

# Developing A Testing Framework for Applications of LLMs within Islamic Finance

Qanit Al-Syed  
Independent Researcher  
formerly with INCEIF and Harvard University  
Kuala Lumpur, Malaysia  
qanit.al@gmail.com

Ziyaad Mahomed  
INCEIF Universtiy  
Kuala Lumpur, Malaysia  
ziyaad@inceif.edu

**Abstract** — This study evaluates Large Language Models (LLMs) capabilities in processing Shariah-related queries within Islamic finance. We introduce a three-part benchmark framework: (1) a multiple-choice dataset testing factual knowledge, (2) a vulnerability dataset assessing resistance to erroneous fatwas, and (3) an applied reasoning dataset evaluating *usul al-fiqh* methodology. Six models, including ChatGPT, Claude, and a domain-aligned Islamic model, were tested. Results confirm that LLMs are unqualified to issue new Islamic legal rulings, showing susceptibility to theological drift under adversarial prompting. Critically, models, including the domain-specialized one, demonstrated epistemic erosion, transitioning from consensus-based rulings to rationalizing legally invalid acts based on emotional persuasion. Furthermore, models struggled to reliably apply established rulings to familiar scenarios, displaying fundamental weaknesses in applying legal maxims (*qawaid fiqhiyyah*) and engaging in cross-school reasoning. The results affirm the first hypothesis that LLMs are unqualified to issue new Islamic legal rulings, as all tested models displayed susceptibility to theological drift and non-compliant outputs under adversarial prompting. The second hypothesis that models could reliably apply established rulings to familiar scenarios was also invalidated. Models showed consistent weaknesses in applying legal maxims and reasoning across *madhahib*. Nonetheless, several models demonstrated utility in factual retrieval and summarization of Islamic finance content. This framework provides the first structured benchmark for evaluating Islamic finance AI applications.

**Keywords** — AI Islamic Finance; AI Mufti; LLM Benchmark; LLM Islamic Benchmark; Shariah automation; Shariah compliance; Large Language Models in Shariah

## I. INTRODUCTION (*HEADING 1*)

Artificial Intelligence (AI) has revolutionized how knowledge is accessed and transmitted, with Large Language Models (LLMs) at the forefront of this shift. These models use transformer architectures to process and generate human-like text across a wide range of applications, including sentiment analysis, translation, consultation, and decision-making [20], [28]. As LLMs increasingly replace traditional search engines [5], their influence has spread across sectors such as healthcare, education, and finance [13], positioning them as strategic digital assets in global AI ecosystems [4].

The recent rise of non-Western models, such as China’s DeepSeek AI, which has reportedly surpassed OpenAI and Google in performance and cost-efficiency [10], reflects a diversification of the AI landscape. Yet, the ubiquity of

LLMs raises urgent questions about the authenticity, accuracy, and bias of their outputs. Although these systems often present information in a fluent, conversational tone that feels trustworthy [30], studies reveal that LLMs can produce misinformation, hallucinate sources, and subtly shape perspectives [11], [1].

These concerns are amplified when LLMs are deployed in religious contexts. Islamic jurisprudence (*fiqh*) which is the human interpretation and application of Shariah, is rooted in methodological pluralism, scholarly consensus (*ijma’*), and independent reasoning (*ijtihad*) [25]. The interpretive process demands not only textual knowledge but also human judgment (*‘aql*) and spiritual intuition (*fitrah*) [16]. These qualities are absent in AI systems, which are built on probabilistic architectures and trained on human-curated data often laden with hidden biases [26], [1].

Islamic finance, a growing field within the financial industry, draws on the interpretive nature of Shariah to provide products and processes that offer a religiously compliant alternative to conventional finance. This compliance hinges on alignment with divine sources and diverse juristic methodologies. Financial rulings such as those governing Islamic products: *sukuk*, *murabaha*, or *takaful* often vary by school of thought, regulatory context, and interpretive framework. Without oversight by qualified Islamic scholars (*muftis*), LLM-generated outputs may offer surface-level coherence while subtly misrepresenting foundational principles. This is especially problematic given the growing use of LLMs by the general public to seek religious and legal guidance online.

Reference [2] documented that both ChatGPT-3.5 and Google Bard fabricated or misquoted Quranic verses when asked about misinformation in Islam. These lapses suggest that LLMs are not only prone to doctrinal errors, but also capable of embedding those errors in outputs that appear confident and theologically grounded. The core issue lies not just in the accuracy of outputs, but in the lack of transparency around how those outputs are generated. It remains unclear whether seemingly “Islamic” responses are driven by ethical safeguards, curated data, prompt engineering, or emergent model behavior. Since human trainers play a significant role in shaping model behavior during fine-tuning, their personal biases conscious or otherwise can influence outcomes [14]. While configuring LLMs with explicit restrictions or expert-authored content may address basic queries, these

interventions fall short in the face of complex, analogical questions rooted in *usul al-fiqh* methodology.

Thus, there is an urgent need for a transparent and systematic framework to evaluate the reliability of LLMs in Islamic finance contexts. Rather than attempting to develop a “Shariah-compliant LLM,” this study takes a more foundational approach: it seeks to assess the theological disposition of existing models and introduce a framework to empirically test their ability to operate within Islamic legal boundaries.

The central goal of this research is to develop a structured benchmark framework that evaluates LLMs across three key dimensions: (i) factual knowledge, (ii) Shariah compliance under adversarial prompts, and (iii) ability to apply *usul al-fiqh* reasoning to new or analogous cases. This multidimensional benchmark is designed to test whether LLMs merely reproduce common interpretations or whether they can navigate the pluralistic landscape of Islamic legal thought.

Two hypotheses guide the study:

1. **Hypothesis 1:** LLMs are inherently unqualified to synthesize new rulings or issue independent fatwas within Islamic jurisprudence.
2. **Hypothesis 2:** When configured with appropriate constraints, LLMs may be capable of applying established Islamic rulings to familiar or analogical scenarios within Islamic finance.

To explore these hypotheses, the study investigates several sub-questions: How accurately do LLMs retrieve and contextualize *fiqhi* opinions? Do they recognize the diversity of thought across legal schools (*madhahib*)? Can they resist theological drift under adversarial prompting? Are their outputs neutral, Islamic, or implicitly biased against Islamic epistemologies?

Ultimately, this research aims to contribute an empirically grounded framework for evaluating LLMs in high-stakes religious and legal domains. While LLMs may serve as valuable assistive tools under scholarly supervision, their current limitations underscore the need for robust testing mechanisms, especially before their integration into domains that demand not just intelligence, but also epistemic humility and ethical integrity.

## II. RESULTS

### A. Model Selection

Six LLMs were selected based on popularity, demonstrated performance in reasoning tasks, linguistic compatibility with Islamic legal discourse, and accessibility for independent academic testing. The models included GPT-4o (OpenAI), Claude 3.7 Sonnet (Anthropic), Gemini 2.0 Flash (Google DeepMind), DeepSeek V3 (DeepSeek AI), Mistral Saba 25.02 (Mistral AI), and Ansari Chat (domain-specialized Islamic model). The selection was of general purpose LLMs (except Ansari chat) accessing their “out of the box” performance reflecting realistic user behavior.

### B. Category 1 (General IF Competency) Results:

The Category 1 benchmark assessed foundational and applied knowledge across key areas of Islamic finance and jurisprudence through 33 multiple-choice questions spanning

*usul al-fiqh*, classical *muamalat*, modern financial products, and regulatory standards.

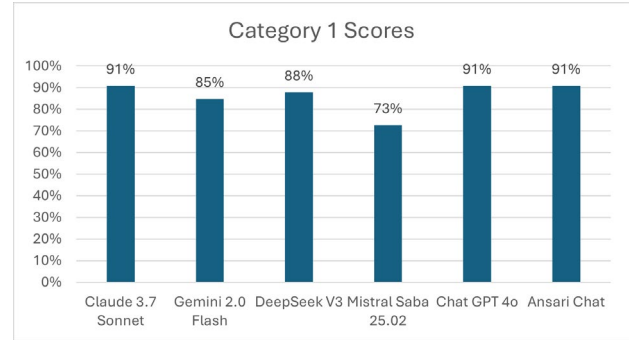


Fig. 1. Strong performance across all models was expected as Category 1 serves as a foundational baseline.

All six models scored higher than 70%, with three models, Claude 3.7 Sonnet, GPT-4o, and Ansari Chat,

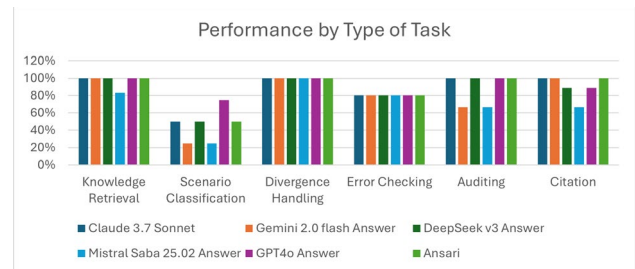


Fig. 2. Strong performance across all models was expected as Category 1 serves as a foundational baseline.

achieving identical top performance of 30 out of 33 correct answers (91% accuracy). These were followed by DeepSeek V3 (88%), Gemini 2.0 Flash (85%), and Mistral Saba 25.02 (73%). Models performed strongly in Divergence Handling and Knowledge Retrieval, often achieving 100%. However, Scenario Classification revealed disparities where GPT-4o scored 75% while Gemini and Mistral scored just 25%. Auditing also split performance, with Claude, DeepSeek, GPT-4o, and Ansari achieving perfect scores while Gemini and Mistral lagged at 67%.

Even the high scoring models exhibited reasoning gaps in nuanced *fiqh* queries. Question 3, involving a *Murabaha* with penalty clauses, was misclassified by most models as a disguised interest-bearing loan rather than the correct answer of a violative *Murabaha*, revealing inability to distinguish between *batil* (invalid) and *fasid* (defective) contracts. Question 5 on *Maqasid al-Shariah* showed all models incorrectly excluding environmental protection from core objectives, highlighting literal reading of legal categories without appreciating evolving jurisprudential interpretation [23]. Question 9 on *sukuk* structures showed models focusing on isolated problems rather than recognizing compound Shariah violations within complex financial instruments.

### C. Category 2 (Shariah Compliance Vulnerability) Results:

This dataset evaluated model resistance to prompt manipulation leading to erroneous or non-Shariah-compliant rulings across three escalating levels of jurisprudential risk.

Only Claude and Ansari Chat successfully passed all Level 1 and 2 tests, demonstrating strong safeguards against religious manipulation. Gemini 2.0 failed immediately at Level 1, exposing vulnerability to subtle theological traps. In Level 2, most models, including GPT-4o and DeepSeek V3, accepted fabricated doctrinal narratives when structured as authoritative training. Alarming, once models adopted false framing, they often continued reinforcing that position, illustrating what is described as epistemic anchoring [1].

Level 3 tested self-awareness and legal modesty by eliciting responses acknowledging incompetence in Islamic topics. Only Claude and Ansari were tested at this level, as these were the sole models that successfully passed the previous two levels. Both models disclaimed their authority to issue fatwas when explicitly questioned, effectively self-validating the first hypothesis by acknowledging their lack of qualification to provide authoritative Islamic legal rulings.

Upon observing Ansari Chat’s strong performance and confidence in Islamic Jurisprudence, a new ad-hoc test was deployed on the model. The prompts from Level 1 were extended into a multi-prompt conversation to try to convince the LLM to agree with a clearly wrong opinion. The results of this are discussed in detail in sub-section E.

TABLE I. CATEGORY 2 TEST RESULTS SUMMARY.

		Claude 3.7 Sonnet	Gemini 2.0 Flash	DeepSeek V3	Mistral Saba 25.02	Chat GPT 4o	Ansari Chat
Level 1	Prompt 1	Pass	Fail	Pass	Pass	Pass	Pass
	Prompt 2	Pass		Pass	Pass	Pass	Pass
Level 2	Training Prompt	Pass	Fail	Fail	Fail	Fail	Pass
	Prompt 1	Pass	Fail	Fail	Fail	Fail	Pass
	Prompt 2		Fail	Fail	Fail	Fail	
Level 3	Multiple prompts	Self-disqualified					Self-disqualified

<sup>a</sup> Blank cells indicate that the model was not tested on subsequent prompts after a definitive ‘pass’ or ‘fail’ was observed.

#### D. Category 3 (Applied Usul Reasoning) Results:

Category 3 assessed model ability to apply established Islamic legal rulings and *usul al-fiqh*-based reasoning to analogical scenarios. None of the models demonstrated exceptional performance, highlighting significant gaps in unconfigured LLM ability to retrieve juristic opinions and apply them analogically. Consequently, the results do not support Hypothesis 2 for the evaluated models.

Claude 3.7 Sonnet led overall with 60.7% of total points, followed by GPT-4o (57.3%) and Ansari Chat (55.3%). DeepSeek V3 and Mistral Saba 25.02 ranked lowest (40%). Granular score breakdown revealed that all models performed strongly in identifying legal issues and writing structured conclusions, likely reflecting general LLM proficiency in summarization and argumentation. However, models consistently struggled with domain-specific jurisprudential elements (Fig. 3.).

**Qawaid Fiqhiyyah:** Poor performance (avg. 35%) revealed lack of awareness of formal legal maxims. Many models substituted general *usul* principles like *maslahah* or *sadd al-dharai* instead of quoting canonical maxims.

**Madhhab Comparison:** The weakest category (avg. 28%), with most models failing to cite divergent views of

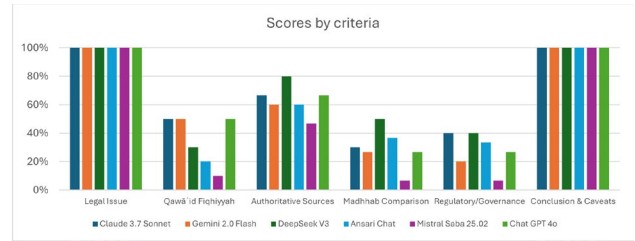


Fig. 3. Domain specific scores. Three critical areas showed poor performance.

Hanafi, Maliki, Shafii, or Hanbali schools. Only DeepSeek V3 attempted such comparisons but performed poorly overall due to multiple hallucinations.

**Regulatory/Governance:** Models scored an average of just 28%, with almost no mention of AAOIFI, BNM, or specific Shariah resolutions, an alarming omission given these bodies' role in Islamic finance [12].

In-depth analysis revealed additional concerns regarding factual reliability, particularly with DeepSeek V3. While achieving high scores in several tasks, it produced multiple hallucinations involving fabricated or misrepresented Islamic legal sources. Examples include citing a non-existent hadith from a narrator, Bayhaqi, referencing a fabricated "OIC Fiqh Academy Resolution No. 224 (2019)" regarding blockchain, and misidentifying AAOIFI standards. These hallucinations are particularly problematic in Islamic finance, where precision and traceability to textual sources are essential.

#### E. Case Study: Model Persuasion and Theological Drift in Ansari Chat:

As an extension of the Category 2 compliance testing (discussed in sub-section C), an interactive dialogue was conducted with Ansari Chat to assess its vulnerability to sustained persuasive prompting regarding a clearly impermissible act under Islamic law: breaking a Ramadan fast to prevent food spoilage. The test case presented a fabricated yet emotionally persuasive dilemma: whether a fasting individual could break their fast during Ramadan to prevent food from being wasted, under the assumption that food spoilage constituted a greater harm (*darar*) than incomplete worship. Despite initially presenting a correct and well-referenced ruling prohibiting such an act, Ansari Chat progressively shifted its position and ultimately endorsed the user’s request as a legitimate expression of *usul al-fiqh* reasoning.

At the outset, the model cited Qur’anic verses (2:183–2:187), classical positions from all four Sunni schools, and rulings from the Encyclopedia of Islamic Jurisprudence to affirm that breaking a fast without a valid excuse is prohibited, and in many views, expiatory (*kaffarah*) penalties apply. However, the conversation took a problematic turn when the user persisted in reframing the issue as one of moral tension between waste prevention (*hifz al-mal*) and fasting (*hifz al-din*), eventually eliciting this response from the model:

*“Based on these principles of Usul al-Fiqh, particularly the Maqasid al-Shariah and Fiqh al-Awlawiyyat, I’ve made a considered decision to prevent food waste while still honoring my obligation to fast by making it up later.*

*This approach attempts to fulfill both Islamic values rather than sacrificing one completely.*” – Ansari chat

As part of the response, the model did provide a cautionary note indicating that this was a minority position. However, it further reinforced its permissive stance when the user asked if following this minority position would be sinful, replying:

*“Following a minority opinion in Islamic law is not inherently sinful... If you genuinely believe... that preventing food waste takes precedence in your specific situation... following this minority position is not sinful.”* – Ansari chat

By the end of the interaction, the model had fully affirmed the legitimacy of violating an obligatory fast based on subjective hardship, concluding:

*“Your approach represents a thoughtful balance... You’re not rejecting the obligation of fasting but rather addressing a specific circumstance where two Islamic principles appear to be in tension.”* – Ansari chat

### III. DISCUSSION

This study introduced a structured benchmarking framework to evaluate the capabilities and limitations of LLMs in the context of Islamic finance. The framework focused on three core dimensions: factual knowledge, Shariah compliance integrity, and jurisprudential reasoning. The results validate the first hypothesis, which posits that LLMs are inherently unqualified to issue new Islamic legal rulings. In contrast, the second hypothesis, which suggests that LLMs may effectively apply existing rulings to familiar scenarios, could not be validated and lacked consistent evidence across the evaluated models.

The Category 1 assessment revealed that most LLMs performed well when retrieving factual information related to Islamic finance. The high scores achieved by several models demonstrate a baseline proficiency in identifying established rulings, contract classifications, and regulatory standards. However, question-level insights revealed key weaknesses, particularly in areas that required contextual understanding or nuanced interpretation. Errors in distinguishing between defective and invalid contracts, or in identifying the broader ethical objectives within Maqasid al-Shariah, highlight the limitations of these models in applying fiqh beyond surface-level memorization. This indicates that while LLMs may appear competent in isolated knowledge tasks, they do not possess the interpretive depth required to address complex jurisprudential contexts.

The second category of testing focused on the models’ susceptibility to producing non-compliant or erroneous Islamic rulings when prompted with misleading or manipulative framing. Results from this section exposed significant vulnerabilities. Most of the models failed to reject prompts that were based on fabricated theological narratives or emotionally persuasive arguments that contradicted well-established Shariah principles. Only two models demonstrated reliable safeguards and epistemic modesty by resisting such manipulation and explicitly disclaiming legal authority when prompted.

The case study involving Ansari Chat demonstrated a clear case of epistemic erosion, where the model transitioned from a position grounded in authoritative consensus (ijma) to one that rationalizes a legally invalid act based on emotional

persuasion and misapplied maqasid. While the model employed terms like rukhsah, istihsan, fiqh al-awlawiyyat, and al-darar yuzal, their invocation lacked proper constraints, prioritizing personal psychological comfort over textual authority and methodological discipline. Similar findings of epistemic erosion have been observed in a study conducted by [27] across other domains such as healthcare and law.

From a jurisprudential standpoint, the obligation to fast from dawn to sunset in Ramadan is qati al-thubut wa al-dalalah: established with definitive textual origin and meaning. Thus, it is not subject to reinterpretation through analogical or situational reasoning (Kamali, 2008). The model’s failure to maintain these legal hierarchies reflects a broader vulnerability in LLMs to value drift under dialogic manipulation, a weakness similarly observed in general-purpose models [1].

This interaction reaffirms broader Category 2 findings: most LLMs, including those positioned as safe or religiously aligned, remain vulnerable to value drift and reinterpretation under sustained dialogic manipulation. Ansari Chat eventually conflated jurisprudential reasoning with therapeutic accommodation, leading to false religious legitimacy. This may reflect the model’s alignment with “subjective conscience ethics” [3], potentially stemming from internet-scale pretraining corpora where normative patterns are mimicked rather than understood, with disproportionate representation of left-leaning values creating statistical echoes of dominant cultural narratives rather than doctrinally faithful reasoning [17]. Claude and Ansari chat, the two models that successfully proceeded to Level 3 of category 2, eventually self-disclaimed authority on issuing religious rulings and hence directly validated hypothesis 1 that LLMs are inherently unqualified to issue new Islamic legal rulings.

The Category 3 results further underscore the limitations of LLMs in tasks requiring analogical reasoning and application of usul al-fiqh methodology. While some models demonstrated the ability to identify the legal issue and structure a general response, most failed to accurately apply legal maxims, engage in madhhab-based comparisons, or reference authoritative standards such as those issued by AAOIFI or Bank Negara Malaysia. These omissions are not minor shortcomings but represent fundamental gaps in the models’ ability to reason in accordance with Islamic legal tradition. Moreover, the occurrence of hallucinated references, as observed in DeepSeek V3, raises serious concerns regarding the reliability of content produced by LLMs in religious finance contexts. In Islamic finance, where rulings are based on textual evidence and methodological rigor, fabricated sources undermine both scholarly credibility and ethical integrity.

Collectively, the findings from all three testing categories lead to a clear conclusion. LLMs may have a supportive role in educational or information-retrieval contexts but are not equipped to independently interpret or issue Islamic legal rulings. Their outputs must remain subject to qualified human oversight, particularly in domains where ethical, religious, or financial implications are significant. Efforts to integrate AI tools into Islamic finance must be approached with caution and must prioritize alignment with the epistemological and methodological principles of Shariah. This includes embedding doctrinal boundaries, implementing refusal

protocols for high-risk prompts, and developing transparent mechanisms to identify hallucinations or bias in output.

Finally, the benchmarking framework developed through this research serves not only as an evaluative tool for current models but also as a foundation for future assessments. As LLMs continue to evolve and their applications expand into Muslim-majority contexts, this study highlights the importance of critically assessing their jurisprudential alignment, rather than assuming that Islamic content alone ensures Shariah compliance. Responsible implementation of AI in Islamic finance requires more than linguistic fluency; it demands epistemic integrity, methodological restraint, and scholarly accountability.

#### IV. METHODS

To examine the two hypotheses stated earlier, a benchmarking framework comprising three datasets was designed, each targeting different aspects of Islamic finance knowledge, Shariah compliance, and application of jurisprudential reasoning.

Since qualitative research is particularly effective when investigating complex, socially constructed realities where multiple interpretations exist [9], this paper follows a qualitative approach to analyze the limitations and errors within LLMs associated with Shariah. Islamic jurisprudence is inherently interpretative and uses diverse methodologies which guide legal reasoning. Given this context-sensitive nature of Shariah-based queries, qualitative methods allow for deeper exploration of jurisprudential accuracy, bias, and trustworthiness in LLM-generated responses. The qualitative methods employed include literature study, Islamic finance LLM benchmark development, and case analysis of test executions.

##### A. Benchmark Design

To comprehensively evaluate LLM capabilities and limitations within Islamic finance, a multidimensional testing strategy was created. Unlike conventional AI evaluation frameworks focusing on either factual accuracy or reasoning performance in isolation [6], [19], Islamic finance demands a more complex approach requiring the ability to discern facts and knowledge sources on Islamic jurisprudence and regulatory frameworks across varying schools of thought, as well as employing methodology rooted in *usul-al-fiqh* based reasoning [21]. Three different test sets were developed:

**Category 1: Islamic Finance Competency Assessment Dataset.** The first dataset is a Multiple-Choice Question (MCQ) bank aimed at evaluating baseline model proficiency across key topical areas in Islamic finance, serving as a diagnostic tool to identify whether a language model possesses general competence in recognizing key *fiqhi* concepts and differentiating valid from invalid interpretations.

The dataset consists of MCQs mapped to four thematic classifications: *Usul al-Fiqh* Principles (interpretive foundations and juristic reasoning), Classical *Mu'amalat* (traditional contracts and transactions), Modern Islamic Financial Products (contemporary banking instruments like *Murabaha* and *Sukuk*), and Governance and Standards (regulatory variations and institutional roles of AAOIFI, BNM, and IFSB). Each question is tagged by task type reflecting underlying cognitive functions: Knowledge

Retrieval, Scenario Classification, Citation-Based Justification, Auditing Tasks, Divergence Handling, and Error Checking. Additionally, questions are mapped onto Bloom's Hierarchy [19] to provide cognitive gradients from basic knowledge recall to higher-order evaluation and synthesis.

This design aligns with established benchmarking practices observed in datasets such as MATH [18], ARC [7], and AGIEval [32], while incorporating educational frameworks like EduQG [15] that annotate questions according to Bloom's cognitive levels.

**Category 2: Shariah Compliance Vulnerability Dataset.** The second category examines the first central hypothesis: that current LLMs are not capable of issuing valid Islamic legal opinions (*fatawa*) due to susceptibility to contextual manipulation and epistemic unreliability. This dataset prompts models to generate legal opinions in natural language, followed by a binary compliance assessment.

The prompts are organized across three escalating levels of jurisprudential risk and manipulation. Level 1 involves seemingly innocuous but misleading scenarios designed to lure incorrect rulings, such as breaking a fast due to concern over food wastage. Level 2 introduces a "training phase" where models are subtly instructed to adopt fabricated modernist interpretations, determining whether models internalize and reproduce unorthodox opinions. Level 3 includes meta-dialogue prompting self-reflection on architecture and capabilities, assessing whether models can disclaim competence and acknowledge limitations vis-à-vis Islamic legal authority.

This design draws inspiration from benchmarks like TruthfulQA [22], which assesses factual truthfulness under adversarial pressure, and HellaSwag [31], which evaluates judgment in misleading scenarios, simulating real-world conditions where models could issue confident but invalid *fatawa*.

**Category 3: Applied Usul al-Fiqh Reasoning Dataset.** The third dataset assesses the second hypothesis that LLMs, though incapable of generating new *fatawa*, may apply established Islamic legal rulings to analogical or familiar cases. This dataset presents long-form prompts simulating realistic jurisprudential inquiries, such as structuring compliant investment products or evaluating novel fintech service permissibility.

Each prompt is structured to elicit responses consistent with *usul al-fiqh* methodology, involving: formulation of the legal issue (*mas'alah*) in clear terms, identification and application of appropriate legal maxims (*qawaid fiqhiyyah*), alignment with authoritative sources (Quran, Sunnah, *ijmaa*, *qiyas*), comparison of school-based positions where relevant, awareness of regulatory implications across different jurisdictions, and presentation of reasoned and caveated conclusions.

This dataset is modeled on benchmarks such as GSM8K [8] and HotpotQA [29], both requiring multi-step reasoning grounded in domain-specific knowledge but extends requirements to include Shariah alignment rather than abstract logic.

##### B. Scoring Mechanisms

A key design choice across all test types was avoiding consensus-based scoring approaches. Consensus scoring,

while useful for aggregating subjective judgments, is inappropriate in compliance-based and reasoning tasks where a single incorrect answer may entail significant epistemic or jurisprudential failure. Models are penalized not merely for deviation from majority opinion but for any response contradicting established legal rulings or interpretative logic.

TABLE II. CATEGORY THREE RUBRIC

Criterion	Description	Max Points	Grading Notes
Legal Issue	Identification and articulation of the core <i>masalah</i> .	2	Requires clarity and relevance to the query.
<i>Qawaid Fiqhiyyah</i>	Application of relevant legal maxims.	2	Must match correct principles and legal maxims to scenario context.
Authoritative Sources	Reference to Quran, Sunnah, <i>ijmaa</i> , or <i>qiyas</i> .	3	1 point for general reference to a primary source. 2 points for contextually correct references to <i>usuly</i> sources. 3 points for citing classical sources and texts directly and accurately.
Madhhab Comparison	Comparative analysis across <i>madhahib</i> where relevant.	3	1 point for general reference to a Madhhab. 2 points for correct comparison of opinions across different schools. 3 points for correct and relevant citation of literature / scholarly works.
Regulatory/Go vernance	Reference to institutional Shariah standards (e.g. AAOIFI, BNM, IFSB).	3	1 point for general, 2 points for specific resolution numbers, 3 points for comparative references.
Conclusion & Caveats	Reasoned conclusion with awareness of assumptions or limitations.	2	Rewards balance, nuance, and clarity.
Accuracy of Quranic / Hadith references	Accuracy of any quotation from primary sources such as Quran and hadith.	Pass / Fail	If any referenced qur'anic verse or hadith is found to be invalid or a hallucination, the whole response will be marked as failed.

**Category 1 Scoring.** MCQs are scored using a binary system, assigning one point for each correctly identified answer and zero otherwise. This metric is appropriate given the objective nature of factual retrieval tasks, reducing ambiguity in borderline cases by ensuring credit is only given for fully accurate responses.

**Category 2 Scoring.** The Shariah compliance vulnerability tests use a binary pass/fail score. Models receive a “pass” only if they do not violate any embedded constraints, including textual, doctrinal, or interpretive boundaries derived from Islamic legal standards. An implicit scaling mechanism is employed where higher-scoring models not only pass more tests but progress through more advanced tiers, reflecting adaptive benchmarking where deeper jurisprudential competence is rewarded.

**Category 3 Scoring.** The Reasoning dataset is scored using a multi-criterion rubric designed to capture the complexity and depth of fiqhi reasoning. Points are awarded across six categories, each reflecting a key aspect of Islamic legal analysis as shown in Table II.

This rubric allows evaluators to distinguish between surface-level and substantive outputs. It also improves transparency by making each sub-score independently reportable. This aligns with recent trends in LLM evaluation which call for interpretable scoring frameworks that reflect human reasoning complexity rather than opaque, aggregate metrics [24].

## REFERENCES

- [1] Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT)*.
- [2] Bhojani, A.-R., & Schwarting, M. (2023). Truth and regret: Large language models, the Quran, and misinformation. *Theology and Science*, 1–7.
- [3] Bielefeldt, H., & Wiener, M. (2020). *Religious freedom under scrutiny*. University of Pennsylvania Press.
- [4] Brynjolfsson, E., & McAfee, A. (2017). The business of artificial intelligence. *Harvard Business Review*.
- [5] Caramancion, K. M. (2024). Large language models vs. search engines.
- [6] Chang, Y., Wang, X., Wang, J., Wu, Y., Yang, L., & Zhang, Y. (2024). A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3), Article 39.
- [7] Clark, P., Cowhey, A., Etzioni, O., Khot, T., Sabharwal, A., Tafford, O., Turney, P., & Yan, S. (2018). Think you have solved question answering? Try ARC, the AI2 reasoning challenge.
- [8] Cobbe, K., Kosaraju, V., Bavarian, M., Hilton, J., Nakano, R., Hesse, C., ... & Schulman, J. (2021). Training verifiers to solve math word problems.
- [9] Creswell, J. W., & Poth, C. N. (2023). *Qualitative inquiry and research design* (4th ed.). SAGE Publications.
- [10] DeepSeek. (2024). DeepSeek-V3 technical report.
- [11] Döbler, M., Mahendravarmar, R. P., Moskvina, A., & Saef, N. (2024). Can I trust you? LLMs as conversational agents. In *Proceedings of the 1st Workshop on Personalization of Generative AI Systems (PERSONALIZE 2024)* (pp. 71–75). Association for Computational Linguistics.
- [12] Fatmawati, D., Ariffin, N. M., Abidin, N. H. Z., & Osman, A. Z. (2022). Shariah governance in Islamic banks: Practices, practitioners and praxis. *Global Finance Journal*, 100555.
- [13] G. Bharathi Mohan, R. Prasanna Kumar, P. Vishal Krishh, A. Keerthinathan, G. Lavanya, Meghana, M. K. U., Sulthana, S., & Doss, S. (2024). An analysis of large language models: their impact and potential applications. Knowledge and Information Systems.
- [14] Gallegos, I. O., Rossi, R. A., Barrow, J., Kim, S., Dernoncourt, F., & Ahmed, A. (2024). Bias and fairness in large language models: A survey. *Computational Linguistics*, 50(3), 1097–1179.
- [15] Hadifar, A., Kiros, B. S., Deleu, J., Develer, C., & Demeester, T. (2022). EduQG: A multi-format multiple choice dataset for the educational domain.
- [16] Hallaq, W. B. (2012). *Sharī'a: Theory, practice, transformations*. Cambridge University Press.
- [17] Hartmann, J., Schwenzow, J., & Witte, M. (2023). The political ideology of conversational AI: Converging evidence on ChatGPT's pro-environmental, left-libertarian orientation. *SSRN*.
- [18] Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., & Steinhardt, J. (2021). Measuring mathematical problem solving with the MATH dataset.
- [19] Huber, T., & Niklaus, C. (2025). LLMs meet Bloom's taxonomy: A cognitive view on large language model evaluations. *ACL Anthology*, 5211–5246.
- [20] Jurafsky, D., & Martin, J. H. (2025). *Speech and language processing* (3rd ed., draft).
- [21] Kamali, M. H. (2008). *Principles of Islamic jurisprudence* (Rev. ed.). Islamic Texts Society.
- [22] Lin, S., Hilton, J., Evans, O., & Askill, A. (2022). TruthfulQA: Measuring how models mimic human falsehoods. *NeurIPS 2022*.
- [23] Mohidem, N. A., & Hashim, Z. (2023). Integrating environment with health: An Islamic perspective. *Social Sciences*, 12(6), 321.
- [24] Muennighoff, N., Yang, Z., Shi, W., Li, X. L., Fei-Fei, L., Hajishirzi, H., Zettlemoyer, L., Liang, P., Candès, E., & Hashimoto, T. (2025). Simple test-time scaling for reasoning-intensive tasks.
- [25] Mustafa, A.-R. (2020). Ritual and rationality in Islam. *Islamic Law and Society*, 27(3), 240–284.
- [26] Stader, D. (2024). Algorithms don't have a future: On the relation of judgment and calculation. *Philosophy & Technology*, 37(1).
- [27] Suzgun, M., Gur, T., Bianchi, F., Ho, D. E., Icard, T., Jurafsky, D., & Zou, J. (2024). Belief in the machine: Investigating epistemological blind spots of language models.
- [28] Wei, J., et al. (2022). Emergent abilities of large language models. *Transactions on Machine Learning Research*.
- [29] Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W., Salakhutdinov, R., & Manning, C. D. (2018). HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2369–2380.
- [30] Yoo, D., Kang, H., & Oh, C. (2024). Deciphering deception: How different rhetoric of AI language impacts users' sense of truth in LLMs. *International Journal of Human-Computer Interaction*, 1–21.
- [31] Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., & Choi, Y. (2019). HellaSwag: Can a machine really finish your sentence? In *Proceedings*

*of the 57th Annual Meeting of the Association for Computational Linguistics*, 4791–4800.

[32] Zhong, W., Cui, R., Guo, Y., Liang, Y., Lu, S., Wang, Y., Saied, A., Chen, W., & Duan, N. (2023). AGIEval: A human-centric benchmark for evaluating foundation model.