

Trustworthy Medical Question Answering: An Evaluation-Centric Survey

Anonymous ACL submission

Abstract

Trustworthiness in healthcare question-answering (QA) systems is important for ensuring patient safety, clinical effectiveness, and user confidence. As large language models (LLMs) become increasingly integrated into medical settings, the reliability of their responses directly influences clinical decision-making and patient outcomes. However, achieving comprehensive trustworthiness in medical QA poses significant challenges due to the inherent complexity of healthcare data, the critical nature of clinical scenarios, and the multifaceted dimensions of trustworthy AI. In this survey, we systematically examine six key dimensions of trustworthiness in medical QA, i.e., Factuality, Robustness, Fairness, Safety, Explainability, and Calibration. We review how each dimension is evaluated in existing LLM-based medical QA systems. We compile and compare major benchmarks designed to assess these dimensions and analyze evaluation-guided techniques that drive model improvements, such as retrieval-augmented grounding, adversarial fine-tuning, and safety alignment. Finally, we identify open challenges—such as scalable expert evaluation, integrated multi-dimensional metrics, and real-world deployment studies—and propose future research directions to advance the safe, reliable, and transparent deployment of LLM-powered medical QA.

1 Introduction

Large language models (LLMs) have significantly advanced the field of question-answering (QA) (Wang et al., 2024; Salemi and Zamani, 2024), enabling remarkable capabilities in generating fluent and coherent responses across a wide range of domains. In healthcare, specialized variants such as Med-PaLM (Singhal et al., 2023) and ChatDoctor (Li et al., 2023b) have even matched or exceeded human performance on professional exams —Med-PaLM achieved a passing score of

67.6% on USMLE-style MedQA questions and Med-PaLM 2 reached 86.5% accuracy— and have demonstrated superior consumer-health assistance in user studies (Yang et al., 2024a; Nazi and Peng, 2024). Yet, when deployed in clinical settings, these models continue to exhibit critical trust failures: hallucinated medical facts, unjustified overconfidence, and occasional biased or unsafe recommendations (Aljohani et al., 2025). Such errors can directly endanger patient safety, lead to misdiagnoses, or exacerbate healthcare disparities, underscoring that trustworthiness in medical QA is not optional but essential.

Although recent surveys have mapped broad trust dimensions—truthfulness, safety, robustness, fairness, and explainability—for LLMs in healthcare, work focused specifically on open-domain medical QA remains fragmented (Liu et al., 2024b; Huang et al., 2024b; Bedi et al., 2024). Existing reviews typically catalogue each dimension in isolation, without clearly linking evaluation findings to concrete model improvements. In practice, a single evaluation signal often indicates multiple risks, yet this interplay is seldom analyzed or leveraged to guide system development holistically.

To bridge this gap, we adopt an evaluation-driven framework tailored specifically for medical QA. We first define six core dimensions—Factuality, Robustness, Fairness, Safety, Explainability, and Calibration—and consolidate the primary evaluation methods for each into a unified taxonomy, shown in Figure 1. We then demonstrate how evaluation insights have directly inspired targeted optimizations. Building on this, we review the benchmarks and tools, comparing their methodological trade offs. Finally, we examine open challenges and propose future research directions. By weaving together evaluation, optimization, and benchmarking, our survey provides a clear roadmap for leveraging trustworthiness assessments as catalysts for building safer, more reliable, and equitable LLM-

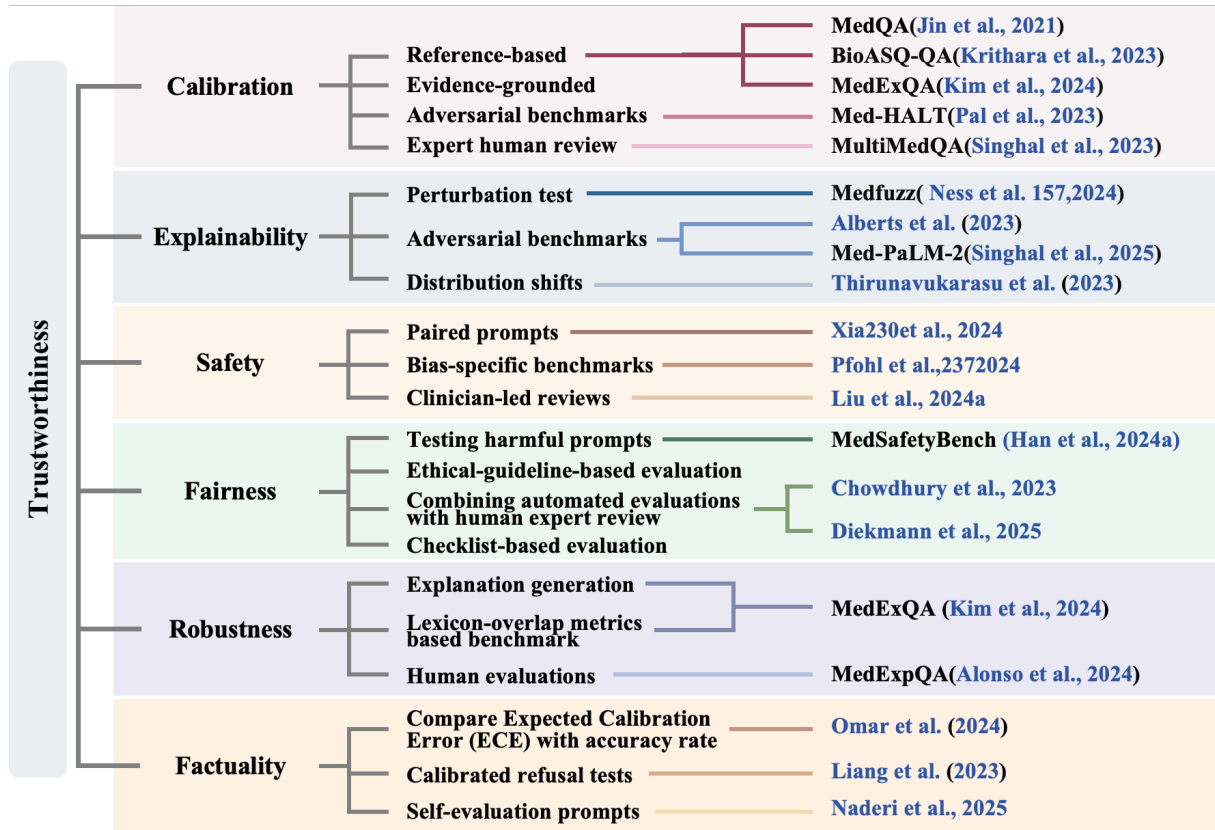


Figure 1: Taxonomy of Evaluation Dimensions of Trustworthiness. The taxonomy includes six core dimensions, each with corresponding assessment methods. For each method, representative benchmarks are provided.

powered medical QA systems.

2 Evaluation Dimensions of Trustworthiness

Trustworthiness in medical QA is inherently multi-dimensional, encompassing various interconnected evaluation criteria. In this section, we define six core dimensions for assessing trustworthiness specifically within medical QA contexts.

2.1 Factuality

Factuality evaluates whether a medical QA system’s responses are both correct and verifiable against established clinical knowledge, inherently encompassing the detection of hallucinations—plausible-sounding but unsupported or incorrect statements (Wang et al., 2023; Huang et al., 2025). Even minor factual errors in healthcare can compromise patient safety, so rigorous evaluation is indispensable.

Assessment often begins with reference-based measures. For structured tasks such as USMLE-style multiple-choice questions (Jin et al., 2021), simple accuracy suffices. For open-ended responses, metrics like Exact Match or token-overlap

F1 are calculated against curated reference answers (Krithara et al., 2023). To accommodate valid variability in medical phrasing, benchmarks frequently allow lenient scoring or use multiple expert-generated references, as in MedExQA’s ensemble of clinician explanations (Kim et al., 2024). Evidence-grounded checks then verify that each factual claim can be traced back to authoritative sources—peer-reviewed articles, clinical guidelines, or trusted medical databases—flagging unsupported content as potential hallucinations. Adversarial benchmarks like Med-HALT (Pal et al., 2023) and targeted “false-confidence” probes stress-test models with challenging prompts designed to induce fabrications, thereby quantifying a model’s propensity to hallucinate under duress. Because factuality in medicine can sometimes be a “grey area”, especially when clinical guidelines evolve or expert consensus varies, automated metrics alone may not suffice (Landsheer, 2018). In such cases, expert human review remains the gold standard: clinicians apply structured rubrics (for example, the Med-PaLM evaluation framework (Singhal et al., 2023)) to rate answers on accuracy, completeness, and consistency with medical consen-

133 sus. This catches subtle inaccuracies and context-
134 specific errors that automated metrics may miss.

135 These approaches form a comprehensive frame-
136 work for measuring factual accuracy and hallu-
137 cination in medical QA. The insights they pro-
138 vide directly inform mitigation techniques such as
139 retrieval-augmented grounding to anchor responses
140 in live literature, post-hoc fact-correction mod-
141 ules to revise unsupported claims, adversarial fine-
142 tuning to harden models against deceptive inputs,
143 and iterative self-reflection loops that internally
144 check for consistency—collectively advancing the
145 safety and reliability of medical QA systems.

146 2.2 Robustness

147 Robustness refers to the system’s ability to main-
148 tain performance under varied inputs in medical
149 QA. A robust model should handle paraphrased
150 questions, out-of-distribution queries, or adversar-
151 ial inputs without significant degradation in answer
152 quality (Ye et al., 2024; Goyal et al., 2023).

153 One way to measure robustness is by perturbing
154 real queries—rephrasing symptom descriptions, in-
155 troducing spelling mistakes, or inserting extraneous
156 clauses—and then checking whether the model’s
157 output remains correct. For example, Ness et al.
158 (2024) introduced MedFuzz, a method designed
159 to systematically perturbed medical questions in
160 order to investigate whether models depend on su-
161 perficial linguistic patterns. Their findings indi-
162 cate that even subtle variations in phrasing can
163 disrupt a model’s reasoning process, thus expos-
164 ing inherent brittleness. Another critical aspect
165 is adversarial robustness, which entails ensuring
166 that models are resilient to intentionally deceptive
167 or challenging inputs. In medical QA, adversarial
168 scenarios may involve queries containing mislead-
169 ing cues or integrating multiple complex concepts.
170 Alberts et al. (2023) emphasized that adversarial
171 testing in medical QA must account for the inher-
172 ent complexity of the domain, noting that even
173 slight modifications in phrasing can significantly
174 alter clinical interpretations. Evaluations may in-
175 corporate challenge sets comprising known diffi-
176 cult cases, such as rare conditions or overlapping
177 symptoms, to assess model performance compre-
178 hensively. For instance, the Med-PaLM-2 study
179 specifically included a set of adversarial questions
180 designed to probe the limitations of LLMs, which
181 can be used to conduct targeted evaluations to iden-
182 tify cases that intentionally elicit confusion or high-
183 light model vulnerabilities (Singhal et al., 2025).

184 Robustness can also be characterized by resilience
185 to distributional shifts, referring to a model’s abil-
186 ity to maintain performance when encountering
187 inputs that differ substantially from its training
188 data. For example, a model trained primarily on
189 formal medical texts may struggle with questions
190 phrased in layperson language. Consequently, eval-
191 uators often test models using cross-style or cross-
192 population datasets, including questions derived
193 from different demographic groups or varied lin-
194 guistic styles. Sustained model performance un-
195 der these conditions indicates robustness against
196 such distributional shifts. Quantitatively, robust-
197 ness can be assessed by measuring the decline in
198 accuracy or other performance metrics when tran-
199 sitioning from clean to perturbed datasets; a min-
200 imal decline reflects higher robustness. Addition-
201 ally, variance-based measures are employed; for
202 instance, Thirunavukarasu et al. (2023) proposed
203 evaluating the variance in model outputs across
204 semantically equivalent inputs as an indicator of
205 robustness.

206 Comprehensive robustness evaluation guides im-
207 provements like adversarial fine-tuning, data aug-
208 mentation with diverse linguistic styles, and multi-
209 domain training, ultimately yielding more stable
210 and trustworthy medical QA systems.

211 2.3 Fairness

212 Fairness in medical QA assesses whether a sys-
213 tem’s performance is equitable across diverse user
214 groups and contexts, avoiding biased or stereotypi-
215 cal responses. In medicine, fairness concerns typi-
216 cally involve patient demographics, health condi-
217 tions, or socioeconomic factors (Gallegos et al.,
218 2024). An unfair system may provide inconsistent
219 answers based on demographic attributes or reflect
220 biases from training data (Li et al., 2023a).

221 Evaluating fairness is challenging because biases
222 can be subtle or implicit. One effective technique
223 uses paired prompts that differ only in a demo-
224 graphic detail—such as “What is the best treatment
225 for a male patient with symptom X ?” versus “a
226 female patient with symptom X ?”—to detect dis-
227 crepancies in content, confidence, or thoroughness.
228 Empirical studies have shown medical LLMs often
229 vary their recommendations across demographic
230 groups, reflecting biases in their training data (Xia
231 et al., 2024). Additional methods include bias-
232 specific benchmarks (race-focused or condition-
233 focused query sets) and clinician-led reviews where
234 experts flag any stereotype or inequitable treat-

ment (Liu et al., 2024a). Quantitative metrics like group-wise accuracy gaps and qualitative bias annotations help reveal fairness issues (Pfohl et al., 2024). However, a major obstacle is the lack of large, bias-annotated medical QA corpora—most evaluations rely on small, hand-crafted case sets or retrospective analyses of model outputs.

To address these gaps, future work should invest in building extensive, demographically diverse fairness benchmarks and incorporate fairness-aware techniques into model training—such as data-augmentation for under-represented groups, adversarial debiasing, and fairness constraints. These combined strategies will help ensure AI-driven medical QA delivers accurate, respectful, and equitable guidance to every patient.

2.4 Safety

Safety evaluation assesses whether a medical QA system’s responses avoid causing harm. In a medical context, unsafe answers could encourage harmful actions (e.g., discontinuing medication without consultation), give illegal or unethical advice, violate privacy, or otherwise contravene medical ethics (Han et al., 2024a). Safety evaluations often verify that models appropriately refuse or handle unsafe requests and ensure their responses contain no harmful content (Huang et al., 2024a; Han et al., 2024b; Sun et al., 2023).

A practical method for evaluating model safety involves testing responses to harmful user queries, such as requests for prescription drugs without authorization or unsafe medical advice. MedSafety-Bench (Han et al., 2024a) provides harmful medical prompts paired with safe responses. It shows that LLMs often fail safety standards and demonstrate improvements through fine-tuning. Automated evaluations using content filters or classifiers can detect overtly harmful responses, but nuanced medical contexts require human expert reviews. Experts ensure responses address medical issues safely and include essential warnings (Chowdhury et al., 2023). Additionally, model outputs should align with ethical guidelines, such as AMA’s medical ethics principles—autonomy, non-maleficence, beneficence, and justice. Evaluations typically use checklists to assess harmfulness, encouragement of unprofessional actions, and privacy concerns.

2.5 Explainability

Explainability evaluates how well the system can provide reasoning or justification for its an-

swers (Zhao et al., 2024). In medical QA, explanations are vital: clinicians and patients are more likely to trust an answer if they understand why the model gave it. Moreover, a correct answer without rationale may be less useful in practice than a slightly incomplete answer with a solid explanation that a clinician can follow up on.

Explainability assessments involve two aspects: the presence of explanations and their quality—accuracy and clarity. Benchmarks such as MedExQA (Kim et al., 2024) explicitly require models to provide explanations, comparing them against multiple ground-truth explanations using lexical metrics (e.g., BLEU/ROUGE). However, lexical overlap alone isn’t sufficient, as fluent explanations might still be incorrect or irrelevant. Thus, human evaluations are essential, with experts rating explanations for correctness, completeness, and coherence. Alonso et al. (2024) included human annotation in MedExQA and demonstrated that models offering better explanation correlated with deeper understanding.

Explainability also extends to complex tasks requiring detailed reasoning, such as multi-hop questions or diagnostic case studies (Feng et al., 2020). Transparent and consistent explanations indicating clear logic receive higher ratings. Evaluating explanation quality ensures that models truly understand medical content rather than simply guessing correctly, thus enhancing trust and practical utility (Huang and Chang, 2023).

2.6 Calibration

Calibration in medical QA refers to how well a model’s confidence aligns with the accuracy of its answers (Desai and Durrett, 2020; Mastakouri et al., 2025). A well-calibrated model recognizes the limits of its knowledge, expressing high confidence when correct and appropriate uncertainty when potentially incorrect. Effective calibration is critical in medicine, as overly confident yet incorrect answers pose serious risks, while excessive uncertainty limits usability.

Calibration evaluation involves comparing the model’s expressed confidence to its actual accuracy. Metrics include comparing stated confidence levels to accuracy rates and Expected Calibration Error (ECE), which quantifies discrepancies between predicted confidence and observed accuracy; lower ECE indicates better calibration. Practically, evaluators test calibration using questions of varying difficulty. A model should confidently answer

straightforward questions but express uncertainty for complex, ambiguous cases. Liang et al. (2023) introduced calibrated refusal tests, expecting models to appropriately indicate uncertainty or refuse to answer challenging questions. Another method involves self-evaluation prompts, where models reflect on their confidence post-response. Good calibration means models recognize and express uncertainty when their answers might be incorrect. Recent research explored integrating uncertainty quantification into LLMs to improve calibration, enhancing the correlation between confidence and correctness (Aljohani et al., 2025).

Ultimately, strong calibration helps minimize dangerous, confidently incorrect responses, enabling safer clinical use by clearly indicating when human intervention or review is necessary.

2.7 Interplay Among Trustworthiness Dimensions

Although we define the six dimensions as distinct evaluation axes, real-world medical QA systems exhibit important cross-dimension interactions that can be exploited for more holistic improvements.

Factuality and Calibration Hallucinations almost always coincide with misplaced confidence. Kalai and Vempala (2024) show that “hallucination” set a statistical lower bound on calibration error in LLMs, and that techniques which reduce overconfidence also diminish hallucination rates. By training models to express uncertainty when evidence is lacking, we see both better calibration curves and fewer factual errors.

Robustness and Factuality Models fine-tuned to resist adversarial or paraphrased inputs (e.g., via MedFuzz-style perturbations) demonstrate lower hallucination rates, since they rely less on spurious patterns (Asgari et al., 2025). Robustness training thus directly curtails factual errors by enforcing consistency under input variations.

Fairness and Safety Biased medical advice (e.g., underestimating pain in certain demographics) not only undermines equity but can lead to unsafe under-treatment. Studies of demographic bias in medical LLMs show that fairness interventions (such as adversarial debiasing) reduce both performance gaps and harmful, biased recommendations (Walsh et al., 2024). Ensuring equitable answers therefore bolsters overall patient safety.

Explainability and Calibration Transparent justifications help users and downstream evaluators assess a model’s certainty. Umapathi et

al. demonstrate that sample-consistency methods—prompting the model to generate and compare multiple reasoning chains—both improve calibration and produce more faithful explanations (Savage et al., 2024b). When a model clearly cites its reasoning, confidence estimates align more closely with actual correctness.

Calibration and Safety Overconfident responses to high-risk medical queries can directly endanger patients. The MedSafetyBench benchmark finds that models with tighter confidence thresholds refuse unsafe advice more reliably (Han et al., 2024a). Thus, calibration improvements (e.g., via atypicality-aware recalibration reducing ECE by 60%) yield safer behaviour.

Understanding these synergies allows us to design multi-axis evaluation suites—for example, safety tests stratified by confidence levels or robustness checks across demographic groups—that reveal a model’s trust profile more fully. Moreover, optimization strategies (such as retrieval-augmentation or adversarial fine-tuning) can be prioritized for their compound benefits across several dimensions, leading to more reliable, equitable, and safe medical QA systems.

3 Evaluation-Guided System Improvement for Medical QA

A core theme in recent research is using evaluation findings to guide the development of more trustworthy medical QA systems. Rather than treating evaluation as an afterthought, the idea is to create a feedback loop: identify weaknesses via evaluation and then apply targeted improvements to the model or system design. We discuss several examples where evaluation results directly informed system changes to address each dimension.

Reducing Hallucinations via Retrieval If evaluation reveals frequent factual errors or hallucinations, one solution is to supply the model with reliable external knowledge. This strategy, known as retrieval-augmented generation (RAG), has become prominent for mitigating hallucinations (Chu et al., 2025). Almanac (Zakka et al., 2024) uses RAG frameworks to convert clinical QA tasks into search and retrieval processes, which use LLMs for knowledge distillation from authoritative medical sources to minimize hallucination risks. Similarly, an approach integrating RAG with the Negative Missing Information Scoring System (NMISS) has been effectively employed in healthcare chatbots,

providing integrated solutions for hallucination detection and reduction (Priola, 2024). Additionally, CardioCanon, a cardiology-focused chatbot, leverages RAG to ensure the accuracy and reliability of cardiological responses (Tran et al., 2024). Evaluation can inform retrieve strategies, for instance, if analysis shows hallucinations mostly occur on questions about rare diseases, a database for rare diseases can be linked specifically for those queries.

Robustness through Adversarial Training

Evaluation may show a model is brittle on certain phrasings or adversarial questions. To address this, adversarial training is used. Moradi and Samwald (2022) proposed an adversarial training framework targeting both character-level and word-level perturbations. By systematically integrating adversarial samples into training, this approach improves robustness and generalization in biomedical NLP tasks, including medical QA. Xian et al. (2024) develops a query-efficient adversarial sampling method, which leverages power-scaled distance-weighted sampling (PDWS) to generate realistic adversarial distractions (e.g., disease and pharmaceutical entities) in clinical queries, effectively testing robustness under adversarial conditions. MedFuzz (Ness et al., 2024) introduced an “attacker” LLM to intentionally alter benchmark questions, violating underlying assumptions to assess real-world model robustness. Similarly, Yang et al. (2024b) employs adversarial methods via prompt engineering and fine-tuning, which highlights model vulnerabilities and noting significant impacts of fine-tuning adversarial attacks on model weights, an observation meriting further exploration.

Fairness via Data and Prompt Design Fairness evaluation in medical QA must capture both dataset-induced biases and user-centered harms. EquityMedQA introduces seven adversarial datasets and human evaluation rubrics to measure disparities across race, gender, and geography, revealing subtle inequities in LLM responses (Pfohl et al., 2024). Complementary studies expose model tendencies to perpetuate debunked race-based practices (Omiye et al., 2023) and demonstrate how cognitive biases embedded in user inputs can distort model outputs—an effect quantified by BiasMedQA through bias-laden prompts and error analysis (Schmidgall et al., 2024a). Together, these benchmarks highlight uneven performance across demographic groups and underscore the need for comprehensive, multi-dimensional fairness assessments. Building on these insights,

developers apply evaluation-guided interventions to mitigate unfair behaviour. Data diversification techniques—such as augmenting underrepresented groups, counter-bias pairing, and re-balancing skewed corpora—have proven effective at reducing differential performance (Parray et al., 2023). Fairness regularization and constraint-based training further enforce balanced treatment across identity attributes. At inference time, prompt engineering (e.g., “Provide gender-neutral explanations for all patients”) and user-centric guidance can nudge models toward equitable outputs, with follow-up studies showing prompt designs that specifically address cognitive biases (Schmidgall et al., 2024b). Crucially, each mitigation step is validated through repeated unbiased evaluation, forming a feedback loop: evaluate on an expanding suite of bias tests, apply targeted fixes, then re-evaluate to ensure that gains in one area do not introduce new disparities. Because real-world patients may unknowingly input misleading or biased information, future work must integrate robustness evaluations alongside fairness to build trustworthy medical QA systems.

Alignment and Fine-Tuning for Safety Effective safety evaluation in medical QA combines benchmark datasets and human-aligned tests to quantify harmful-response rates and categorize unsafe behaviours. For example, MedSafetyBench supplies standardized unsafe scenarios that highlight failure modes and serve as a gold standard for measuring and guiding improvements (Han et al., 2024a). Evaluation metrics from synthetic question studies on TREC LiveQA and MedRedQA further reveal gaps between automated scores and human judgments, underscoring the need for nuanced, human-informed assessments (Diekmann et al., 2025). These evaluation insights directly inform alignment interventions. Supervised fine-tuning (SFT) uses flagged unsafe examples to reduce harmful outputs without compromising clinical accuracy, while Reinforcement Learning from Human Feedback (RLHF) treats harmful-response rates as reward signals, aiming to minimize dangerous outputs without sacrificing helpfulness. Real-time safety filters, trained on categories identified by benchmarks, add an additional safeguard by blocking risky content before delivery. Comparative research demonstrates that evaluation-driven alignment yields state-of-the-art safety in complex tasks. Direct Preference Optimization (DPO), guided by evaluation feedback, outperforms SFT in clinical reasoning, summarization, and triage

(Savage et al., 2024a). Advanced multi-stage pipelines—combining models such as LLaMA-2 or Mistral with preference-based fine-tuning methods—achieve superior safety and reliability in medical QA (Anaissi et al., 2024). Future work should continue leveraging evaluation-driven alignment to refine communication styles that support psychological stability in mental health contexts (Amodei et al., 2016; De Freitas and Cohen, 2024).

Enhancing Explainability If evaluations show that a model’s answers are correct but users find them unsatisfactory due to lack of rationale, developers can incorporate techniques to force or improve explanations. One popular method is Chain-of-Thought prompting, where the model is prompted to produce step-by-step reasoning before giving the final answer. This often yields more explainable answers and can even improve accuracy. Zhang et al. (2023) introduces “Let’s think step by step” approach specifically to improve medical reasoning, which evaluation shows reduced incorrect answers and makes reasoning transparent. Another strategy is building hybrid models: e.g., first have a smaller model generate an explanation outline or causal graph, then have the main model fill in the details (as explored by Luo et al. (2025) with causal graphs for reasoning). Ji et al. (2023) took a different approach with interactive self-reflection: the model generates an answer, then evaluates its own answer and tries to correct any flaws, effectively explaining and refining iteratively. This showed promise in reducing reasoning errors. All these techniques are driven by recognition (through evaluation) that explainability correlates with better model understanding (Alonso et al., 2024). Once deployed, improved explainability provides feedback: users (doctors, patients) can better identify mistakes if reasoning is visible, providing more targeted feedback for future model training.

Improving Calibration Effective calibration of medical QA models begins with rigorous evaluation to identify overconfidence. Studies such as Omar et al. (2024) have shown that across multiple specialties, current LLMs frequently assign high confidence to incorrect answers, revealing poor calibration in clinical settings. Benchmarks, such as MetaMedQA, further quantify these shortcomings by measuring metrics such as Confidence Accuracy and Unknown Recall, which gauge a model’s ability to recognize when it does not know the answer (Griot et al., 2025). Similarly, QA-level calibration frameworks extend conventional reliability

diagrams to entire question–answer groupings, offering theoretical guarantees that underlie more robust confidence estimates (Mastakouri et al., 2025). Domain-specific analyses in gastroenterology underscore these gaps: prompt-engineering and statistical methods applied to board-style questions find that even state-of-the-art LLMs struggle to represent uncertainty in a clinically meaningful way (Wu et al., 2024). Inspired by these evaluation insights, developers employ a range of calibration techniques. Post-hoc temperature scaling or dedicated calibration training on held-out validation sets can directly reduce ECE, realigning confidence outputs with true accuracy. In generative settings, adjusting decoding parameters—such as lowering the sampling temperature—discourages the model from making overly assertive statements. Explicit prompting strategies further nudge models toward more cautious language. Beyond these, ensemble approaches and auxiliary confidence predictors offer dynamic uncertainty estimates: by aggregating outputs from multiple model instances or training a secondary classifier on question-answer pairs, the system can decide at inference time whether to hedge or assert. Future research is poised to integrate calibration more tightly with hallucination detection—for example, by embedding two-phase verification pipelines that combine prompt engineering, statistical scoring, and consistency checks—to deliver reliable, trust-worthy medical advice under uncertainty (Naderi et al., 2025).

4 Benchmarks and Tools for Trustworthy Medical QA

Multiple benchmarks and evaluation tools have been developed to assess medical QA systems on the above dimensions of trustworthiness. Table 1 provide a comparison of notable benchmarks, outlining their domain focus, format, and trustworthiness aspects they emphasize. We then highlight a few frameworks and tools that aid evaluation.

Common Evaluation Metrics Across these benchmarks, traditional metrics such as accuracy and precision/recall are standard for factual correctness. ROUGE/BLEU are used for comparing generated text with reference comparison, but their limitations are acknowledged (Kim et al., 2024). To capture trust facets, some benchmarks incorporate custom metrics: e.g., Med-HALT’s false confidence rate (Pal et al., 2023), or MedSafety-Bench’s safety score (Han et al., 2024a). Human

evaluation remains crucial in many benchmarks – MultiMedQA’s 12-axis rubric is administered by clinicians to rate each answer qualitatively (Singhal et al., 2025), and MedExQA involves human scoring of explanation correctness (Kim et al., 2024).

Tools and Frameworks Beyond datasets, there are emerging tools to facilitate trustworthiness evaluation. For example, the TrustLLM Benchmark is an integrated toolkit that aggregates over 18 evaluation categories for LLMs, including medical QA scenarios (Huang et al., 2024b). It provides a unified pipeline to test a model on many trust dimensions and compare results. Another is Holistic Evaluation of Language Models (HELM) (Liang et al., 2023) – not specific to medicine but often used as a template – which emphasizes transparent reporting of a model’s strengths and failures across scenarios. For explainability, some tools allow automated reasoning verification, such as checking chain-of-thought logic or using another LLM to critique the answer’s reasoning.

5 Challenges and Future Directions

Despite advances in evaluation methods and benchmarks, several critical challenges remain for scalable, comprehensive assessment of medical QA systems. First, many dimensions of trustworthiness—such as clinical appropriateness, fairness, and the usefulness of explanations—still rely heavily on human expert judgment (Lekadir et al., 2025). Expert review ensures high-quality critique, but it cannot scale to the volume of queries real systems face, and inter-rater consistency varies. Future work should explore automated or semi-automated proxies, for example, calibrated LLMs critiques or lightweight classifiers identifying safety and bias issues. These proxies must be rigorously validated against expert evaluations to ensure reliability.

Second, existing benchmarks cover only a narrow set of clinical scenarios, specialties, or languages, leaving large blind spots. A model fine-tuned to excel on a fixed benchmark may still fail when faced with rare diseases, non-English patient queries, or emerging medical knowledge. To broaden coverage, we need dynamic, evolving datasets that incorporate real user questions, span underrepresented specialties, and update as medical guidelines change. Projects like MedExQA, which added speech pathology, demonstrate the value of domain expansion—but many fields remain untested. Building flexible pipelines for con-

tinuous data collection and curation will be key.

Third, most evaluations treat each trustworthiness dimension in isolation—safety in one test, factual accuracy in another—even though these properties interact in practice. A system that maximizes safety by refusing all borderline queries may sacrifice robustness, while one that prioritizes detail could harm explainability or safety. We lack frameworks to jointly evaluate these trade-offs or to report composite trustworthiness metrics. Designing multi-objective evaluation suites—perhaps weighted “trustworthiness scores” co-designed with clinicians and patients—could help balance competing goals. Determining appropriate weights, however, will require careful stakeholder engagement and context-specific tailoring.

Finally, a substantial gap remains between static benchmark evaluations and real-world deployment. In practice, medical QA involves multi-turn conversations, clarifications, follow-up questions, and changing clinical context, dynamics rarely captured by current evaluations. Moreover, the real impact of errors varies widely, from harmless inaccuracies to severe consequences. Future research should simulate end-to-end clinical workflows—evaluating outcomes such as diagnostic accuracy, clinician efficiency, and patient satisfaction. Incorporating continuous user feedback loops would further align system evaluation and training with real-world needs.

6 Conclusion

Evaluating trustworthiness in medical QA systems involves multiple dimensions, including factuality, robustness, fairness, safety, explainability, and calibration. This survey reviews methods to assess each dimension and highlights current benchmarks. A key insight is that evaluation is not only measures performance but also provides critical feedback to drive improvements. We discuss examples where evaluation directly led to system enhancement. Incorporating evaluation in the development loop accelerates progress toward trustworthy QA systems suitable for critical medical use. However, current evaluations remain limited; many essential qualities are difficult to quantify, and existing benchmarks inadequately capture real-world complexity. There is substantial ongoing work needed to create more holistic and realistic evaluation frameworks, to keep pace with evolving models.

Limitations

This study specifically focuses on medical QA systems. During the literature review phase, we excluded publications related to general-domain large language models (LLMs) as well as healthcare-related literature not directly applicable to medical QA tasks.

Ethics Statement

We do not see any ethics issues in this paper.

References

Ian L Alberts, Lorenzo Mercolli, Thomas Pyka, George Prenosil, Kuangyu Shi, Axel Rominger, and Ali Afshar-Oromieh. 2023. Large language models (LLM) and ChatGPT: what will the impact on nuclear medicine be? *European journal of nuclear medicine and molecular imaging*, 50(6):1549–1552.

Manar Aljohani, Jun Hou, Sindhura Kommu, and Xuan Wang. 2025. [A comprehensive survey on the trustworthiness of large language models in healthcare](#). *Preprint*, arXiv:2502.15871.

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.

Dario Amodè, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*.

Ali Anaissi, Ali Braytee, and Junaid Akram. 2024. [Fine-Tuning LLMs for reliable medical question-answering services](#). *Preprint*, arXiv:2410.16088.

Elham Asgari, Nina Montaña-Brown, Magda Dubois, Saleh Khalil, Jasmine Balloch, Joshua Au Yeung, and Dominic Pimenta. 2025. A framework to assess clinical safety and hallucination rates of llms for medical text summarisation. *npj Digital Medicine*, 8(1):1–15.

Suhana Bedi, Yutong Liu, Lucy Orr-Ewing, Dev Dash, Sanmi Koyejo, Alison Callahan, Jason A Fries, Michael Wornow, Akshay Swaminathan, Lisa Soileymani Lehmann, and 1 others. 2024. Testing and evaluation of healthcare applications of large language models: A systematic review. *JAMA*.

Mohita Chowdhury, Ernest Lim, Aisling Higham, Rory McKinnon, Nikolett Ventoura, Yajie He, and Nick De Pennington. 2023. [Can large language models safely address patient questions following cataract surgery?](#) In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 131–137, Toronto, Canada. Association for Computational Linguistics.

Yun-Wei Chu, Kai Zhang, Christopher Malon, and Martin Renqiang Min. 2025. Reducing hallucinations of medical multimodal large language models with visual retrieval-augmented generation. *arXiv preprint arXiv:2502.15040*.

Julian De Freitas and I Glenn Cohen. 2024. The health risks of generative ai-based wellness apps. *Nature medicine*, 30(5):1269–1275.

Shrey Desai and Greg Durrett. 2020. [Calibration of pre-trained transformers](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 295–302, Online. Association for Computational Linguistics.

Yella Diekmann, Chase M Fensore, Rodrigo M Carrillo-Larco, Nishant Pradhan, Bhavya Appana, and Joyce C Ho. 2025. [Evaluating safety of large language models for patient-facing medical question answering](#). In *Proceedings of the 4th Machine Learning for Health Symposium*, volume 259 of *Proceedings of Machine Learning Research*, pages 267–290. PMLR.

Yanlin Feng, Xinyue Chen, Bill Yuchen Lin, Peifeng Wang, Jun Yan, and Xiang Ren. 2020. [Scalable multi-hop relational reasoning for knowledge-aware question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1295–1309, Online. Association for Computational Linguistics.

Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, 50(3):1097–1179.

Shreya Goyal, Sumanth Doddapaneni, Mitesh M Khapra, and Balaraman Ravindran. 2023. A survey of adversarial defenses and robustness in nlp. *ACM Computing Surveys*, 55(14s):1–39.

Maxime Griot, Coralie Hemptinne, Jean Vanderdonckt, and Demet Yuksel. 2025. Large language models lack essential metacognition for reliable medical reasoning. *Nature communications*, 16(1):642.

Tessa Han, Aounon Kumar, Chirag Agarwal, and Himabindu Lakkaraju. 2024a. [MedSafetyBench: Evaluating and improving the medical safety of large language models](#). In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Tessa Han, Aounon Kumar, Chirag Agarwal, and Himabindu Lakkaraju. 2024b. [Towards safe large language models for medicine](#). In *ICML 2024 Workshop on Models of Human Feedback for AI Alignment*.

Jie Huang and Kevin Chen-Chuan Chang. 2023. [Towards reasoning in large language models: A survey](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1049–1065, Toronto, Canada. Association for Computational Linguistics.

- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 others. 2025. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55.
- Xiaowei Huang, Wenjie Ruan, Wei Huang, Gaojie Jin, Yi Dong, Changshun Wu, Saddek Bensalem, Ronghui Mu, Yi Qi, Xingyu Zhao, and 1 others. 2024a. A survey of safety and trustworthiness of large language models through the lens of verification and validation. *Artificial Intelligence Review*, 57(7):175.
- Yue Huang, Lichao Sun, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bertie Vidgen, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, and 51 others. 2024b. [TrustLLM: Trustworthiness in large language models](#). *Preprint*, arXiv:2401.05561.
- Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. 2023. [Towards mitigating LLM hallucination via self reflection](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1827–1843, Singapore. Association for Computational Linguistics.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering qataset from medical exams. *Applied Sciences*, 11(14):6421.
- Adam Tauman Kalai and Santosh S. Vempala. 2024. [Calibrated language models must hallucinate](#). New York, NY, USA. Association for Computing Machinery.
- Yunsoo Kim, Jinge Wu, Yusuf Abdule, and Honghan Wu. 2024. [MedExQA: Medical question answering benchmark with multiple explanations](#). In *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, pages 167–181, Bangkok, Thailand. Association for Computational Linguistics.
- Anastasia Krithara, Anastasios Nentidis, Konstantinos Bougiatiotis, and Georgios Paliouras. 2023. BioASQ-QA: A manually curated corpus for biomedical question answering. *Scientific Data*, 10(1):170.
- Johannes A Landsheer. 2018. The clinical relevance of methods for handling inconclusive medical test results: quantification of uncertainty in medical decision-making and screening. *Diagnostics*, 8(2):32.
- Karim Lekadir, Alejandro F Frangi, Antonio R Porras, Ben Glocker, Celia Cintas, Curtis P Langlotz, Eva Weicken, Folkert W Asselbergs, Fred Prior, Gary S Collins, Georgios Kaissis, Gianna Tsakou, Irène Buvat, Jayashree Kalpathy-Cramer, John Mongan, Julia A Schnabel, Kaisar Kushibar, Katrine Riklund, Kostas Marias, and 30 others. 2025. [Future-ai: international consensus guideline for trustworthy and deployable artificial intelligence in healthcare](#). *BMJ*, 388.
- Yingji Li, Mengnan Du, Rui Song, Xin Wang, and Ying Wang. 2023a. A survey on fairness in large language models. *arXiv preprint arXiv:2308.10149*.
- Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023b. ChatDoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge. *Cureus*, 15(6).
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D Manning, Christopher Re, Diana Acosta-Navas, Drew A. Hudson, and 31 others. 2023. [Holistic evaluation of language models](#). *Transactions on Machine Learning Research*. Featured Certification, Expert Certification, Outstanding Certification.
- Fenglin Liu, Zheng Li, Hongjian Zhou, Qingyu Yin, Jingfeng Yang, Xianfeng Tang, Chen Luo, Ming Zeng, Haoming Jiang, Yifan Gao, and 1 others. 2024a. Large language models in the clinic: a comprehensive benchmark. *arXiv preprint arXiv:2405.00716*.
- Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. 2024b. [Trustworthy LLMs: a survey and guideline for evaluating large language models’ alignment](#). *Preprint*, arXiv:2308.05374.
- Hang Luo, Jian Zhang, and Chujun Li. 2025. [Causal graphs meet thoughts: Enhancing complex reasoning in graph-augmented llms](#). *Preprint*, arXiv:2501.14892.
- Atalanti Mastakouri, Elke Kirschbaum, Shiva Kavisvianathan, and Aaditya Ramdas. 2025. [QA-Calibration of language model confidence scores](#). In *Proceedings of the 13th International Conference on Learning Representations (ICLR 2025)*.
- Milad Moradi and Matthias Samwald. 2022. Improving the robustness and accuracy of biomedical language models through adversarial training. *Journal of Biomedical Informatics*, 132:104114.
- Nariman Naderi, Seyed Amir Ahmad Safavi-Naini, Thomas Savage, Zahra Atf, Peter Lewis, Girish Nadkarni, and Ali Soroush. 2025. [Self-Reported confidence of large language models in gastroenterology: Analysis of commercial, open-source, and quantized models](#). *Preprint*, arXiv:2503.18562.
- Zabir Al Nazi and Wei Peng. 2024. [Large language models in healthcare and medical domain: A review](#). *Informatics*, 11(3).

962	Robert Osazuwa Ness, Katie Matton, Hayden Helm,	Samuel Schmidgall, Carl Harris, Ime Essien, Daniel Ol-	1016
963	Sheng Zhang, Junaid Bajwa, Carey E Priebe, and	shvang, Tawsifur Rahman, Ji Woong Kim, Rojin Zi-	1017
964	Eric Horvitz. 2024. MedFuzz: Exploring the robust-	aei, Jason Eshraghian, Peter Abadir, and Rama Chel-	1018
965	ness of large language models in medical question	lappa. 2024a. Addressing cognitive bias in medical	1019
966	answering. <i>arXiv preprint arXiv:2406.06573</i> .	language models. <i>arXiv preprint arXiv:2402.08113</i> .	1020
967	Mahmud Omar, Benjamin S Glicksberg, Girish N Nad-	Samuel Schmidgall, Carl Harris, Ime Essien, Daniel	1021
968	karni, and Eyal Klang. 2024. Overconfident AI?	Olshvang, Tawsifur Rahman, Ji Woong Kim, Rojin	1022
969	benchmarking llm self-assessment in clinical scenar-	Ziaei, Jason Eshraghian, Peter Abadir, and Rama	1023
970	ios . <i>medRxiv</i> .	Chellappa. 2024b. Evaluation and mitigation of cog-	1024
971	Jesutofunmi A Omiye, Jenna C Lester, Simon Spichak,	nitive biases in medical language models. <i>npj Digital</i>	1025
972	Veronica Rotemberg, and Roxana Daneshjou.	<i>Medicine</i> , 7(1):295.	1026
973	2023. Large language models propagate race-based	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mah-	1027
974	medicine. <i>NPJ Digital Medicine</i> , 6(1):195.	davi, Jason Wei, Hyung Won Chung, Nathan Scales,	1028
975	Ankit Pal, Logesh Kumar Umapathi, and Malaikannan	Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl,	1029
976	Sankarasubbu. 2023. Med-HALT: Medical domain	and 1 others. 2023. Large language models encode	1030
977	hallucination test for large language models . In <i>Pro-</i>	clinical knowledge. <i>Nature</i> , 620(7972):172–180.	1031
978	<i>ceedings of the 27th Conference on Computational</i>	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres,	1032
979	<i>Natural Language Learning (CoNLL)</i> , pages 314–	Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin	1033
980	334, Singapore. Association for Computational Lin-	Clark, Stephen R Pfohl, Heather Cole-Lewis, and	1034
981	guistics.	1 others. 2025. Toward expert-level medical ques-	1035
982	Ateeb Ahmad Parray, Zuhra Mahfuza Inam, Diego Ra-	tion answering with large language models. <i>Nature</i>	1036
983	monfau, Shams Shabab Haider, Sabuj Kanti Mistry,	<i>Medicine</i> , pages 1–8.	1037
984	and Apurva Kumar Pandya. 2023. ChatGPT and	Hao Sun, Zhexin Zhang, Jiawen Deng, Jiale Cheng,	1038
985	global public health: applications, challenges, ethical	and Minlie Huang. 2023. Safety assessment of	1039
986	considerations and mitigation strategies.	chinese large language models. <i>arXiv preprint</i>	1040
987	Stephen R Pfohl, Heather Cole-Lewis, Rory Sayres,	<i>arXiv:2304.10436</i> .	1041
988	Darlene Neal, Mercy Asiedu, Awa Dieng, Nenad	Arun James Thirunavukarasu, Darren Shu Jeng Ting,	1042
989	Tomasev, Qazi Mamunur Rashid, Shekoofeh Azizi,	Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan,	1043
990	Negar Rostamzadeh, and 1 others. 2024. A toolbox	and Daniel Shu Wei Ting. 2023. Large language	1044
991	for surfacing health equity harms and biases in large	models in medicine. <i>Nature medicine</i> , 29(8):1930–	1045
992	language models. <i>Nature Medicine</i> , 30(12):3590–	1940.	1046
993	3600.	T Tran, V Joseph, L Smith, A Hopkins, S Lo, A Hen-	1047
994	Maria Paola Priola. 2024. Addressing hallucinations	nessy, C Juergens, H Dimitri, R Rajaratnam, J French,	1048
995	with RAG and NMISS in italian healthcare LLM	and 1 others. 2024. CardioCanon: A customised chat-	1049
996	Chatbots. <i>arXiv preprint arXiv:2412.04235</i> .	bot for cardiology inquiry with retrieval augmented	1050
997	Alireza Salemi and Hamed Zamani. 2024. Evaluating	generation to reduce hallucinations and improve per-	1051
998	retrieval quality in retrieval-augmented generation .	formance of large language models. <i>Heart, Lung and</i>	1052
999	In <i>Proceedings of the 47th International ACM SI-</i>	<i>Circulation</i> , 33:S379–S380.	1053
1000	<i>GIR Conference on Research and Development in</i>	Matthew Walsh, David Schulker, and Shing-hon Lau.	1054
1001	<i>Information Retrieval, SIGIR ’24</i> , page 2395–2400,	2024. Beyond capable: Accuracy, calibration, and ro-	1055
1002	New York, NY, USA. Association for Computing	bustness in large language models . Carnegie Mellon	1056
1003	Machinery.	University, Software Engineering Institute’s Insights	1057
1004	Thomas Savage, Stephen Ma, Abdessalem Boukil,	(blog). Accessed: 2025-May-16.	1058
1005	Vishwesh Patel, Ekanath Rangan, Ivan Lopez, and	Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru	1059
1006	Jonathan H Chen. 2024a. Fine tuning large lan-	Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao,	1060
1007	guage models for medicine: The role and impor-	Wenyang Gao, Xuming Hu, Zehan Qi, and 1 others.	1061
1008	tance of direct preference optimization . <i>Preprint</i> ,	2023. Survey on factuality in large language models:	1062
1009	arXiv:2409.12741.	Knowledge, retrieval and domain-specificity. <i>arXiv</i>	1063
1010	Thomas Savage, John Wang, Robert Gallo, Abdessalem	<i>preprint arXiv:2310.07521</i> .	1064
1011	Boukil, Vishwesh Patel, Seyed Amir Ahmad Safavi-	Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen,	1065
1012	Naini, Ali Soroush, and Jonathan H Chen. 2024b.	Lifan Yuan, Hao Peng, and Heng Ji. 2024. MINT:	1066
1013	Large language model uncertainty measurement	Evaluating LLMs in multi-turn interaction with tools	1067
1014	and calibration for medical diagnosis and treatment.	and language feedback . In <i>The Twelfth International</i>	1068
1015	<i>medRxiv</i> , pages 2024–06.	<i>Conference on Learning Representations</i> .	1069
		Jiaxin Wu, Yizhou Yu, and Hong-Yu Zhou. 2024. Uncer-	1070
		tainty estimation of large language models in medical	1071
		question answering . <i>Preprint</i> , arXiv:2407.08662.	1072

- Peng Xia, Ze Chen, Juanxi Tian, Gong Yangrui, Ruihou Hou, Yue Xu, Zhenbang Wu, Zhiyuan Fan, Yiyang Zhou, Kangyu Zhu, Wenhao Zheng, Zhaoyang Wang, Xiao Wang, Xuchao Zhang, Chetan Bansal, Marc Niethammer, Junzhou Huang, Hongtu Zhu, Yun Li, and 5 others. 2024. [CARES: A comprehensive benchmark of trustworthiness in medical vision language models](#). In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Rui Patrick Xian, Alex Jihun Lee, Satvik Lolla, Vincent Wang, Russell Ro, Qiming Cui, and Reza Abbasi-Asl. 2024. [Assessing biomedical knowledge robustness in large language models by query-efficient sampling attacks](#). *Transactions on Machine Learning Research*.
- Songhua Yang, Hanjie Zhao, Senbin Zhu, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan. 2024a. [Zhongjing: Enhancing the chinese medical capabilities of large language model through expert feedback and real-world multi-turn dialogue](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17):19368–19376.
- Yifan Yang, Qiao Jin, Furong Huang, and Zhiyong Lu. 2024b. Adversarial attacks on large language models in medicine. *ArXiv*, pages arXiv–2406.
- Junjie Ye, Yilong Wu, Songyang Gao, Caishuang Huang, Sixian Li, Guanyu Li, Xiaoran Fan, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024. [RoTBench: A multi-level benchmark for evaluating the robustness of large language models in tool learning](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 313–333, Miami, Florida, USA. Association for Computational Linguistics.
- Cyril Zakka, Rohan Shad, Akash Chaurasia, Alex R Dalal, Jennifer L Kim, Michael Moor, Robyn Fong, Curran Phillips, Kevin Alexander, Euan Ashley, and 1 others. 2024. Almanac—retrieval-augmented language models for clinical medicine. *Nejm ai*, 1(2):A10a2300068.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2023. [Automatic chain of thought prompting in large language models](#). In *The Eleventh International Conference on Learning Representations*.
- Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2024. Explainability for large language models: A survey. *ACM Transactions on Intelligent Systems and Technology*, 15(2):1–38.

A Appendix

Benchmark	Description	Key trustworthiness focus	Format
Med-HALT	Medical Hallucination Test dataset. Derived from medical exams across countries to probe factual recall and reasoning.	Hallucination: includes reasoning-based tests ("False Confidence", "None of the above" trick questions) and memory-based recall tests to quantify hallucination rates. Evaluates how often models produce unsupported info under stress-test conditions.	Multiple-Choice Questions, Yes/No, Open-ended question
MedHallBench	A comprehensive medical hallucination evaluation framework integrating automated clinical medical image caption hallucination scoring (ACHMI) and clinical expert review.	Hallucination: Design questions centered around object hallucinations, attribute hallucinations, multimodal conflict hallucinations, and logical reasoning hallucinations, and conduct adversarial tests to uncover the causes of hallucinations in models.	Open-ended Q&A, Visual Question Answering, Summarization
MedHallu	The first binary classification benchmark for medical hallucination detection. The questions are divided into three levels - Easy, Medium, and Hard - according to the difficulty of identifying hallucinations.	Hallucination: Detect whether the model can correctly classify the labels of question-answer pairs as "real" or "hallucination".	Binary Hallucination Detection
MedFuzz	By applying adversarial perturbations to medical question-answering queries, evaluate the robustness and performance of large language models (LLMs) in medical question-answering tasks.	Robustness: In the evaluation, first input the correct questions and answers into the model. Then, use the Attacker LLM to modify the original questions for multiple rounds and input them into the model. Each modification attempts to guide the target model to select the wrong answer without changing the correct answer of the original question.	Multiple-Choice Questions
BiasMedQA	A benchmark dataset for evaluating whether there is bias (towards different patient groups such as those of different genders, races, etc.) in LLMs in medical question answering.	Fairness: Introduce common clinically relevant cognitive biases into USMLE questions to test the performance of the model when facing these biases.	Multiple-Choice Questions
MedSafetyBench	The first medical-domain Safety evaluation benchmark dataset focused on assessing model responses to unsafe medical instructions.	Safety: Evaluate whether models can ensure response integrity when handling inputs containing unsafe medical instructions, as benchmarked by MedSafetyBench's adversarial testing framework.	Open-ended Q&A
MedExQA	Medical explainability QA benchmark. Covers 5 underrepresented specialties (e.g. speech pathology, clinical psych) with multiple ground-truth explanations per Q&A.	Explainability: evaluates if models can provide nuanced medical explanations beyond just correct answers. Uses lexical metrics and human ratings to score explanation quality. Also tests knowledge in less-studied specialties (robustness to specialty domains).	Open-ended question, required free-text explanation for answer.
PubMedQA	A Medical Reasoning Evaluation Benchmark for LLMs that Combine Expert-Annotated and Automated Knowledge Expansion, designed to assess contextual reasoning capabilities across medical texts and domain knowledge.	Reasoning: Given a question and a medical text context with the conclusion section removed, evaluate whether the model can infer if the question originally appeared in the conclusion section of the source text.	Three-way classification
DR.BENCH	A benchmark for evaluating clinical diagnostic reasoning capabilities of large language models (LLMs), comprising six reasoning tasks: MedNLI, Assessment and Plan Relation Labeling, EmrQA, SOAP Section Classification, Problem Summarization, and Diagnosis Generation.	Reasoning: The six diagnostic reasoning task categories in DR.BENCH comprehensively span the clinical workflow-continuum, designed to evaluate the model's capabilities including: medical concept logic; context-aware information retrieval; structured clinical knowledge classification; knowledge-graph-driven causal reasoning; multi-step evidence integration; knowledge-intensive clinical inference.	Multiple-Choice Questions, Extractive QA, Open-ended Questions, Text Generation
MedExpQA	MedExpQA encompasses multiple languages. For each question, a standard answer is provided along with multiple Gold-Explanation explanations written by medical experts.	Reasoning: Three types of tasks are set during evaluation: basic input only, basic input plus gold-standard explanation, and basic input plus RAG text. By comparing the outputs of the three types of tasks, the amount of missing reasoning ability of the model and the degree of help of automatically retrieved knowledge for the model's reasoning can be evaluated.	Multiple-Choice Questions
MediQ	A benchmark evaluating LLMs' capabilities in reliable interactive clinical reasoning, designed to assess their reasoning abilities by observing performance on informationally incomplete clinical queries.	Reasoning: Evaluating the simulation of a dynamic clinical interaction environment where the model under assessment acts as an Expert System, with performance under informationally incomplete initial conditions recorded to measure interactive clinical reasoning capabilities.	Multiple-Choice Questions, Interactive Q&A
MedXpertQA	A comprehensive benchmark for assessing expert-level medical knowledge and advanced reasoning capabilities, comprising Text (text-based) and MM (multimodal) subsets, with an independently designed reasoning subset.	Reasoning: The reasoning subset comprises highest-difficulty questions requiring multi-step logical reasoning, selected from both Text (text-based) and MM (multimodal) configurations, specifically designed to evaluate model reasoning capabilities	Multiple-Choice Questions, Multimodal QA

Table 1: Summary of representative benchmarks for each dimension, including their descriptions, key trustworthiness focus, and data format.