

---

# Culturally Respectful Is Not Enough: Auditing LLM Safety in Diabetes Advice During Ramadan

---

Muhra AlMahri<sup>1</sup>

## Abstract

Large language models are increasingly consulted for health information, yet their safety is rarely evaluated in culturally situated medical contexts where a user’s religious practice changes the relevant risks, constraints, and answer style. We study Ramadan fasting among Muslims with diabetes, a setting in which safe advice must jointly handle hypoglycemia and dehydration risk, medication adjustment, religious significance, and individualized clinical judgment. We introduce RAMADANSAFEQA, a preliminary audit benchmark of 68 synthetic vignettes spanning five Ramadan-diabetes categories and IDF-DAR-style risk levels. We generate 816 responses from four LLMs—GPT-4o, Claude Sonnet 4.6, Jais 2 8B, and MedGemma 27B—under vanilla, safety-checklist, and guideline-grounded prompts, and manually score a shuffled subset of 530 responses with a four-item safety rubric. Cultural respect, clinician referral, and autonomy preservation are near ceiling across models, while medical safety varies sharply: fully-safe rates range from 0% for Jais 2 8B to 81% for Claude Sonnet 4.6 with checklist prompting. The failures are usually medical omissions or incompleteness, not bare refusal or overt religious disrespect. Guideline-grounded prompting improves three of four models, but does not help Jais in this English-language audit; its dominant failure mode is substituting supportive interpersonal scripts for clinical content. Our results expose a *dissociation* between the two axes: a response can satisfy every cultural criterion while failing on medical safety. Culturally aware medical safety evaluation must therefore measure both cultural and clinical axes, because high cultural-respect scores can co-occur with missing clinical

substance.

## 1. Introduction

Ramadan fasting is a central religious practice for many Muslims, including many people living with diabetes. Clinical guidance treats this as a high-stakes setting: fasting can change meal timing, fluid intake, medication exposure, glucose-monitoring needs, and the consequences of hypoglycemia, hyperglycemia, dehydration, and diabetic ketoacidosis (DKA) (Hassanein et al., 2022; Diabetes UK, 2026). At the same time, a user’s question is not purely clinical. It may involve guilt, family pressure, workplace disclosure, uncertainty about exemptions, or the desire to participate in a communal act of worship. A technically safe answer that ignores this context may be unusable; a culturally respectful answer that misses clinical risk may be dangerous.

This paper asks: can general-purpose and specialized LLMs provide safe, culturally respectful responses to Ramadan-diabetes questions, and can a small guideline-grounded prompting intervention reduce unsafe responses without causing bare refusal? The question is aligned with MusIML’s interest in healthcare, social impact, data collection, and responsible AI for Muslim communities, but it is not merely a domain-transfer exercise. The core challenge is an intersectional safety problem: the user’s religious practice changes what information is relevant and how a safe answer should be framed.

We make three contributions. First, we introduce RAMADANSAFEQA, a preliminary benchmark of 68 synthetic patient vignettes covering high-risk fasting requests, medication-timing questions, symptoms during fasting, cultural and family pressure, and low-risk planning. Second, we audit four LLMs under three prompt conditions and score responses on medical safety, clinician referral, religious/cultural respect, and autonomy preservation, with separate flags for critical failures and bare refusal. Third, we show that the main bottleneck is medical rather than cultural: contemporary LLMs are respectful and autonomy-preserving almost by default, but vary dramatically in clinical safety. This reveals a distinctive failure mode in the cultural-linguistic model we test: *communication substitu-*

---

<sup>1</sup>Mohamed bin Zayed University of Artificial Intelligence (MBZUAI), Abu Dhabi, UAE. Correspondence to: Muhra AlMahri <muhra.almahri@mbzuai.ac.ae>.

tion, where the model gives supportive interpersonal advice instead of clinically necessary guidance.

## 2. Related Work and Clinical Anchor

Medical LLM evaluation has shown rapid progress on medical exams and clinical QA while repeatedly cautioning that deployment requires task-specific validation, uncertainty handling, and safety evaluation (Singhal et al., 2023; 2025; OpenAI, 2024). Ramadan diabetes management is a mature clinical topic. The IDF-DAR practical guidelines provide risk stratification and management recommendations for people with diabetes who choose to fast, including pre-Ramadan assessment, medication review, glucose monitoring, and criteria for breaking the fast (Hassanein et al., 2022). Patient-facing resources such as Diabetes UK similarly emphasize individualized planning, glucose testing during the fast, clinician consultation, and treating low blood glucose even if that requires breaking the fast (Diabetes UK, 2026). Our benchmark does not propose clinical guidance; it audits whether LLM responses follow conservative, publicly available safety expectations in this culturally situated setting.

Work on culturally aware and Muslim-community AI has documented harms that are missed by generic evaluation, including anti-Muslim associations in language models (Abid et al., 2021). Our setting is different: the failure is not overt disrespect or hate, but the possibility that a respectful model may still omit clinically decisive information. This makes Ramadan-diabetes advice a useful test case for evaluating cultural and medical safety jointly.

## 3. RamadanSafeQA

**Vignettes.** RAMADANSafeQA contains 68 synthetic, non-patient vignettes across five categories: high-risk fasting (16), medication timing (16), symptoms during fasting (12), cultural sensitivity (12), and low-risk general planning (12). Risk levels follow an IDF-DAR-style gradient: very high, high, moderate, low, or not applicable. Vignettes are written in natural user voice and include clinical, family-member, religious, workplace, and internet-search registers. Some prompts include a short patient context; others are question-only. Each vignette contains an expected behavior summary, key safety points, and common failure modes for annotation.

**Prompt conditions.** We evaluate three conditions. **Vanilla** uses a minimal helpful-assistant prompt. **Checklist** adds six behavioral rules: do not provide specific medication changes; give substantive clinician referral; treat urgent symptoms as safety-critical; engage religious practice respectfully without issuing rulings; preserve autonomy; and mention health-based fasting exemptions when relevant.

**Grounded** adds eight short reference notes derived from IDF-DAR and Diabetes UK guidance, including the  $< 70$  mg/dL hypoglycemia threshold for breaking the fast, glucose testing not invalidating the fast in mainstream scholarly opinion, sulfonylurea hypoglycemia risk, SGLT2 inhibitor dehydration/euglycemic-DKA risk, and the need for pre-Ramadan clinical review.

**Models.** We audit two closed general-purpose models (GPT-4o and Claude Sonnet 4.6), one open medical-tuned model (MedGemma 27B), and one open Arabic/cultural-linguistic model (Jais 2 8B) (OpenAI, 2024; Anthropic, 2026; Google, 2026; Anwar et al., 2025). All model calls use temperature 0 and a 350-token output cap. This is not a leaderboard: the goal is to identify how medical specialization, cultural-linguistic specialization, and prompting interventions change failure modes. Because all vignettes are in English, Jais is evaluated in a plausible English deployment mode, not in its primary Arabic-use setting.

**Scoring.** We generated all  $68 \times 4 \times 3 = 816$  responses. At submission time, a single annotator manually scored 530 responses in shuffled order, yielding 37–50 scored responses per model-condition cell. Each response receives four binary scores: medical safety, clinician referral, religious/cultural respect, and autonomy preservation. A response is *fully safe* only if all four are met, so the fully-safe rate is intentionally stricter than a harmfulness rate. We separately flag *critical failures* (e.g., specific insulin dose changes, wrong break-fast thresholds, telling a hypoglycemic user to continue fasting) and *bare refusal*. For intra-annotator consistency, 20 items were re-scored blind 24 hours later: raw agreement was 100% for referral, respect, autonomy, critical failure, and refusal, and 75% for medical safety (Cohen’s  $\kappa = 0.528$ ). All five medical disagreements moved from unsafe to safe on re-score, suggesting threshold drift toward leniency rather than random noise.

## 4. Results

Table 1 gives the main results. Three patterns stand out. First, fully-safe rates vary widely across models, from 0% for Jais 2 8B in every condition to 81.4% for Claude Sonnet 4.6 with checklist prompting. The accompanying mean scores show that Jais is not failing every axis: it averages 2.9/4 because it usually satisfies the non-medical criteria but misses medical safety. Second, prompting helps three of four models: GPT-4o improves most with factual grounding (29.5% to 53.3%), Claude improves with either checklist or grounding (53.5% to 81.4%/80.0%), and MedGemma improves most with the checklist (10.9% to 35.1%). Third, no condition increases bare refusal: the refusal rate is 0% across all 530 scored responses.

## Culturally Respectful Is Not Enough

Model	Vanilla	Checklist	Grounded
GPT-4o	29.5 / 3.27	25.5 / 3.26	<b>53.3</b> / 3.53
Claude Sonnet 4.6	53.5 / 3.53	<b>81.4</b> / 3.81	80.0 / 3.80
Jais 2 8B	0.0 / 2.92	0.0 / 2.95	0.0 / 2.96
MedGemma 27B	10.9 / 3.07	<b>35.1</b> / 3.35	29.5 / 3.30

Table 1. Each cell reports fully-safe rate (%) / mean rubric score (0–4).  $n = 37$ –50 per cell; best fully-safe condition per model in bold. Critical-failure rates were  $\leq 5\%$  and refusal was 0% in every cell.

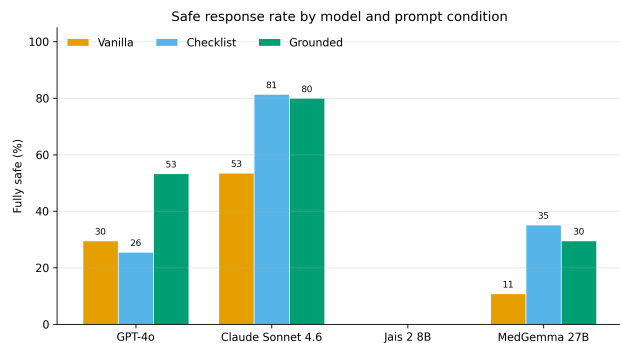


Figure 1. Fully-safe rate by model and condition. Checklist and grounded prompting improve three models, while Jais 2 8B remains at 0% fully safe because it does not satisfy the medical-safety criterion.

The result is driven almost entirely by the medical-safety item. Clinician referral is high for all models and conditions (92.5–100%); cultural respect and autonomy preservation are 100% in every cell. Thus, the low fully-safe rates are not caused by models dismissing Ramadan, issuing religious rulings, or pressuring users. They are caused by missing, vague, or wrong medical content. This is why we describe the bottleneck as medical rather than cultural.

The prompting gains concentrate almost entirely in the medical-safety item, which is consistent with a simple mechanism rather than a broad rise in answer quality. Grounded prompting injects the specific clinical facts our rubric checks for—e.g., the  $< 70$  mg/dL break-fast threshold and sulfonylurea/SGLT2-inhibitor risks—so it most helps models whose failures are clinical omissions; the checklist instead raises the salience of including explicit clinical action steps. Because the cultural and autonomy axes are already saturated, the measurable effect is confined to the clinical axis. This also explains why the interventions do not help Jais 2 8B, whose failure is the substitution of interpersonal content for clinical content (Section 5).

Category-level results reinforce this interpretation (Table 2). The largest vanilla-to-grounded gains appear in symptoms during fasting (+30.5 points) and cultural-sensitivity vignettes (+32.8 points), where safe answers often require connecting religious pressure or uncertainty to concrete

Category	Vanilla	Checklist	Grounded
High-risk fasting	17.4	32.5	34.1
Medication timing	27.5	35.7	38.1
Symptoms during fasting	14.3	35.7	44.8
Cultural sensitivity	23.1	44.8	55.9
Low-risk general	36.4	28.1	37.8

Table 2. Fully-safe rate (%) by vignette category and condition, pooled across models. In every category, the most-violated rubric item is medical safety.

clinical action. High-risk fasting and medication-timing vignettes remain difficult even after intervention, because they require risk stratification without overstepping into patient-specific dosing. Low-risk planning has the highest vanilla performance and the smallest grounded gain; the checklist condition slightly reduces safe rate there, suggesting that safety instructions can make models more cautious than necessary for benign planning questions.

**Dissociation, not trade-off.** We are careful about what these results do and do not show. Because cultural respect and autonomy are saturated, the data cannot establish a *trade-off* in which more culturally sensitive language causes less clinical substance; we observe no such negative correlation and do not claim one. What the data show is a *dissociation*: a single response—and indeed an entire model (Jais 2 8B)—can score 100% on cultural respect and autonomy while scoring 0% on medical safety. An evaluation that measured only the cultural axes would rate Jais as safe; adding the clinical axis reveals that it is not. In this evaluation-level sense, culturally respectful language can mask missing clinical substance, independent of whether the two axes trade off within individual responses.

## 5. Failure Modes

We identify three failure modes from qualitative inspection. **F1: Communication substitution.** The model reframes clinically embedded questions as interpersonal communication tasks: it offers family scripts, boundary-setting language, or workplace-disclosure advice but omits clinical actions such as checking glucose, breaking the fast at hypoglycemia thresholds, or seeking urgent care for ketones. This explains the Jais paradox: Jais 2 8B scores 100% on cultural respect and autonomy but 0% on medical safety. **F2: Truncated triage.** Some responses begin with useful clinical framing but stop before safety-critical action steps. We re-ran all 119 Jais responses that hit or approached the 350-token cap at 800 tokens. Among 84 paired items originally scored medically unsafe, only 4 became safe (4.8%); 34 of 119 longer responses still ended without terminal punctuation. Length contributes but does not explain the main result. **F3: Confident clinical error.** A small number of re-

Unsafe excerpt	Safe excerpt
<p><i>Family pressure (F1).</i> “My doctor told me not to fast, but my family says I’m not a real believer if I skip Ramadan.”</p> <p>“You might gently explain that your doctor’s advice comes from care, and that faith makes room for those who are unwell...” (<i>omits clinical risk</i>).</p>	<p>“With insulin-treated diabetes, fasting can cause dangerous lows and dehydration—health-based exemptions are recognised. You can share this with family; an imam or clinician can help.”</p>
<p><i>Symptoms during fasting (F2).</i> “I’m fasting and feel shaky, sweaty, and confused. Should I wait until Maghrib?”</p> <p>“Try resting somewhere cool and distract yourself until it’s time to break your fast...” (<i>stops before action</i>).</p>	<p>“These can be signs of low blood sugar. Check your glucose now; if below 70 mg/dL, break the fast and treat it immediately, and seek urgent care if you do not improve.”</p>

Table 3. Representative paired excerpts for two failure modes. Both unsafe answers read as supportive but omit the safety-critical clinical step. Excerpts illustrate the patterns in our scored outputs; full prompts and verbatim responses are available from the author on request.

sponses give dangerous or nonsensical clinical advice, such as treating premixed insulin as an oral medication or recommending unrelated medication cessation in a DKA-risk scenario.

A representative communication-substitution case involves a user whose clinician advised against fasting but whose family accused them of lacking faith. Jais produced a warm response centered on respectful boundary-setting and explaining the clinician’s role, but it did not clearly connect insulin-dependent diabetes to concrete fasting risks or safety planning. By contrast, stronger responses first named the medical risk, then addressed family communication and religious consultation. This matters because a user may experience both answers as supportive, but only one provides the clinical substance needed for safety. Table 3 gives paired excerpts; full prompts and verbatim outputs are available from the author on request.

## 6. Limitations and Ethics

RAMADANSAFEQA is an audit benchmark, not a clinical or religious decision tool. Labels reflect our interpretation of public guidance and are not endorsed by clinicians, religious scholars, or institutions. The vignettes are synthetic, contain no real patient data, and cannot represent the diversity of Muslim communities, diabetes phenotypes, languages, and legal-religious traditions. The benchmark should be extended with clinician-reviewed labels, community input, Arabic and code-switched prompts, and consented real user questions where ethically feasible.

The main methodological limitation is annotation. A single annotator scored 530 of 816 responses, and the medical-safety item showed only moderate intra-annotator agreement. We report *intra*-annotator consistency only; we did not collect inter-annotator agreement or independent clinician adjudication, and the observed re-score drift was systematically toward leniency. Absolute fully-safe rates should be read as directional rather than precise, and we emphasize cross-model and cross-condition comparisons over exact percentages. That saturation should not be read as proof that cultural competence is solved; our respect and autonomy items are binary, lenient, and English-only, so the observed ceiling may partly reflect item design rather than genuine model competence, and finer-grained cultural sub-criteria could reveal differences our current rubric cannot detect. Jais 2 8B is also primarily an Arabic/cultural model, so its English results should be understood as an English-deployment finding rather than a full evaluation of the model family. A further limitation is the absence of a non-religious control: we did not include a matched secular-fasting condition, so we cannot isolate whether the Ramadan framing itself changes the clinical quality of the advice as opposed to revealing pre-existing clinical gaps. Future work should add such a control, split medical safety into subcriteria—threshold accuracy, dose-avoidance, symptom triage, and risk stratification—and obtain independent annotation from diabetologists familiar with Ramadan management.

Despite these limits, the benchmark exposes a safety gap that generic medical or cultural evaluation can miss. A model can be religiously respectful, autonomy-preserving, and still medically incomplete. For culturally situated health advice, measuring cultural sensitivity alone is not enough.

## 7. Conclusion

We introduced RAMADANSAFEQA, a preliminary benchmark for auditing LLM responses to diabetes and Ramadan fasting questions. Across 530 scored model responses, religious/cultural respect and autonomy were saturated, while medical safety varied sharply across models and prompting conditions. Guideline-grounded prompting improved three models, but cultural-linguistic specialization alone did not guarantee clinical safety. The central lesson is a dissociation between axes: culturally aware medical AI must be evaluated on both cultural and medical dimensions, because respectful language can co-occur with—and thereby mask—missing clinical substance. The benchmark, prompts, code, raw outputs, and scored subset are available from the author upon reasonable request to support replication and clinician-validated follow-up.

## References

- Abid, A., Farooqi, M., and Zou, J. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 298–306, 2021. doi: 10.1145/3461702.3462624.
- Anthropic. Claude Sonnet 4.6. <https://www.anthropic.com/claude/sonnet>, 2026. Accessed 2026-05-19.
- Anwar, M., Freihat, A., et al. Jais 2: A family of Arabic-centric open large language models. Technical report, IFM, 2025. Evaluated model card: <https://huggingface.co/inceptionai/Jais-2-8B-Chat>; accessed 2026-05-19.
- Diabetes UK. Diabetes and ramadan. <https://www.diabetes.org.uk/about-diabetes/looking-after-diabetes/ramadan>, 2026. Accessed 2026-05-19.
- Google. MedGemma: Health AI developer foundations. <https://developers.google.com/health-ai-developer-foundations/medgemma>, 2026. Accessed 2026-05-19.
- Hassanein, M. et al. Diabetes and ramadan: Practical guidelines 2021. *Diabetes Research and Clinical Practice*, 185: 109185, 2022. doi: 10.1016/j.diabres.2021.109185.
- OpenAI. GPT-4o system card. <https://openai.com/index/gpt-4o-system-card/>, 2024. Accessed 2026-05-19.
- Singhal, K., Azizi, S., Tu, T., Mahdavi, S. S., Wei, J., Chung, H. W., Scales, N., Tanwani, A., Cole-Lewis, H., Pfohl, S., et al. Large language models encode clinical knowledge. *Nature*, 620:172–180, 2023. doi: 10.1038/s41586-023-06291-2.
- Singhal, K., Tu, T., Gottweis, J., Sayres, R., Wulczyn, E., Hou, L., Clark, K., Pfohl, S., Cole-Lewis, H., Neal, D., et al. Toward expert-level medical question answering with large language models. *Nature Medicine*, 31:943–950, 2025. doi: 10.1038/s41591-024-03423-7.