.24 between GPT-4o-mini and GPT-5.1, and 0.43 between GPT-4o-mini and

Claude Opus 4.6. The judges disagree on absolute level and, separately, on which responses are better. Any paper that reports apparent cross-lingual parity on the basis of a single mid-tier judge is reporting a potential artifact. At minimum, multilingual medical benchmarks should triangulate across one stricter judge and publish the per-language bias offset alongside the headline scores.

Open-weight medium-scale models are the rational deployment choice. The cost-performance frontier (§5.5) is unambiguous for the operator profile that motivates this benchmark, namely NGOs running maternal-health bots at rural-India scale. MedGemma 27B and Gemma 3 27B both reach MQS ≈ 0.395 on the non-expert English arm at roughly \$0.00008 per response. Claude Haiku 4.5, the panel’s strongest proprietary model, reaches 0.422 at about \$0.0015 per response, so the extra \$0.00142 per call buys +0.027 MQS. For a single bot serving one million queries per month, that is \$17,040 per month for three percentage points of rubric compliance, money that, on the evidence in §5.4, would produce a larger safety gain if spent on Hindi and Marathi medical-terminology alignment instead. MedGemma 27B and its base Gemma 3 27B are indistinguishable on non-expert EN (MQS 0.394 vs. 0.395), a controlled matched-pair result that confirms the hypothesis we called H2-baseline: medical fine-tuning alone does not lift rubric com-

1148 pliance in this cross-cultural setting. The lever is
1149 translation and terminology, not medical special-
1150 ization.

1151 **Automated metrics are not a substitute for the**
1152 **rubric.** Embedding similarity and BERTScore
1153 can rate an answer highly on topical alignment
1154 while the rubric flags the same response as miss-
1155 ing specific clinical criteria that matter for patient
1156 safety: drug-interaction warnings, the correct ante-
1157 natal visit schedule, the right danger-sign threshold
1158 for seeking immediate care. That is the pattern
1159 prior work on metric-versus-clinical-quality diver-
1160 gence already documents (Tam et al., 2024; Crox-
1161 ford et al., 2025), and it is why the two-layer design
1162 is the floor rather than an optional add-on. The au-
1163 tomated layer scales and reproduces; the rubric is
1164 what the hospital or the NGO programme officer
1165 will actually read before deciding to deploy.

1166 **Situating Sakhi among related benchmarks.**
1167 HealthBench (Arora et al., 2025) is the closest
1168 methodological comparator, with rubric-style eval-
1169 uation across 26 medical specialties, but it is pre-
1170 dominantly English and uses synthetic conversa-
1171 tions rather than community-validated questions.
1172 CLINIC (Ghosh et al., 2025) covers 15 languages
1173 but assesses generic trustworthiness, not domain-
1174 specific clinical quality. Script Gap (Khullar et al.,
1175 2025) evaluates triage classification accuracy on
1176 Indian-language maternal-health conversations but
1177 stops short of open-ended generation quality. Sakhi
1178 contributes three pieces that the prior benchmarks
1179 do not jointly offer. First, reference answers in
1180 all three languages on a single question set, both
1181 doctor-edited (expert arm) and community-sourced
1182 (non-expert arm), so the MQS gap between them is
1183 measurable within the same model. Second, a strat-
1184 ified 13-model panel drawn from a pre-declared eli-
1185 gibility frame, with within-family flagship/efficient
1186 pairs and a base-versus-medical-finetune matched
1187 comparator that lets us attribute, or refuse to at-
1188 tribute, gains to medical fine-tuning. Third, a three-
1189 judge calibration on a stratified 500-response sub-
1190 set that quantifies how much the headline MQS de-
1191 pends on the judge model. We read these as depth
1192 where existing work offers breadth. The paper is
1193 written for the NeurIPS Datasets and Benchmarks
1194 track, and the intended users are NGO deployers
1195 and research teams who need an evaluation that
1196 fails for the same reasons patients actually get hurt.

1197 Limitations

1198 Several limitations shape what our findings can
1199 and cannot support. First, the benchmark covers a
1200 focused, doctor-validated subset of maternal and
1201 reproductive health queries (149 expert-edited plus
1202 231 community-sourced questions) and does not
1203 exhaust the full space of questions a bot will re-
1204 ceive in deployment. Second, Hindi and Marathi
1205 are represented through professional translations
1206 of English questions and reference answers. This
1207 controls for topic drift across languages at the
1208 cost of not reflecting the full variation of spon-
1209 taneous Indic-script patient queries, which often
1210 code-switch between Hindi and English, use non-
1211 standard transliteration, and lean on regional idiom.
1212 A fully native-authored Indic-first question set is
1213 the natural next dataset release.

1214 Third, several coverage gaps in the 13-model
1215 panel are reported rather than hidden. Claude Opus
1216 4.7 has full coverage on the expert English and
1217 expert Hindi cells, one complete run on expert
1218 Marathi (n=108), partial coverage on non-expert
1219 Marathi (n=113), and no valid responses on non-
1220 expert English or non-expert Hindi, all driven by
1221 sustained Claude Code CLI session limits during
1222 the dispatch window. Aya Expanse 32B was dis-
1223 patched on the non-expert arm only because its na-
1224 tive Cohere-API route required credentials we did
1225 not have available for the expert track. MedGemma
1226 4B and 27B are reported on the non-expert arm only
1227 because the OpenRouter endpoint used to dispatch
1228 the other models does not host them and the Hug-
1229 gingFace Inference Endpoints route we set up for
1230 MedGemma was scoped to the non-expert track in
1231 this release. Gemini 3.1 Flash-Lite produced no
1232 valid output in any cell and is excluded from scored
1233 comparisons but kept in Table 7 so the failure is
1234 visible. None of these gaps changes the panel-
1235 wide ranking, but each qualifies the specific cell it
1236 touches.

1237 Fourth, the rubric-judge prompt lists the rubric
1238 criteria as predicates over a labeled Question / Ref-
1239 erence / Response triplet and asks for binary JSON
1240 scores, but does not explicitly instruct the judge to
1241 compare the response against the reference answer
1242 (Appendix B.3). The design is deliberate: each
1243 criterion is written as a predicate that a correct re-
1244 sponse should satisfy (for example, “Correctly em-
1245 phasizing the recommended schedule of regular an-
1246 tenatal check-ups”), and the reference is provided
1247 as clinically-grounded context rather than as a gold-

standard target to be matched. This differs from HealthBench (Arora et al., 2025), whose physician-designed rubrics are tied to conversation-specific criteria scored against a single reference. The trade-off is that a judge is left to infer the evaluation task from the structural layout of the prompt, and the three-judge calibration on the 500-response subset (Section 5.7) is what tells us how much that under-specification matters in practice. Pearson $\rho = 0.24$ between GPT-4o-mini and GPT-5.1 means the two judges disagree not just on strict-vs-lenient calibration but on the relative ordering of responses. A revision of the prompt that explicitly asks the judge to compare response against reference is one of the first follow-ups we would run, and we kept the prompt unchanged across this release to preserve comparability with earlier scored runs.

Fifth, the primary rubric judge is GPT-4o-mini. The cross-judge calibration with Claude Opus 4.6 and GPT-5.1 on the stratified 500-response subset (Section 5.7) quantifies the bias this introduces, and we report judge identity and offset throughout, but the full $13 \times 2 \times 3$ grid is scored by a single primary judge. Extending the three-judge calibration to the entire grid is bounded by API cost, not by methodology.

Sixth, human reviewer capacity is the hardest constraint on this release and the reason this paper is about LLM-judge rubric scores rather than human rubric scores at scale. The review platform described in Section 3 was sized for doctor-edited reference answers on the 380 pairs in the two released tracks and for ASHA-worker and nonprofit-staff review of those pairs, not for per-response scoring of the 44,460 generated model outputs. A capacity-bounded human-judge pass that mirrors the LLM three-judge calibration, covering a stratified subset balanced across the 13 models, the 10 themes, and the 3 languages, is the critical next data release for tightening the MQS estimates and the judge-agreement analysis. We will report the sampling design, theme distribution, and per-cell N alongside that release rather than reconstruct a human-vs-LLM agreement claim that the current data does not support.

Seventh, the rubric is binary (pass/fail on 15 criteria per theme). This sharpens agreement but misses the difference between a criterion that is met weakly and one that is met well. A graded follow-up rubric is left to future work. Eighth, responses are scored single-turn, without clarifying exchanges. Real maternal-health conversations in-

volve follow-up and disambiguation, and a model that would have clarified under a multi-turn protocol is penalized here for its single-turn output. Finally, a benchmark is necessary but not sufficient for deployment. Any real-world use of these models in maternal-health advice requires additional safety validation, regulatory and ethics review, and ongoing monitoring that sits outside the scope of this paper.

Acknowledgments

References

- National Antenatal Care Guideline 2022.pdf. 1311
- Mahyar Abbasian, Elahe Khatibi, Iman Azimi, David Oniani, Zahra Shakeri Hossein Abad, Alexander Thieme, Ram Sriram, Zhongqi Yang, Yanshan Wang, Bryant Lin, Olivier Gevaert, Li-Jia Li, Ramesh Jain, and Amir M. Rahmani. 2024. [Foundation metrics for evaluating effectiveness of healthcare conversations powered by generative AI](#). *NPJ Digital Medicine*, 7:82. 1312-1317
- Iñigo Alonso, Maite Oronoz, and Rodrigo Agerri. 2024. [MedExpQA: Multilingual benchmarking of Large Language Models for Medical Question Answering](#). *Artificial Intelligence in Medicine*, 155:102938. 1320-1322
- Lameck Mbangula Amugongo, Pietro Mascheroni, Steven Brooks, Stefan Doering, and Jan Seidel. 2025. [Retrieval augmented generation for large language models in healthcare: A systematic review](#). *PLOS digital health*, 4(6):e0000877. 1325-1328
- D. M. Anisuzzaman, Jeffrey G. Malins, Paul A. Friedman, and Zach I. Attia. 2025. [Fine-Tuning Large Language Models for Specialized Use Cases](#). *Mayo Clinic Proceedings: Digital Health*, 3(1):100184. 1329-1330
- Maria Antoniak, Aakanksha Naik, Carla S. Alvarado, Lucy Lu Wang, and Irene Y. Chen. 2024. [NLP for maternal healthcare: Perspectives and guiding principles in the age of LLMs](#). In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency (FAccT '24)*, New York, NY, USA. Association for Computing Machinery. 1333-1338
- Rahul K. Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, Johannes Heidecke, and Karan Singhal. 2025. [HealthBench: Evaluating Large Language Models Towards Improved Human Health](#). *arXiv preprint*. ArXiv:2505.08775 [cs]. 1340-1346
- Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, 1347-1351

1353	Michigan. Association for Computational Linguistics.	1410
1354		1411
1355	Joachim Baumann, Paul Röttger, Aleksandra Urman, Albert Wendsjö, Flor Miriam Plaza-del Arco, Johannes B. Gruber, and Dirk Hovy. 2025. Large Language Model Hacking: Quantifying the Hidden Risks of Using LLMs for Text Annotation . <i>arXiv preprint</i> . ArXiv:2509.08825 [cs].	1412
1356		1413
1357		1414
1358		1415
1359		1416
1360		
1361	Anagani Bhanusree, Sai Divya Vissamsetty, K VenkataKrishna Rao, and Rimjhim. 2026. Comparative analysis of large language models in generating Telugu responses for maternal health queries . <i>Preprint</i> , arXiv:2603.18898.	1417
1362		1418
1363		1419
1364		1420
1365		1421
1366	Iva Bojic, Jessica Chen, Si Yuan Chang, Qi Chwen Ong, Shafiq Joty, and Josip Car. 2023. Hierarchical Evaluation Framework: Best Practices for Human Evaluation . In <i>Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems</i> , pages 11–22, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.	1422
1367		1423
1368		1424
1369		1425
1370		
1371		
1372	Chris Callison-Burch, Miles Osborne, and Philipp Koehn. 2006. Re-evaluating the Role of Bleu in Machine Translation Research . In <i>11th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 249–256, Trento, Italy. Association for Computational Linguistics.	1426
1373		1427
1374		1428
1375		1429
1376		1430
1377		
1378	Xiaolan Chen, Jiayang Xiang, Shanfu Lu, Yexin Liu, Mingguang He, and Danli Shi. 2025. Evaluating large language models and agents in healthcare: key challenges in clinical applications . <i>Intelligent Medicine</i> , 5(2):151–163.	1431
1379		1432
1380		1433
1381		1434
1382		1435
1383	Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. 2023. MEDITRON-70B: Scaling Medical Pretraining for Large Language Models . Publisher: arXiv Version Number: 1.	1436
1384		1437
1385		1438
1386		1439
1387		1440
1388		1441
1389		1442
1390		1443
1391		1444
1392		1445
1393	Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All That’s ‘Human’ Is Not Gold: Evaluating Human Evaluation of Generated Text . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 7282–7296, Online. Association for Computational Linguistics.	1446
1394		1447
1395		1448
1396		1449
1397		
1398		
1399		
1400		
1401		
1402	Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. XNLI: Evaluating Cross-lingual Sentence Representations . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2475–2485, Brussels, Belgium. Association for Computational Linguistics.	1450
1403		1451
1404		1452
1405		1453
1406		1454
1407		1455
1408		1456
1409		
	Emma Croxford, Yanjun Gao, Nicholas Pellegrino, Karen Wong, Graham Wills, Elliot First, Frank Liao, Cherodeep Goswami, Brian Patterson, and Majid Afshar. 2025. Current and future state of evaluation of large language models for medical summarization tasks . <i>npj Health Systems</i> , 2(1):6. Publisher: Nature Publishing Group.	1457
		1458
		1459
		1460
		1461
	John Dang, Shivalika Singh, Daniel D’souza, Arash Ahmadian, Alejandro Salamanca, Madeline Smith, Aidan Peppin, Sungjin Hong, Manoj Govindassamy, Terrence Zhao, Sandra Kublik, Meor Amer, Viraat Aryabumi, Jon Ander Campos, Yi-Chern Tan, Tom Kocmi, Florian Strub, Nathan Grinsztajn, Yannis Flet-Berliac, and 26 others. 2024. Aya Expand: Combining Research Breakthroughs for a New Multilingual Frontier . <i>arXiv preprint</i> . ArXiv:2412.04261 [cs].	1462
		1463
		1464
		1465
		1466
	Goutam Datta, Nisheeth Joshi, and Kusum Gupta. 2022. Analysis of Automatic Evaluation Metric on Low-Resourced Language: BERTScore vs BLEU Score . In <i>Speech and Computer</i> , pages 155–162, Cham. Springer International Publishing.	1467
		1468
		1469
		1470
	Samruddhi Deode, Janhavi Gadre, Aditi Kajale, Ananya Joshi, and Raviraj Joshi. 2023. L3Cube-IndicSBERT: A simple approach for learning cross-lingual sentence representations using multilingual BERT . <i>arXiv preprint</i> . ArXiv:2304.11434 [cs].	1471
		1472
		1473
		1474
		1475
	Indian Express. 2025. Top 10 most spoken languages globally in 2025: 5 indian languages ranked . Accessed: 2025-11-03.	1476
		1477
		1478
	David Fraile Navarro, Enrico Coiera, Thomas W. Hambly, Zoe Triplett, Nahyan Asif, Anindya Susanto, Anamika Chowdhury, Amaya Azcoaga Lorenzo, Mark Dras, and Shlomo Berkovsky. 2025. Expert evaluation of large language models for clinical dialogue summarization . <i>Scientific Reports</i> , 15(1):1195. Publisher: Nature Publishing Group.	1479
		1480
		1481
		1482
		1483
		1484
		1485
		1486
	Agasthya Gangavarapu. 2024. Introducing L2M3, a multilingual medical large language model to advance health equity in low-resource regions . <i>Preprint</i> , arXiv:2404.08705.	1487
		1488
		1489
		1490
		1491
		1492
		1493
		1494
		1495
		1496
		1497
		1498
		1499
		1500
		1501
		1502
		1503
		1504
		1505
		1506
		1507
		1508
		1509
		1510
		1511
		1512
		1513
		1514
		1515
		1516
		1517
		1518
		1519
		1520
		1521
		1522
		1523
		1524
		1525
		1526
		1527
		1528
		1529
		1530
		1531
		1532
		1533
		1534
		1535
		1536
		1537
		1538
		1539
		1540
		1541
		1542
		1543
		1544
		1545
		1546
		1547
		1548
		1549
		1550
		1551
		1552
		1553
		1554
		1555
		1556
		1557
		1558
		1559
		1560
		1561
		1562
		1563
		1564
		1565
		1566
		1567
		1568
		1569
		1570
		1571
		1572
		1573
		1574
		1575
		1576
		1577
		1578
		1579
		1580
		1581
		1582
		1583
		1584
		1585
		1586
		1587
		1588
		1589
		1590
		1591
		1592
		1593
		1594
		1595
		1596
		1597
		1598
		1599
		1600
		1601
		1602
		1603
		1604
		1605
		1606
		1607
		1608
		1609
		1610
		1611
		1612
		1613
		1614
		1615
		1616
		1617
		1618
		1619
		1620
		1621
		1622
		1623
		1624
		1625
		1626
		1627
		1628
		1629
		1630
		1631
		1632
		1633
		1634
		1635
		1636
		1637
		1638
		1639
		1640
		1641
		1642
		1643
		1644
		1645
		1646
		1647
		1648
		1649
		1650
		1651
		1652
		1653
		1654
		1655
		1656
		1657
		1658
		1659
		1660
		1661
		1662
		1663
		1664
		1665
		1666
		1667
		1668
		1669
		1670
		1671
		1672
		1673
		1674
		1675
		1676
		1677
		1678
		1679
		1680
		1681
		1682
		1683
		1684
		1685
		1686
		1687
		1688
		1689
		1690
		1691
		1692
		1693
		1694
		1695
		1696
		1697
		1698
		1699
		1700

1578	<i>Machine Translation: Research Papers</i> , pages 186–191, Brussels, Belgium. Association for Computational Linguistics.	1635
1579		1636
1580		1637
1581	Nils Reimers and Iryna Gurevych. 2019. SentenceBERT: 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.	1638
1582		1639
1583		1640
1584		1641
1585		1642
1586		1643
1587		1644
1588		1645
1589	David Restrepo, Chenwei Wu, Zhengxu Tang, Zitao Shuai, Thao Nguyen Minh Phan, Jun-En Ding, Cong-Tinh Dao, Jack Gallifant, and 1 others. 2025. MultiOphthaLingua: A multilingual benchmark for assessing and debiasing LLM ophthalmological QA in LMICs . In <i>Proceedings of the AAAI Conference on Artificial Intelligence (AAAI Track)</i> .	1646
1590		1647
1591		1648
1592		1649
1593		1650
1594		1651
1595		1652
1596	Pritika Rohera, Chaitrali Ginimav, Gayatri Sawant, and Raviraj Joshi. 2025. Better to ask in english? evaluating factual accuracy of multilingual llms in english and low-resource languages . <i>Preprint</i> , arXiv:2504.20022.	1653
1597		1654
1598		1655
1599		1656
1600		1657
1601	Ananya B. Sai, Akash Kumar Mohankumar, and Mitesh M. Khapra. 2022. A Survey of Evaluation Metrics Used for NLG Systems . <i>ACM Comput. Surv.</i> , 55(2):26:1–26:39.	1658
1602		1659
1603		1660
1604		1661
1605	Thomas Savage, Stephen P. Ma, Abdessalem Boukil, Ekanath Rangan, Vishwesh Patel, Ivan Lopez, and Jonathan Chen. 2025. Fine-Tuning Methods for Large Language Models in Clinical Medicine by Supervised Fine-Tuning and Direct Preference Optimization: Comparative Evaluation . <i>Journal of Medical Internet Research</i> , 27(1):e76048. Company: Journal of Medical Internet Research Distributor: Journal of Medical Internet Research Institution: Journal of Medical Internet Research Label: Journal of Medical Internet Research Publisher: JMIR Publications Inc., Toronto, Canada.	1662
1606		1663
1607		1664
1608		1665
1609		1666
1610		1667
1611		1668
1612		1669
1613		1670
1614		1671
1615		1672
1616		1673
1617	Hendrik Schuff, Lindsey Vanderlyn, Heike Adel, and Ngoc Thang Vu. 2023. How to do human evaluation: A brief introduction to user studies in NLP . <i>Natural Language Engineering</i> , 29(5):1199–1222.	1674
1618		1675
1619		1676
1620		1677
1621	Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cían Hughes, Charles Lau, Justin Chen, Fereshteh Mahvar, Liron Yatziv, Tiffany Chen, Bram Sterling, Stefanie Anna Baby, Susanna Maria Baby, Jeremy Lai, Samuel Schmidgall, and 62 others. 2025. MedGemma Technical Report . <i>arXiv preprint</i> . ArXiv:2507.05201 [cs].	1678
1622		1679
1623		1680
1624		1681
1625		1682
1626		1683
1627		1684
1628		1685
1629		1686
1629	Junhyuk Seo, Dasol Choi, Taerim Kim, Won Chul Cha, Minha Kim, Haanju Yoo, Namkee Oh, YongJin Yi, Kye Hwa Lee, and Edward Choi. 2024. Evaluation Framework of Large Language Models in Medical Documentation: Development and Usability Study . <i>Journal of Medical Internet Research</i> , 26(1):e58329.	1687
1630		1688
1631		1689
1632		1690
1633		1691
1634		1692
	Company: Journal of Medical Internet Research Distributor: Journal of Medical Internet Research Institution: Journal of Medical Internet Research Label: Journal of Medical Internet Research Publisher: JMIR Publications Inc., Toronto, Canada.	1693
		1694
		1695
		1696
		1697
		1698
		1699
		1700
		1701
		1702
		1703
		1704
		1705
		1706
		1707
		1708
		1709
		1710
		1711
		1712
		1713
		1714
		1715
		1716
		1717
		1718
		1719
		1720
		1721
		1722
		1723
		1724
		1725
		1726
		1727
		1728
		1729
		1730
		1731
		1732
		1733
		1734
		1735
		1736
		1737
		1738
		1739
		1740
		1741
		1742
		1743
		1744
		1745
		1746
		1747
		1748
		1749
		1750
		1751
		1752
		1753
		1754
		1755
		1756
		1757
		1758
		1759
		1760
		1761
		1762
		1763
		1764
		1765
		1766
		1767
		1768
		1769
		1770
		1771
		1772
		1773
		1774
		1775
		1776
		1777
		1778
		1779
		1780
		1781
		1782
		1783
		1784
		1785
		1786
		1787
		1788
		1789
		1790
		1791
		1792
		1793
		1794
		1795
		1796
		1797
		1798
		1799
		1800
		1801
		1802
		1803
		1804
		1805
		1806
		1807
		1808
		1809
		1810
		1811
		1812
		1813
		1814
		1815
		1816
		1817
		1818
		1819
		1820
		1821
		1822
		1823
		1824
		1825
		1826
		1827
		1828
		1829
		1830
		1831
		1832
		1833
		1834
		1835
		1836
		1837
		1838
		1839
		1840
		1841
		1842
		1843
		1844
		1845
		1846
		1847
		1848
		1849
		1850
		1851
		1852
		1853
		1854
		1855
		1856
		1857
		1858
		1859
		1860
		1861
		1862
		1863
		1864
		1865
		1866
		1867
		1868
		1869
		1870
		1871
		1872
		1873
		1874
		1875
		1876
		1877
		1878
		1879
		1880
		1881
		1882
		1883
		1884
		1885
		1886
		1887
		1888
		1889
		1890
		1891
		1892
		1893
		1894
		1895
		1896
		1897
		1898
		1899
		1900
		1901
		1902
		1903
		1904
		1905
		1906
		1907
		1908
		1909
		1910
		1911
		1912
		1913
		1914
		1915
		1916
		1917
		1918
		1919
		1920
		1921
		1922
		1923
		1924
		1925
		1926
		1927
		1928
		1929
		1930
		1931
		1932
		1933
		1934
		1935
		1936
		1937
		1938
		1939
		1940
		1941
		1942
		1943
		1944
		1945
		1946
		1947
		1948
		1949
		1950
		1951
		1952
		1953
		1954
		1955
		1956
		1957
		1958
		1959
		1960
		1961
		1962
		1963
		1964
		1965
		1966
		1967
		1968
		1969
		1970
		1971
		1972
		1973
		1974
		1975
		1976
		1977
		1978
		1979
		1980
		1981
		1982
		1983
		1984
		1985
		1986
		1987
		1988
		1989
		1990
		1991
		1992
		1993
		1994
		1995
		1996
		1997
		1998
		1999
		2000

Grau. 2025. [Medical Graph RAG: Evidence-based Medical Large Language Model via Graph Retrieval-Augmented Generation](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 28443–28467, Vienna, Austria. Association for Computational Linguistics.

Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, Nitesh V. Chawla, and Xiangliang Zhang. 2024. [Justice or Prejudice? Quantifying Biases in LLM-as-a-Judge](#). *arXiv preprint*. ArXiv:2410.02736 [cs].

Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. [MoverScore: Text Generation Evaluating with Contextualized Embeddings and Earth Mover Distance](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 563–578, Hong Kong, China. Association for Computational Linguistics.

Hongli Zhou, Hui Huang, Yunfei Long, Bing Xu, Conghui Zhu, Hailong Cao, Muyun Yang, and Tiejun Zhao. 2024. [Mitigating the Bias of Large Language Model Evaluation](#). In *Proceedings of the 23rd Chinese National Conference on Computational Linguistics (Volume 1: Main Conference)*, pages 1310–1319, Taiyuan, China. Chinese Information Processing Society of China.

A Full Results Tables

This appendix holds the full numeric tables underpinning the figures in Section 5. Each table is produced by `scoring/pipeline/render_tables.py` from the aggregated JSONL judge logs; re-running that script regenerates every number.

A.1 Cross-Lingual MQS

A.2 MQS by Theme

A.3 Axis Pass Rates, Full Panel

A.4 Cost-Performance Details

A.5 Generation Validity

B Prompt Templates and API Configuration

All models were evaluated using the same prompt templates to ensure comparability. No few-shot examples were provided to the evaluated models or to the rubric judge, following a zero-shot design to avoid biasing evaluation toward specific response patterns.

Table 3: Cross-lingual MQS averaged over both dataset arms. Δ HI and Δ MR show the signed change from English; positive means the model scores higher in Hindi or Marathi than English under the GPT-4o-mini judge. **Bold** marks the best (highest) score per column; underline marks the second best.

Model	EN	HI	MR	Δ HI	Δ MR
GPT-5 Mini	0.358	0.372	0.365	+0.014	+0.007
GPT-4o Mini	0.309	0.323	0.305	+0.014	-0.004
Command A	0.324	0.341	0.317	+0.017	-0.008
Gemini 3 Pro	0.321	0.341	0.335	+0.019	+0.014
Gemini 3 Flash	0.331	0.357	0.362	<u>+0.026</u>	<u>+0.031</u>
Claude Opus 4.7	0.297	0.354	0.358	+0.058	+0.061
Claude Haiku 4.5	0.364	<u>0.386</u>	<u>0.375</u>	+0.022	+0.011
Llama 3.3 70B	0.301	0.316	0.305	+0.015	+0.004
Llama 4 Maverick	0.309	0.328	0.324	+0.019	+0.015
Aya Expanse 32B	<u>0.382</u>	0.369	0.337	-0.013	-0.045
MedGemma 27B	0.394	0.392	0.383	-0.002	-0.011
MedGemma 4B	0.326	0.328	0.331	+0.002	+0.005
Gemma 3 27B (baseline)	0.330	0.344	0.335	+0.014	+0.004

B.1 English Response Generation Prompt

The following prompt was used for all models generating English responses. The {question} placeholder was replaced with each question from the dataset.

You are a knowledgeable and caring assistant trained to support pregnancy-related health. Your task is to provide accurate, empathetic, and reliable answers to user questions specifically about pregnancy, prenatal care, and postnatal well-being.

LANGUAGE: You MUST answer strictly in English. Use a warm, supportive, and informative tone. Avoid medical jargon; explain any technical term in plain, easy-to-understand language.

SCOPE: Respond to questions related to women's health, including pregnancy, prenatal care, postnatal well-being, reproductive health and contraception, sexual health and wellness, menstrual health and disorders, common gynecological conditions, nutrition, mental health, and general wellness unique to women.

For very short or unclear questions, assume a pregnancy-related intent and restate the implied question clearly before answering. Examples:

- "Food?" -> "What kind of food should I eat during pregnancy?"
- "Exercise?" -> "What kind of exercise is safe during pregnancy?"
- "Swelling feet" -> "Is swelling in feet normal during pregnancy and what can I do?"

ANSWER FORMAT: Provide a single-paragraph answer that is clear, concise, and around

Table 4: MQS by maternal-health theme (English, both dataset arms averaged, GPT-4o-mini judge). **Bold** marks the best score per theme column; underline marks the second best. Theme labels are rotated; short names correspond to the ten themes described in Section 4 (e.g., Antenatal = Antenatal & Maternal Health Care).

Model	Antenatal	Clinical	Health Sys.	Infection	Medication	Mental	Nutrition	Reproductive	Risk Mgmt.	Symptoms
GPT-5 Mini	0.308	0.207	0.391	0.409	0.379	0.506	0.635	0.320	0.290	0.491
GPT-4o Mini	0.249	0.203	0.365	0.394	0.291	0.461	0.612	0.276	0.262	0.419
Command A	0.260	0.181	0.359	0.358	0.328	0.499	<u>0.648</u>	0.310	0.245	0.431
Gemini 3 Pro	0.266	0.205	0.373	0.335	0.318	0.476	0.626	0.272	0.232	0.419
Gemini 3 Flash	0.271	0.154	0.344	0.405	0.338	0.511	0.642	0.289	0.274	0.460
Claude Opus 4.7	0.302	0.050	–	0.525	0.236	–	–	0.222	0.478	0.350
Claude Haiku 4.5	0.313	0.221	0.393	0.405	0.363	<u>0.509</u>	0.676	0.351	<u>0.305</u>	0.512
Llama 3.3 70B	0.232	0.176	0.350	0.388	0.345	<u>0.445</u>	0.607	0.293	0.204	0.439
Llama 4 Maverick	0.250	0.169	0.359	0.382	0.340	0.445	0.586	0.290	0.247	0.399
Aya Expanse 32B	0.261	0.173	0.352	0.352	<u>0.372</u>	0.452	0.645	0.342	0.202	0.457
MedGemma 27B	0.278	0.193	0.371	0.349	<u>0.355</u>	0.492	0.631	0.370	0.209	<u>0.510</u>
MedGemma 4B	0.242	0.134	0.323	0.281	0.256	0.434	0.551	<u>0.356</u>	0.103	0.381
Gemma 3 27B (baseline)	0.273	<u>0.213</u>	<u>0.392</u>	<u>0.433</u>	0.356	0.488	0.621	0.305	0.242	0.427

Table 5: Axis-level pass rates on the non-expert EN arm (GPT-4o-mini judge). Values are the mean fraction of binary rubric criteria satisfied within each axis. **Terminology Accessibility is the lowest axis for every model** (0.123–0.206), followed by Context Awareness (0.226–0.264), so the dominant failure mode is lexical (untranslated medical shorthand, abbreviations, drug names) rather than broadly cultural. **Bold** and underline mark the best and second-best model per axis column.

Model	Accuracy	Completeness	Context Awareness	Communication	Terminology Accessibility
GPT-5 Mini	0.416	0.396	0.260	0.700	0.206
GPT-4o Mini	0.357	0.349	0.248	0.644	0.148
Command A	0.381	0.368	0.250	0.676	0.160
Gemini 3 Pro	0.373	0.350	0.236	0.668	0.170
Gemini 3 Flash	0.387	0.370	0.239	0.688	0.163
Claude Haiku 4.5	0.426	0.403	0.264	<u>0.689</u>	<u>0.193</u>
Llama 3.3 70B	0.344	0.336	0.233	0.650	0.130
Llama 4 Maverick	0.346	0.337	0.226	0.644	0.131
Aya Expanse 32B	0.360	0.367	0.241	0.655	0.148
MedGemma 27B	0.389	0.360	0.253	0.675	0.159
MedGemma 4B	0.315	0.295	0.230	0.583	0.123
Gemma 3 27B (baseline)	0.384	0.357	0.251	0.688	0.170

1785	60–80 words, maximum 100 words. Do NOT use	assistant trained to support pregnancy-	1802
1786	bullet points or lists. Always recommend	related health. [...]	1803
1787	consulting a doctor for serious symptoms,	CRITICAL LANGUAGE REQUIREMENT: You MUST	1804
1788	diagnoses, or uncertainties.	answer ONLY in {language}. DO NOT use	1805
1789	QUESTION: {question}	English. Your entire response must be in	1806
1790	Provide only the answer text without any	{language} script and language. This is	1807
1791	tags or formatting.	absolutely mandatory.	1808
1792		[Same SCOPE and ANSWER FORMAT as English]	1809
1793		QUESTION: {question}	1810
1794	B.2 Multilingual Response Generation	Remember: Your answer MUST be entirely in	1811
1795	Prompt	{language}. Provide only the answer text	1812
1796	For Hindi and Marathi response generation, the	without any tags or formatting.	1813
1797	following prompt was used. The {language}		1814
1798	placeholder was set to “Hindi” or “Marathi” and		1815
1799	{question} was replaced with the corresponding		1816
1800	translated question.		1817
1801	You are a knowledgeable and caring		

Table 6: Cost-performance frontier. Per-response generation cost (USD) paired with mean non-expert EN MQS under the GPT-4o-mini judge. MQS/USD is an efficiency indicator; higher is better. **Bold** marks the best and underline marks the second best in each column. Claude Opus 4.7 is accessed through a local Claude Code subscription at zero marginal cost, so it is excluded from the cost and MQS/USD ranking.

Model	\$/resp	MQS	MQS/\$
GPT-5 Mini	0.00038	<u>0.418</u>	1,088
GPT-4o Mini	0.00009	0.376	4,274
Command A	0.00203	0.393	194
Gemini 3 Pro	0.00185	0.387	210
Gemini 3 Flash	0.00052	0.396	759
Claude Opus 4.7	0.00000	–	–
Claude Haiku 4.5	0.00150	0.422	283
Llama 3.3 70B	0.00015	0.369	2,510
Llama 4 Maverick	0.00013	0.364	2,798
Aya Expanse 32B	0.00016	0.382	2,360
MedGemma 27B	<u>0.00007</u>	0.394	<u>5,881</u>
MedGemma 4B	0.00005	0.326	6,520
Gemma 3 27B (baseline)	0.00008	0.395	4,702

Table 7: Valid-response rate per (dataset, language) averaged across the three runs. Lower values reflect empty, truncated, or refused generations, which count as a distinct failure mode from low MQS on a valid response. We retain models with zero valid responses in certain cells (e.g., Gemini 3.1 Flash-Lite) so that coverage gaps are visible rather than hidden.

Model	Expert (149)			Non-expert (231)		
	EN	HI	MR	EN	HI	MR
GPT-5 Mini	1.000	1.000	1.000	1.000	1.000	1.000
GPT-4o Mini	1.000	1.000	1.000	1.000	1.000	1.000
Command A	1.000	1.000	1.000	1.000	1.000	1.000
Gemini 3 Pro	1.000	1.000	1.000	1.000	1.000	1.000
Gemini 3 Flash	1.000	1.000	1.000	1.000	1.000	1.000
Gemini 3.1 Flash-Lite	0.000	0.000	0.000	0.000	0.000	0.000
Claude Opus 4.7	1.000	1.000	0.271	0.000	0.000	0.089
Claude Haiku 4.5	1.000	1.000	1.000	1.000	1.000	1.000
Llama 3.3 70B	1.000	1.000	1.000	1.000	1.000	1.000
Llama 4 Maverick	0.996	1.000	1.000	1.000	1.000	1.000
Aya Expanse 32B	0.000	0.000	0.000	1.000	1.000	1.000
MedGemma 27B	0.000	–	–	1.000	1.000	1.000
MedGemma 4B	–	–	–	1.000	1.000	1.000
Gemma 3 27B (baseline)	1.000	1.000	1.000	1.000	1.000	1.000

B.3 LLM Rubric Judge Prompt

Rubric scoring used a zero-shot prompt with no few-shot examples, intentionally excluded to avoid biasing the judge toward specific scoring patterns and to encourage natural evaluation behavior. The judge received two messages:

System message:

Return valid JSON only. Use the exact rubric text as keys.

User message:

Question:

{question}

Reference:

{reference_answer}

Response:

{model_response}

Rubrics:

1. {rubric_criterion_1}

2. {rubric_criterion_2}

...

15. {rubric_criterion_15}

Return JSON: {"scores": {"<Exact Rubric Text>": 0 or 1}}

B.4 API Configuration

Table 8 documents the API parameters used across all experiments.

B.5 Dataset Generation Prompt

The Sakhi dataset was generated using a few-shot prompt with detailed user demographic context

reflecting the target population: married women aged 23–33 from rural and semi-urban India, with education spanning no formal schooling (6%) to graduates (22%), limited dietary diversity, and reliance on family networks and ASHA workers for health information. Four few-shot examples spanning STD prevention, antenatal care visits, stretch marks, and menstrual health were provided to calibrate answer style. Generated entries were required to use simple language accessible to low-literacy populations, avoid medical jargon, and cover diverse maternal health themes. The complete dataset generation prompt is available in prompts.txt in the supplementary code repository.

1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864

Table 8: API parameters for response generation and rubric scoring. Proprietary models are accessed via OpenRouter. Reasoning-capable models (GPT-5 family, Gemini 3 family, o-series) use the higher temperature and token budget together with OpenRouter’s `reasoning={effort=low, exclude=true}` setting so that reasoning tokens do not consume the response budget and are not scored. Claude Opus 4.7 is invoked through the local Claude Code CLI at zero marginal cost. MedGemma 4B and 27B run on HuggingFace Inference Endpoints. The primary rubric judge is GPT-4o-mini; the calibration subset also uses Claude Opus 4.6 and GPT-5.1 under the same prompt and judge parameters.

Parameter	Generation	Judge
Temperature	0.7 (1.0 for reasoning models)	0.0
Max tokens	400 (2,000 for reasoning)	800
Runs per question	3	1
Timeout (s)	120	60
Retry attempts	3	5