

Evaluating Medical LLMs by Levels of Autonomy: A Survey Moving from Benchmarks to Applications

Anonymous ACL submission

Abstract

Medical large language models achieve strong scores on standard benchmarks; however, transferring those results to safe and reliable performance in clinical workflows remains a challenge. This survey reframes evaluation through a levels-of-autonomy lens (L0–L3), spanning informational tools, information transformation and aggregation, decision support, and supervised agents. We align existing benchmarks and metrics with the actions permitted at each level and their associated risks, making the evaluation targets explicit. This motivates a level-conditioned blueprint for selecting metrics, assembling evidence, and reporting claims, alongside directions that link evaluation to oversight. By centering autonomy, the survey moves the field beyond score-based claims toward credible, risk-aware evidence for real clinical use.

1 Introduction

Large language models (LLMs) have advanced rapidly on medical benchmarks (Singhal et al., 2023; Hendrycks et al., 2020; Nazi and Peng, 2024). Both general-purpose (Qwen Team, 2025; DeepSeek-AI, 2025; OpenAI, 2025) and domain-specialized (Singhal et al., 2025; Bolton et al., 2024) LLMs now achieve high scores on licensing-style examinations and medical Q&A benchmarks (Oliveira et al., 2025; Tang et al., 2023). These headline results suggest that LLMs could provide meaningful assistance to clinicians and patients across a range of information-centric workflows.

However, benchmark correctness alone is not sufficient for clinical use (Hager et al., 2024; Ma et al., 2025). Clinical deployment requires consistency, fairness, auditability, calibrated uncertainty, and demonstrably safe clinical reasoning (Omiye et al., 2023; Fehr et al., 2024). Most benchmarks are Q&A-centric, which rarely probe these aspects, allowing unsafe reasoning and missing context to go undetected (Soroush et al., 2024; Asgari et al.,

2025). A recent systematic review of 519 evaluations of healthcare applications of LLMs shows that 44.5% evaluated medical knowledge (like exam questions) and 19.5% focused on diagnosis, but only 5% used real patient data. (Bedi et al., 2025b). A rigorous evaluation that spans factual grounding, reasoning quality, uncertainty calibration, safety, and human preferences is required (Tam et al., 2024; Shool et al., 2025).

This survey reviews the state of LLM evaluation in the medical field and identifies its limitations. We first summarize how LLMs are being applied in the medical field and what current benchmarks measure. While these scores are informative, they only provide quick capability snapshots and overlook integration into real workflows, calibration, and traceable provenance. Therefore, we move from scores to applications, treating evaluation as a means of showing that a system is sufficient for a defined purpose and scope at a specific autonomy level. For each level, we define the scope, typical applications, evaluation focus, boundaries, and challenges: **L0** Inform (no personalized advice); **L1** Information Transformation & Aggregation (structure, summarize); **L2** Decision Support (recommendations and personalized advice); and **L3** Agents Under Human Supervision (plan + invoke tools/APIs to enqueue actions under explicit review). This organization makes evaluation targets explicit: factual grounding at L0/L1; calibrated reasoning and coverage at L2; tool-use safety and auditability at L3. Regarding the challenges in each level, they are cumulative, not isolated: each higher autonomy level inherits unresolved issues from lower levels. As autonomy and permitted actions expand, new risks emerge that are specific to that level’s capabilities. Finally, we outline future work and recommendations for developing more robust evaluation frameworks to ensure that LLM-based systems can be trusted in clinical practice.

| | L0-Inform | L1-Transformation & Aggregation | L2-Decision Support | L3-Agents under Supervision |
|--------------------|---------------------------------|--|---|---|
| <i>Scope</i> | Explain / Inform; NO advice | Structured extraction and aggregation | Recommendation; NO action | Plan, act under supervision |
| <i>Application</i> | HealthSearchQA, Patient Summary | EHR Transformation, Retrieve & Aggregate | Diagnostic Reasoning, Patient Education | EHR co-pilot, Clinical Simulation |
| <i>Example</i> | What is MRI? | From this discharge note, extract the medication list. | Given the summary, is MRI appropriate? | Assemble prior results, draft an order, queue for sign-off. |
| <i>Challenge</i> | Hallucination, Bias | Attribution, Privacy | Reasoning Consistency, Fairness | Tool-use, Human-AI Interaction |

Table 1: Overview of autonomy levels (L0–L3) for medical LLMs, showing for each level the scope, typical applications, an example question/task, and key challenges. The rightward arrow indicates an increase in autonomy and risk. A detailed worked example/case study of an EHR copilot instantiated at different autonomy levels is provided in Appendix A.5. Example autonomy escalation criteria that specify when evidence at one level is sufficient to justify a move to the next are provided in Appendix A.3.

2 Related Work

Contemporary medical-LLM surveys mostly list datasets (e.g., USMLE/OKAP), task scenarios, and evaluation modes, but they rarely organize evaluation targets by autonomy level or permitted clinical actions. Representative examples include a review contrasting closed and open-ended tasks and discussing agentic settings (Chen et al., 2025c), and a systematic review highlighting reliance on general-purpose models and accuracy metrics with limited calibration and safety assessment (Shool et al., 2025). Parallel strands propose conversation-quality and safety evaluations (Abbasian et al., 2024), human-rater rubrics such as QUEST (Tam et al., 2024), multi-dimensional criteria in SCORE (Tan et al., 2024), and reporting guidance in TRIPOD-LLM (Gallifant et al., 2025). We address this gap by mapping evaluation objectives and metrics to autonomy levels (L0–L3). This mapping clarifies what evidence is sufficient for a system’s intended role, which risks must be tested and where human oversight is required. A detailed comparison with existing survey papers is summarized in Tables 5 and 6 in the Appendix.

Autonomy scales outside medical-LLM evaluation exist, but they target system design or human factors, rather than evaluation blueprints. In clinical decision research, levels of autonomy delineate who acts and who bears responsibility (Festor et al., 2021). In the broader agent literature, five-level design frameworks define autonomy via user roles

(operator → observer) (Feng et al., 2025a), and industry taxonomies describe L0–L5 agentic behavior (Kirkovska et al., 2025). None of these prescribe autonomy-conditioned metrics or align evaluation with healthcare oversight. Our survey fills that gap by integrating autonomy scales with concrete measurement choices for medical LLMs.

3 Evaluation Methodologies

We start by listing current benchmarks and metrics for medical LLMs. This sets clear measurement boundaries before we move to the applications. For LLMs in the medical field, benchmarks provide low-cost, repeatable evidence about specific capabilities within a bounded scope; they surface failure modes early, track progress over time, and help align claims with a system’s intended role (L0–L3). First, we group benchmarks by task and summarize the usefulness of each, as well as the L0–L3 level(s) they inform (§3.1). We then summarize both automated and human metrics, along with their limitations (§3.2). Table 2 provides a compact illustration of our taxonomy with a few representative rows; the full expanded table appears in Appendix A.1.

3.1 Benchmarks

Exam Q&A. Answer multiple-choice and short-answer questions from professional medical examinations across specialties. Typical benchmarks include MedQA (Jin et al., 2021) and MedMCQA (Pal et al., 2022). This task provides a low-cost

| Task | Usefulness | Level | Representative benchmarks | Example |
|--------------------------------|---|-------|-------------------------------------|---|
| Knowledge recall / exam Q&A | Low-cost unit tests of basic clinical knowledge. | L0 | MedMCQA; MultiMedQA | Is it safe to take ibuprofen in early pregnancy? |
| Summarization / transformation | Checks whether models produce faithful summaries. | L1 | MIMIC-IV-BHC; MTS-Dialog; ACI-Bench | Condense a multi-paragraph ICU discharge note into a 2–3 sentence brief hospital course |
| Decision support | Probes decision-making in simulated cases. | L2 | HealthBench; MedHELM; DexBench | Given 7 days of CGM traces and meal logs, generate some advice. |

Table 2: Condensed illustration of benchmark classes; the full expanded version appears in Appendix A.1

way to test factual breadth and specialty coverage. It is useful as a unit test for L0–L1 informational tools and for identifying coarse knowledge gaps.

Summarization. Generate faithful summaries, patient-friendly explanations, or structural rewrites of clinical text. Typical benchmarks include MS² (Wang et al., 2022) and PubMed long-document summarization (Cohan et al., 2018). This task evaluates whether outputs are faithful to the source (no important omissions) and structurally complete for explanation-oriented use cases at L0–L1.

Retrieval-augmented generation (RAG). Answer questions grounded in retrieved documents with explicit attribution to evidence passages. Typical benchmarks include BEIR (Thakur et al., 2021) and TREC-COVID (Roberts et al., 2021). This task demonstrates whether responses are supported by cited sources and avoid contradicting them, enabling checks on evolving topics at L0–L1.

Information extraction. Extract and normalize clinical entities, relations, and codes from notes and reports. Typical benchmarks include CBLUE (Zhang et al., 2021). This task establishes reliable structuring and normalization to support downstream aggregation in L1 systems.

Decision support. Make thresholded recommendations or triage decisions in vignettes or multi-turn clinical scenarios scored by clinician rubrics. Typical benchmarks include HealthBench (Arora et al., 2025) and MedHELM (Bedi et al., 2025a). This task targets selective reliability for L2 systems by quantifying calibrated behavior at thresholds and surfacing harm-proximal errors in simulation. Despite strong exam-style performance, overall success rates on HealthBench remain modest: GPT-3.5-Turbo satisfies only about 16% of physician-authored behavioral criteria, GPT-4o around 32%, and OpenAI’s o3 roughly 60%, with more challenging “HealthBench Hard” subset capped at about

32% even for the best model.

Clinical dialogue. Conduct conversations with patients or clinicians to achieve goals while communicating safely and clearly. Typical benchmarks include MedDialog (Zeng et al., 2020) and MedDG (Liu et al., 2022). This task assesses communication effectiveness and safety across L0–L2.

3.2 Automated vs. Human Metrics

Automated Metrics. Automated metrics quantify answer correctness, calibration, faithfulness, and retrieval quality. For Q&A, Exact Match and F1 summarize accuracy, while ECE, Brier, and NLL capture probability calibration (Guo et al., 2017; Brier, 1950; Manning et al., 2008). In selective prediction, risk-coverage curves relate error to the answered fraction and justify abstention on uncertain cases (Geifman and El-Yaniv, 2017; Geifman and El-Yaniv, 2019; Traub et al., 2024). For summarization, lexical overlap (ROUGE/chrF) and embedding-based similarity (BERTScore/BLEURT) are commonly complemented with source-grounded checks for omissions and contradictions (Lin, 2004; Popović, 2015; Zhang et al., 2020; Sellam et al., 2020). In RAG settings, retrieval ranking is assessed with Recall@k, MRR, and nDCG (Manning et al., 2008).

Human Evaluation. Clinician raters typically apply behaviorally anchored rubrics (BARS-style) to score clinical correctness, coverage, contextualization, reasoning transparency, uncertainty handling, readability, actionability/safety, and empathy; inter-rater agreement is summarized with Cohen’s κ or ICC (Holland et al., 2022; McHugh, 2012). For patient-facing text, readability is assessed using the Flesch-Kincaid or SMOG indices (Singh et al., 2024; Badarudeen and Sabharwal, 2010).

LLM-as-Judge (LAJ). To scale open-ended assessments, many studies adopt LLM-as-judge:

strong models render pairwise preferences or rubric-based scores, often with prompts that request quoted evidence. To control bias, studies randomize order or blind positions and, when feasible, ensemble multiple judges. (Zheng et al., 2023; Li et al., 2024c; Gu et al., 2024; Zhu et al., 2025; Tan et al., 2025). LAJ can track human preferences well in aggregate but is vulnerable to position and verbosity biases and self-preference effects (Chen et al., 2024; Ye et al., 2025a; Wataoka et al., 2024).

4 Applications and Autonomy Levels

From Scores to Applications. Static benchmarks provide an overview of the model capabilities, but they do not capture real workflow context, calibration, provenance, and audit trails, or role-specific needs. We therefore treat evaluation as showing that a system is **sufficient for a defined purpose and scope** at a given autonomy level. Concretely, each autonomy level is organized into five parts: (1) **Definition & scope**; (2) **Typical applications**; (3) **Evaluation focus**; (4) **Scope boundaries**; and (5) **Challenges**. Benchmarks remain useful as ingredients in this assessment, not the destination.

4.1 L0 - Inform

Definition and Scope. At autonomy level L0, the system functions purely as an informational tool: it explains medical concepts and provides a general background in plain language. It neither tailors advice to an individual patient nor initiates clinical decisions. Outputs are educational in nature and include an explicit non-advice disclaimer.

Typical Applications. Representative L0 tasks include answering common health questions (e.g., “What is MRI?”), producing lay summaries of technical passages, and simplifying lab reports. Public datasets used as proxies include consumer Q&A corpora and patient-facing summarization sets: HealthSearchQA within the MultiMedQA (Singhal et al., 2023) suite, the TREC 2017 LiveQA Medical task (Abacha et al., 2017), PubMedQA for abstract-level research comprehension, and consumer-question summarization datasets such as MeQSum and MEDIQA’21 (Ben Abacha and Demner-Fushman, 2019; Ben Abacha et al., 2021).

Evaluation Focus. L0 evaluation should primarily focus on accuracy, completeness, and readability (Srinivasan et al., 2025). Some benchmarks also include structured rubrics or provide reference contexts that enable additional checks: HealthBench

(Arora et al., 2025) grades free-text answers with physician-written rubrics, MultiMedQA reports human ratings for long consumer answers, and PubMedQA uses the abstract as the reference context. We reference these only to motivate the axes here, not to require source citation as a core L0 practice. **Scope Boundaries.** L0 outputs are intentionally non-patient-specific and involve minimal reasoning: they recall and lightly synthesize facts but do not conduct case workups, triage, or make recommendations (those belong to L2 Decision Support). As a result, L0 content may omit person-specific contraindications or time-sensitive context and should avoid language that could be interpreted as advice. Readability should target plain-language norms to reduce misunderstanding, and common safety disclaimers may be shown without tailoring. Some L0 benchmarks include rubrics or reference contexts for scoring, but we do not treat source citation as a core requirement for L0.

Challenges. Hallucination: Even at L0, models often produce fluent but unfaithful text: Broad surveys separate hallucinations into intrinsic and extrinsic, link them to training/decoding choices and weak grounding, and recommend source- and task-aware evaluations instead of generic overlap scores (Ji et al., 2023). In summarization specifically, human studies show that systems invent unsupported details and that standard n-gram metrics miss these errors, motivating faithfulness-oriented checks instead (Maynez et al., 2020). For patient-facing summaries, clinical audits advise using medical rubrics that check whether each claim is supported by the underlying notes and whether important facts are missing—rather than relying on readability alone. (Asgari et al., 2025). However, because these tools audit rather than eliminate hallucinations, the problem persists. Recent clinical summarization audits quantify how often hallucinations occur in practice. In psychiatric discharge summaries, human-written summaries were rated higher than GPT-4 summaries (3.78 vs 3.12 on a 5-point scale), and hallucinations appeared in 40% of AI summaries (Schmidgall et al., 2024). In an emergency-department study of 100 encounters, only 33% of GPT-4 summaries were error-free, while 42% contained hallucinations and 47% omitted clinically relevant information (Williams et al., 2025). **Bias:** We use bias to mean systematic tendencies that push outputs away from truth or intended scope (e.g., sensitivity to superficial cues), such that the same input intent

can yield different explanations. These tendencies plausibly arise from the pretraining objective and are further shaped by alignment procedures that optimize against human preferences (Ouyang et al., 2022; Sharma et al., 2023; Xu et al., 2025b; Shen et al., 2025). In L0 applications, such biases surface as assistants mirroring user beliefs rather than correcting them (Sharma et al., 2023) and shifting outputs under small, meaning-preserving prompt changes such as formatting tweaks, option order in multiple-choice settings, or early anchoring hints (Sclar et al., 2023; Pezeshkpour and Hruschka, 2023; Lou et al., 2024; Zhou et al., 2024; Li et al., 2024b; Ye et al., 2025b; Li et al., 2024a; RRV et al., 2025). They persist in practice because preference-optimized objectives and prompt conventions are integral to how LLMs are used (Kadavath et al., 2022; Steyvers et al., 2025). Consequently, mitigations such as reporting performance ranges across prompt formats, anchor-aware templates, and a conservative tone can reduce bias at L0, but they do not reliably remove it. (Sclar et al., 2023; Pezeshkpour and Hruschka, 2023; Lou et al., 2024; Steyvers et al., 2025).

4.2 L1 - Information Transformation & Aggregation

Definition and Scope. This stage turns raw, heterogeneous clinical data into standardized, computable representations and then combines them with external evidence to produce grounded outputs. In practice, health systems map local EHR fields to an interoperability standard or a research schema, attach machine-readable provenance for auditability, and operate over de-identified corpora such as MIMIC-IV that illustrate the target tables (encounters, labs, medications) and privacy constraints (Abacha et al., 2021; Alsentzer et al., 2023; Zhang et al., 2024a; Liu et al., 2021).

Typical Applications. EHR data transformation: Typical pipelines mix schema harmonization (FHIR/OMOP ETL), clinical NLP to extract entities/attributes from notes, concept normalization to standard vocabularies, and de-identification. i2b2/n2c2 shared tasks supply widely used de-identified note sets for de-identification, concept extraction, relation labeling, and medication-change context, enabling objective measurement of span-level and mapping accuracy (Li et al., 2016; Abacha et al., 2021; Nowak et al., 2023; Mahajan et al., 2023; Henry et al., 2020). **Retrieve & Aggregate:** On top of the transformed corpus, sys-

tems index structured facts (problem lists, meds, labs) and unstructured notes, then pair them with external sources (guidelines, reviews) via retrieval-augmented generation (RAG). Retrieval metrics (Recall@k, nDCG, MRR) assess whether the right evidence is fetched; generation metrics (faithfulness/attribution, grounded-answer rate) check that outputs rely on retrieved passages rather than model priors (Cohan et al., 2020; Zhang et al., 2024a; Tang et al., 2023). Representative resources include BEIR for generalizable retrieval evaluation, TREC-COVID for high-stakes, rapidly evolving topics, and domain-specific retrievers such as MedCPT for biomedical search (Thakur et al., 2021; Roberts et al., 2021; Jin et al., 2023).

Evaluation Focus. Transformation: report extraction/normalization scores (e.g., span-level F1, concept mapping accuracy), coverage/completeness of key fields, and lineage completeness. **Retrieval & Aggregation:** report Recall@k/nDCG/MRR; grounded-answer and attribution rates to retrieved passages; contradiction-to-source; and selective prediction/abstention rates under uncertainty (Alsentzer et al., 2023; Abacha et al., 2021; Zhang et al., 2024a).

Scope Boundaries. L1 improves structure, traceability, and access to evidence, but it does not perform the patient-specific reasoning required for diagnosis, test selection, or treatment trade-offs. Real clinical workups must integrate temporality, comorbidities, contraindications, and uncertainty-capabilities not captured by schema or Recall@k alone. Studies also show that (i) retrieval can surface conflicting or outdated sources and (ii) LLMs may misattribute or over-trust citations; recent medical RAG evaluations highlight these gaps even when retrieval quality is strong (Roberts et al., 2021; Peng et al., 2023; Hueber and Kleyer, 2023; Golan et al., 2023; Dhanvijay et al., 2023; Sezgin et al., 2023; Team et al., 2024). Thus, L0+L1 should be paired with higher-level (L2+) decision-focused assessments before deployment.

Challenges. Attribution: Attribution concerns whether aggregated outputs actually rely on and are traceable to the retrieved sources. Automated audits in the medical domain report that models often cite papers that are only loosely relevant or that do not fully support the generated claims (Wu et al., 2025). Even with explicit attribution metrics and fine-grained factuality scoring, models can pass retrieval tests while still incorporating priors or blending multiple sources in ways that obscure

provenance (Yang et al., 2025). Techniques like Self-RAG improve on-demand retrieval and self-critique, yet do not eliminate misattribution when sources disagree, are low quality, or when prompts nudge the model toward fluent synthesis over faithful quotation (Asai et al., 2023; Jung et al., 2024). Thus, attribution still remains a challenge. **Completeness:** Completeness addresses whether transformed corpora and their aggregations cover all clinical facts (problems, meds, labs, temporality, context) without omissions. Comparative studies find that concept-recognition tools are inconsistent and often miss negations, misread abbreviations, and struggle with ambiguity or misspellings, which leads to missing or distorted facts downstream (Lossio-Ventura et al., 2023). Beyond extraction, mapping text spans to standard concepts is fragile because benchmark datasets contain many ambiguous terms and do not align well with UMLS coverage, so the ‘correct’ code is often unclear from the beginning (Newman-Griffis et al., 2021). In practice, long EHR narratives, events spread across notes, and uneven local coding leave gaps that retrieval can’t fill when the structured substrate is incomplete. These issues persist because gold standards underrepresent edge cases, annotation guidelines vary, and many L1 evaluations emphasize span-level F1 or Recall@k over end-to-end coverage of clinically critical fields (Lossio-Ventura et al., 2023). **Privacy & lineage:** It evaluates whether L1 transformations are governable and safe to share. De-identification of clinical notes reduces direct identifiers but does not preclude membership inference against downstream models (Sarkar et al., 2024). At the model layer, training data extraction and related attacks demonstrate that LLMs can memorize and regurgitate snippets of their training corpora (Carlini et al., 2021). These realities make rigorous lineage essential: machine-readable provenance (i.e., what data, which ETL/normalizers) should be recorded using established schemas, allowing organizations to audit and reproduce outputs (Lebo et al., 2013; Mitchell et al., 2019; Gebru et al., 2021). Yet provenance remains challenging in practice because pipelines are multi-hop and frequently updated; components are swapped or fine-tuned; and evidence bases evolve.

4.3 L2 - Decision Support

Definition and Scope. At autonomy level L2, the system provides patient-specific recommendations

that can assist clinical decision-making. By design, L2 depends on upstream EHR transformation and retrieval/aggregation (L1) but adds reasoning over the individual’s data and clinical context.

Typical Applications. Diagnostic reasoning: it reads a patient’s history, exam, labs, and imaging to propose differentials with brief rationales and possible next steps; evaluations commonly use vignette-based case sets and prompting schemes that elicit stepwise clinical reasoning (Goh and colleagues, 2024; Savage et al., 2024a). **Medication decision support:** it integrates a patient’s active medications, allergies, problems, and recent labs to surface potential adverse drug events, drug–drug interactions, and dosing/contraindication checks; widely used resources include the 2018 n2c2 ADE shared task and the SemEval-2013 DDIE extraction task for literature-based DDI detection (Henry et al., 2020; Segura-Bedmar et al., 2013). **Patient education:** Patient education tools support patient–clinician communication by turning medical information into clear, usable messages and scaffolding two-way conversations. Typical functions include generating plain-language explanations of conditions, tests, and procedures; prompting patients to ask key questions; and using teach-back so patients restate key points to confirm understanding (Shoemaker et al., 2013; Centers for Disease Control and Prevention, 2020; Agency for Healthcare Research and Quality, 2023; Institute for Healthcare Improvement, n.d.; IPDAS Collaboration, 2024; Stacey et al., 2021). Taken together, these tasks are L2 because they require reasoning over an individual patient’s context to generate recommendations or tailored explanations.

Evaluation Focus. For L2 decision support, evaluation should foreground three axes beyond L0’s accuracy/completeness/readability: reasoning consistency, calibration and abstention, and safety. For reasoning, score the process, not just final answers: stepwise diagnostic logic and evidence use should be judged against clinician-authored rubrics or case rationales. For calibration, report reliability at the case level, plus selective prediction curves (risk–coverage) with a tunable “don’t know / escalate” option. For safety, it is essential to track rates of contraindications and guideline violations.

Scope Boundaries. At L2, systems provide decision support by combining a patient’s EHR context with external evidence to offer guidance. However, they do not plan tasks, call tools or APIs, or change the record, so a clinician must review and decide.

This preserves clinical authority but requires clinicians to convert text into orders and messages, which adds to their workload and increases the risk of transcription errors. These limits motivate a shift to agentic L3 configurations that keep a human in the loop while reducing cognitive burden.

Challenges. Reasoning consistency & faithfulness: It concerns whether patient-specific answers are stable across prompt phrasings and whether the rationales actually support the recommendation. Comparative guideline evaluations demonstrate that modifying formats or instructions can significantly impact model outputs and agreement, underscoring the sensitivity of prompts in medical settings (Wang et al., 2024). Beyond prompt sensitivity, in a study of seven LLMs on 2,400 abdominal-pain cases, clinicians achieved around 90% diagnostic accuracy, whereas model performance was on average 16–25 percentage points lower and more frequently violated guideline-concordant workup recommendations (Hager et al., 2024). Even when stepwise prompts are used to elicit reasoning, studies find that the generated “explanations” can be unfaithful (Turpin et al., 2023; Madсен et al., 2024; Kuang et al., 2025; RRV et al., 2025). In diagnostic contexts, structured prompting can improve transparency but does not guarantee faithful causal grounding of the final answer (Savage et al., 2024b; Dineen et al., 2025). The result in practice is that two seemingly careful L2 prompts may yield different plans with rationales that read well but do not reliably reflect the model’s decision process, which keeps consistency and faithfulness an open problem for deployment (Turpin et al., 2023). **Confidence calibration:** Here the question is whether stated or implicit confidence tracks correctness so that uncertain cases can be flagged or deferred. Recent medical evaluations show that simple proxies correlate only weakly with error (Bentegeac et al., 2025). Semantic-entropy approaches detect confabulations and better prioritize abstention, but their reliability varies by task and domain (Farquhar et al., 2024; Kossen et al., 2024; Penny-Dimri et al., 2025). In real L2 use, these trade-offs mean systems may sound certain on incorrect recommendations or abstain too rarely on edge cases (Bentegeac et al., 2025; Farquhar et al., 2024; Feng et al., 2025b). **Fairness:** The core issue is whether recommendations generalize equitably across patient subgroups and clinical contexts. Specialized benchmarks and audit tools reveal that biases are present in long-form answers and clinical

recommendations, which plain accuracy scores fail to capture (Prakash et al., 2024). Purpose-built bias benchmarks for clinical LLMs report subgroup-linked shifts in outputs, while triage studies using counterfactual tests reveal intersectional differences across sex and race (Zhang et al., 2024b; Lee et al., 2025). A recent evaluation in healthcare showed that changing a patient’s race or gender caused GPT-4 to miss the top correct diagnosis in 37% of cases. (Zack et al., 2024) A recent scoping review highlights uneven coverage across medical fields and limited clinician-in-the-loop evaluation-leaving blind spots in where and how biases manifest (Cheng et al., 2025). Consequently, subgroup reliability remains a persistent challenge (Prakash et al., 2024; Cheng et al., 2025).

4.4 L3 - Agents under human supervision

Definition and Scope. We define L3 agents as systems that plan and invoke tools to initiate actions in clinical workflows while keeping a clinician explicitly “in the loop” for review, modification, and sign-off. Core capabilities are task planning, retrieval, and safe tool use, with human oversight enforced at key checkpoints (e.g., order “draft” states).

Typical Applications. Clinical Copilot: Clinical copilots plan tasks, fetch chart context, and then draft orders, messages, referrals, or care plans via tool calls, pausing for human sign-off (e.g., Polaris, Rx Strategist, Almanac Copilot) (Mukherjee et al., 2024; Van et al., 2024; Zakka et al., 2024). **Sandboxed Simulation:** It embeds agents in controlled clinics to probe plan–act–check loops with audit trails-benchmarks like AgentClinic, AI Hospital, ClinicalLab, and 3MDBench (Schmidgall et al., 2024; Fan et al., 2024; Yan et al., 2024; Sviridov et al., 2025; Yue et al., 2025; Chen et al., 2025a; Xu et al., 2025a; Jiang et al., 2025). **Operation & EHR Automation:** It coordinates non-diagnostic workflows (e.g., prior authorization) and return proposed actions for approval (RxLens; multi-agent prior-auth pipelines) (Jagatap et al., 2025). In contrast to L2 (textual recommendations), these L3 systems invoke tools to initiate actions but keep a clinician in the loop for review and sign-off.

Evaluation Focus. At L3, evaluation shifts from answer quality to supervised action quality: studies should first demonstrate end-to-end task success under human oversight. Second, they must verify tool-use correctness: every API call, order, and parameter matches clinical intent. Finally, they should require auditability via machine-readable

630 provenance of data, prompts, models, tools, param- 682
631 eters, and approvals, leveraging standards such as 683
632 FHIR Provenance and AuditEvent so that actions 684
633 can be traced for post-hoc review and governance.
634 **Scope Boundaries.** L3 agents are limited to draft- 685
635 and-queue actions under explicit human oversight: 686
636 they may plan tasks and call approved tools/APIs, 687
637 but execution requires a clinician’s independent 688
638 review. Looking ahead, closed-loop system deploy- 689
639 ments can enable agents to execute pre-approved, 690
640 low-risk steps when explicit policies and fine- 691
641 grained access controls are met, rather than requir- 692
642 ing individual approval for each action.
643 **Challenges. Tool-use failures:** These arise when 693
644 an agent plans correctly but issues the wrong API 694
645 call or constructs malformed parameters, so the 695
646 action fails even if the reasoning was sound. In 696
647 one benchmark of 34 complex tasks, only 50% of 697
648 tasks were successfully completed, with failures at- 698
649 tributed to improper planning, non-functional code 699
650 generation, or inadequate self-correction loops (Lu 700
651 et al., 2025). **Human–AI interaction dynamics:** 701
652 These refer to how clinicians review and sign off on 702
653 queued actions. Studies of decision support reveal 703
654 an overreliance on automated suggestions and alert 704
655 fatigue, where users either ignore or overaccept sys- 705
656 tem outputs, resulting in oversight gaps even when 706
657 accuracy is reasonable (Khera et al., 2023; Abdel- 707
658 wanis et al., 2024). The problem persists because 708
659 busy workflows, lengthy sessions, and inconsis- 709
660 tent interface cues make it challenging to calibrate 710
661 attention and sustain critical review at every step.
662 **Auditability and accountability:** These require a 711
663 traceable record of what data were accessed, which 712
664 tools were called with what parameters, and who 713
665 approved the final action. Standards such as FHIR 714
666 AuditEvent and Provenance define the necessary 715
667 primitives, and governance frameworks emphasize 716
668 logging and traceability for post-hoc review (HL7 717
669 International, 2025a,b; National Institute of Stan- 718
670 dards and Technology, 2023; World Health Organi- 719
671 zation, 2025). Yet multi-tool, multi-service agent 720
672 stacks often lack end-to-end lineage across steps, 721
673 so reconstructing a failure or near-miss remains dif- 722
674 ficult in practice. **Scalability of training:** Medical 723
675 agents cannot typically be trained directly in live 724
676 care settings because unconstrained exploration is 725
677 ethically and logistically infeasible, pushing work 726
678 toward simulated environments (Schmidgall et al., 727
679 2024; Fan et al., 2024). However, strong perfor- 728
680 mance in such synthetic setups does not guarantee 729
681 reliable transfer to real-world workflows: bench-

marks often simplify states and constraints, leaving 682
significant gaps between simulation and reality (Ja- 683
yaraman et al., 2024; Hargrave et al., 2024). 684

5 Future Work 685

Closed-looped System We expect clinical deploy- 686
ments to shift from single-model helpers to closed- 687
loop, hospital-scale systems composed of cooper- 688
ating, role-specialized agents that escalate to clini- 689
cians at predefined gates (Schmidgall et al., 2024; 690
Borkowski and Ben-Ari, 2025). We call for simple, 691
auditable protocols for handoff, disagreement reso- 692
lution, and safety gating of tool use, preceded by 693
simulation and shadow deployments. We recom- 694
mend reporting operational outcomes (e.g., deferral 695
rates) so the community can compare architectures 696
and converge on safe patterns (Liu et al., 2020). 697

Guarantees. We call for reframing “good” perfor- 698
mance around risk-controlled selectivity: in clini- 699
cal settings, systems should act only when a target 700
risk can be met and otherwise defer, evaluated by 701
risk–coverage rather than accuracy. In this agenda, 702
reliable high accuracy on a subset with explicit 703
deferral (e.g., 99% accuracy on 20% task) is prefer- 704
able to broad moderate accuracy (e.g., 80% on 705
80%), because the former enables out-of-scope de- 706
tection and smooth human routing while the latter 707
obscures which cases are incorrect. We recommend 708
that future evaluations report calibrated confidence 709
on the acted-on subset and coverage at target risk, 710
along with slice breakdowns (Guo et al., 2017). 711

6 Conclusion 712

This survey reframes the evaluation of medical 713
LLMs around levels of autonomy (L0–L3). We 714
emphasize risk coverage rather than average accu- 715
racy: safe systems act only when they can meet 716
a target risk and defer otherwise. At L0–L1, the 717
focus is factual fidelity, bias, and structural correct- 718
ness with clear grounding to sources. At L2–L3, 719
the bar rises to uncertainty, subgroup robustness, 720
tool-use correctness, and verified human oversight. 721

Taken together, this level-conditioned lens turns 722
benchmark scores into clinically relevant claims. 723
Evaluations are most persuasive when they make 724
the level explicit, pair target risk with achieved cov- 725
erage, report performance across levels, and verify 726
both confidence calibration and oversight check- 727
points. We hope this provides a clear language for 728
building medical LLMs that are not only capable 729
but also reliably useful and safe in practice. 730

731 Limitations

732 This survey is necessarily selective and time-
733 bounded; model releases, datasets, and guidance
734 evolve rapidly, so some details may become out-
735 dated. Evidence in the literature remains un-
736 even: many studies emphasize lab benchmarks over
737 prospective or randomized evaluations, and report-
738 ing quality is inconsistent.

739 References

- 740 Asad Aali, Dave Van Veen, Yamin Arefeen, Jason
741 Hom, Christian Bluethgen, Eduardo Pontes Reis, Ser-
742 gios Gatidis, Namuun Clifford, Joseph Daws, Arash
743 Tehrani, Jangwon Kim, and Akshay Chaudhari. 2025.
744 [Mimic-iv-ext-bhc: Labeled clinical notes dataset for
745 hospital course summarization](#). PhysioNet Reposi-
746 tory.
- 747 A. B. Abacha, E. Agichtein, Y. Pinter, and D. Demner-
748 Fushman. 2017. Overview of the medical question
749 answering task at trec 2017 liveqa. *TREC*, pages
750 1–12.
- 751 A. B. Abacha, Y. M'rabet, Y. Zhang, C. Shivade, C. Lan-
752 glotz, and D. Demner-Fushman. 2021. Overview of
753 the mediqa 2021 shared task on summarization in the
754 medical domain. *Proceedings of the 20th Workshop
755 on Biomedical Language Processing*, pages 74–85.
- 756 Mahyar Abbasian, Elahe Khatibi, Iman Azimi, David
757 Oniani, Zahra Shakeri Hossein Abad, Alexander
758 Thieme, Ram Sriram, Zhongqi Yang, Yanshan Wang,
759 Bryant Lin, Olivier Gevaert, Li-Jia Li, Ramesh Jain,
760 and Amir M. Rahmani. 2024. [Foundation metrics
761 for evaluating effectiveness of healthcare conversa-
762 tions powered by generative ai](#). *npj Digital Medicine*,
763 7(1):82.
- 764 Moustafa Abdelwanis, Hamdan Khalaf Alarafati,
765 Maram M. S. Tammam, and Mecit Can Emre Sim-
766 sekler. 2024. [Exploring the risks of automation bias
767 in healthcare artificial intelligence applications: A
768 bowtie analysis](#). *Journal of Safety Science and Re-
769 siliance*, 5:460–469.
- 770 Agency for Healthcare Research and Quality. 2023.
771 [Tool: Teach-back](#).
- 772 Ahmed M. Alkalbani, Ahmed S. Alrawahi, Ahmed
773 Salah, Vafa Haghghi, Yong Zhang, Saif Alkindi,
774 and Quan Z. Sheng. 2025. [A systematic review of
775 large language models in medical specialties: Appli-
776 cations, challenges and future directions](#). *Informa-
777 tion*, 16(6):489.
- 778 E. Alsentzer, M. J. Rasmussen, R. Fontoura, A. L. Cull,
779 B. Beaulieu-Jones, K. J. Gray, D. W. Bates, and V. P.
780 Kovacheva. 2023. Zero-shot interpretable phenotyp-
781 ing of postpartum hemorrhage using large language
782 models. *NPJ Digital Medicine*, 6(1).

- Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Pre-
ston Bowman, Joaquin Quiñero-Candela, Foivos
Tsimpourlas, Michael Sharman, Meghan Shah, An-
drea Vallone, Alex Beutel, and 1 others. 2025.
Healthbench: Evaluating large language models
towards improved human health. *arXiv preprint
arXiv:2505.08775*. 783
784
785
786
787
788
789
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and
Hannaneh Hajishirzi. 2023. [Self-RAG: Learning to
retrieve, generate, and critique through self-reflection](#).
arXiv preprint arXiv:2310.11511. 790
791
792
793
- 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):274. 794
795
796
797
798
799
- Sameer Badarudeen and Sanjeev Sabharwal. 2010. [As-
sessing readability of patient education materials:
Current role in orthopaedics](#). *Clinical Orthopaedics
and Related Research*, 468(10):2572–2580. 800
801
802
803
- Suhana Bedi, Hejie Cui, Miguel Fuentes, Alyssa Unell,
Michael Wornow, Juan M Banda, Nikesh Kotecha,
Timothy Keyes, Yifan Mai, Mert Oez, and 1 oth-
ers. 2025a. [Medhelm: Holistic evaluation of large
language models for medical tasks](#). *arXiv preprint
arXiv:2505.23802*. 804
805
806
807
808
809
- Suhana Bedi, Yutong Liu, Lucy Orr-Ewing, Dev Dash,
Sanmi Koyejo, Alison Callahan, Jason A. Fries,
Michael Wornow, Akshay Swaminathan, Lisa So-
leymani Lehmann, and 1 others. 2025b. [Testing
and evaluation of health care applications of large
language models: A systematic review](#). *JAMA*,
333(4):319–328. 810
811
812
813
814
815
816
- Asma Ben Abacha and Dina Demner-Fushman. 2019.
[On the summarization of consumer health questions](#).
In *Proceedings of ACL*, pages 2228–2234. 817
818
819
- Asma Ben Abacha, Dina Demner-Fushman, Shweta
Yadav, and Deepak Gupta. 2021. [MEDIQA 2021:
Consumer health question summarization, multi-
answer summarization, and radiology report sum-
marization](#). [https://sites.google.com/view/
mediqa2021](https://sites.google.com/view/mediqa2021). NAACL-BioNLP shared task site. 820
821
822
823
824
825
- Asma Ben Abacha, Wen-wai Yim, Yadan Fan, and
Thomas Lin. 2023. [An empirical study of clinical
note generation from doctor–patient encounters](#). In
*Proceedings of the 17th Conference of the European
Chapter of the Association for Computational Lin-
guistics (EACL)*. 826
827
828
829
830
831
- Raphaël Bentegeac, Bastien Le Guellec, Grégory
Kuchcinski, Philippe Amouyel, and Aghiles Ham-
roun. 2025. [Token probabilities to mitigate large
language models overconfidence in answering medi-
cal questions: Quantitative study](#). *Journal of Medical
Internet Research*, 27:e64348. 832
833
834
835
836
837

| | | | |
|-----|--|-----|--|
| 838 | Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, and Christopher D. Manning. 2024. Biomedlm: A 2.7b parameter language model trained on biomedical text. <i>Preprint</i> , arXiv:2403.18421. | | |
| 839 | | | |
| 840 | | | |
| 841 | | | |
| 842 | | | |
| 843 | | | |
| 844 | Andrew A. Borkowski and Alon Ben-Ari. 2025. Multiagent ai systems in health care: Envisioning next-generation intelligence. <i>Federal Practitioner</i> , 42(5):188–193. | | |
| 845 | | | |
| 846 | | | |
| 847 | | | |
| 848 | Glenn W. Brier. 1950. Verification of forecasts expressed in terms of probability. <i>Monthly Weather Review</i> , 78(1):1–3. | | |
| 849 | | | |
| 850 | | | |
| 851 | Nathan Brown, Marco Fiscato, Marwin Segler, and Alain C. Vaucher. 2019. Guacamol: Benchmarking models for de novo molecular design. <i>Journal of Chemical Information and Modeling</i> , 59(3):1096–1108. | | |
| 852 | | | |
| 853 | | | |
| 854 | | | |
| 855 | | | |
| 856 | Maria Ana Cardei, Josephine Lamp, Mark Derdzinski, and Karan Bhatia. 2025. Dexbench: Benchmarking llms for personalized decision making in diabetes management. <i>Preprint</i> , arXiv:2510.00038. | | |
| 857 | | | |
| 858 | | | |
| 859 | | | |
| 860 | Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. Extracting training data from large language models. In <i>USENIX Security Symposium</i> , pages 2633–2650. | | |
| 861 | | | |
| 862 | | | |
| 863 | | | |
| 864 | | | |
| 865 | | | |
| 866 | Centers for Disease Control and Prevention. 2020. Clear communication index user guide. | | |
| 867 | | | |
| 868 | Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. 2024. Humans or llms as the judge? a study on judgement bias. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 8301–8327. Association for Computational Linguistics. | | |
| 869 | | | |
| 870 | | | |
| 871 | | | |
| 872 | | | |
| 873 | | | |
| 874 | Shan Chen, Pedro Moreira, Yuxin Xiao, Sam Schmidgall, Jeremy Warner, Hugo Aerts, Thomas Hartvigsen, Jack Gallifant, and Danielle S. Bitterman. 2025a. Medbrowsecomp: Benchmarking medical deep research and computer use. <i>Preprint</i> , arXiv:2505.14963. | | |
| 875 | | | |
| 876 | | | |
| 877 | | | |
| 878 | | | |
| 879 | | | |
| 880 | Tong Chen, Zimu Wang, Yiyi Miao, Haoran Luo, Yuanfei Sun, Wei Wang, Zhengyong Jiang, Procheta Sen, and Jionglong Su. 2025b. Medfact: A large-scale chinese dataset for evidence-based medical fact-checking of llm responses. <i>Preprint</i> , arXiv:2509.17436. | | |
| 881 | | | |
| 882 | | | |
| 883 | | | |
| 884 | | | |
| 885 | | | |
| 886 | Xinyu Chen, Jiayi Xiang, Shudong Lu, Yifan Liu, Meng He, and Dan Shi. 2025c. Evaluating large language models and agents in healthcare: Key challenges in clinical applications. <i>Intelligent Medicine</i> , 5(2):151–163. | | |
| 887 | | | |
| 888 | | | |
| 889 | | | |
| 890 | | | |
| | Lionel Tim-Ee Cheng, Jasmine Chiat Ling Ong, Zhen Ling Teo, Ting Fang Tan, Narrendar RaviChandran, Fei Wang, Leo Anthony Celi, Marcus Eng Hock Ong, and Nan Liu. 2025. A scoping review and evidence gap analysis of clinical AI fairness. <i>npj Digital Medicine</i> . PMID: PMC12167363. | 891 | |
| | | 892 | |
| | | 893 | |
| | | 894 | |
| | | 895 | |
| | | 896 | |
| | Jan Clusmann, Franziska R. Kolbinger, Hannah S. Muti, Zoe I. Carrero, Jens-Niklas Eckardt, Nima Ghaffari Laleh, Carsten M. L. Löffler, Stefan-Christian Schwarzkopf, Matthias Unger, Gabriëlle P. Veldhuizen, Sabine J. Wagner, and Jakob N. Kather. 2023. The future landscape of large language models in medicine. <i>Communications Medicine</i> , 3:141. | 897 | |
| | | 898 | |
| | | 899 | |
| | | 900 | |
| | | 901 | |
| | | 902 | |
| | | 903 | |
| | A. Cohan, S. Feldman, I. Beltagy, D. Downey, and D. S. Weld. 2020. Specter: Document-level representation learning using citation-informed transformers. <i>arXiv preprint arXiv:2004.07180</i> . | 904 | |
| | | 905 | |
| | | 906 | |
| | | 907 | |
| | Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. 2018. A discourse-aware attention model for abstractive summarization of long documents. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 615–621, New Orleans, Louisiana. Association for Computational Linguistics. | 908 | |
| | | 909 | |
| | | 910 | |
| | | 911 | |
| | | 912 | |
| | | 913 | |
| | | 914 | |
| | | 915 | |
| | | 916 | |
| | | 917 | |
| | DeepSeek-AI. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>Preprint</i> , arXiv:2501.12948. | 918 | |
| | | 919 | |
| | | 920 | |
| | Jay DeYoung, Iz Beltagy, Madeleine van Zuylen, Bailey Kuehl, and Lucy Lu Wang. 2021. Ms2: Multi-document summarization of medical studies. <i>Preprint</i> , arXiv:2104.06486. | 921 | |
| | | 922 | |
| | | 923 | |
| | | 924 | |
| | A. K. D. Dhanvijay, M. J. Pinjar, N. Dhokane, S. R. Sorte, A. Kumari, H. Mondal, and A. K. Dhanvijay. 2023. Performance of large language models (chatgpt, bing search, and google bard) in solving case vignettes in physiology. <i>Cureus</i> , 15(8). | 925 | |
| | | 926 | |
| | | 927 | |
| | | 928 | |
| | | 929 | |
| | Jacob Dineen, Aswin RRV, Qin Liu, Zhikun Xu, Xiao Ye, Ming Shen, Zhaonan Li, Shijie Lu, Chitta Baral, Muhao Chen, and Ben Zhou. 2025. Qa-lign: Aligning llms through constitutionally decomposed qa. <i>Preprint</i> , arXiv:2506.08123. | 930 | |
| | | 931 | |
| | | 932 | |
| | | 933 | |
| | | 934 | |
| | Zhihao Fan, Jialong Tang, Wei Chen, Siyuan Wang, Zhongyu Wei, Jun Xi, Fei Huang, and Jingren Zhou. 2024. Ai hospital: Benchmarking large language models in a multi-agent medical interaction simulator. <i>Preprint</i> , arXiv:2402.09742. | 935 | |
| | | 936 | |
| | | 937 | |
| | | 938 | |
| | | 939 | |
| | Faiha Fareez, Tishya Parikh, Christopher Wavell, Saba Shahab, Meghan Chevalier, Scott Good, Isabella De Blasi, Rafik Rhouma, Christopher McMahan, Jean-Paul Lam, Thomas Lo, and Christopher W. Smith. 2022. A dataset of simulated patient–physician medical interviews with a focus on respiratory cases. <i>Scientific Data</i> , 9:313. | 940 | |
| | | 941 | |
| | | 942 | |
| | | 943 | |
| | | 944 | |
| | | 945 | |
| | | 946 | |

| | | |
|------|---|------|
| 947 | Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. Detecting hallucinations in large language models using semantic entropy . <i>Nature</i> , 630(8017):625–630. | 1002 |
| 948 | | 1003 |
| 949 | | |
| 950 | | |
| 951 | Jana Fehr, Brian Citro, Rohit Malpani, Christoph Lippert, and Vince I. Madai. 2024. A trustworthy ai reality-check: the lack of transparency of artificial intelligence products in healthcare . <i>Frontiers in Digital Health</i> , 6:1267290. | 1004 |
| 952 | | 1005 |
| 953 | | 1006 |
| 954 | | 1007 |
| 955 | | 1008 |
| 956 | | 1009 |
| 957 | Kevin J. Feng, David W. McDonald, and Amy X. Zhang. 2025a. Levels of autonomy for ai agents . arXiv:2506.12469. | 1010 |
| 958 | | 1011 |
| 959 | | 1012 |
| 960 | Yu Feng, Ben Zhou, Weidong Lin, and Dan Roth. 2025b. Bird: A trustworthy bayesian inference framework for large language models . <i>Preprint</i> , arXiv:2404.12494. | 1013 |
| 961 | | 1014 |
| 962 | | 1015 |
| 963 | | 1016 |
| 964 | Paul Festor, Ibrahim Habli, Yan Jia, Anthony Gordon, A. Aldo Faisal, and Matthieu Komorowski. 2021. Levels of autonomy and safety assurance for ai-based clinical decision systems . In <i>SAFECOMP 2021 Workshops</i> , pages 291–296. Springer. | 1017 |
| 965 | | 1018 |
| 966 | | 1019 |
| 967 | | 1020 |
| 968 | Jack Gallifant, Majid Afshar, Saleem Ameen, Yindalon Aphinyanaphongs, Sherri Chen, Giovanni E. Cacciamani, Dina Demner-Fushman, Dmitriy Dligach, Roxana Daneshjou, Christopher Fernandes, and 1 others. 2025. The tripod-llm reporting guideline for studies using large language models . <i>Nature Medicine</i> , 31(1):60–69. | 1021 |
| 969 | | 1022 |
| 970 | | 1023 |
| 971 | | 1024 |
| 972 | | 1025 |
| 973 | | 1026 |
| 974 | | 1027 |
| 975 | Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2021. Datasheets for datasets . <i>Communications of the ACM</i> , 64(12):86–92. | 1028 |
| 976 | | 1029 |
| 977 | | 1030 |
| 978 | | 1031 |
| 979 | | 1032 |
| 980 | Yonatan Geifman and Ran El-Yaniv. 2017. Selective classification for deep neural networks . In <i>Advances in Neural Information Processing Systems</i> , volume 30. | 1033 |
| 981 | | 1034 |
| 982 | | 1035 |
| 983 | | 1036 |
| 984 | Yonatan Geifman and Ran El-Yaniv. 2019. Selectivenet: A deep neural network with an integrated reject option . In <i>Proceedings of the 36th International Conference on Machine Learning (ICML 2019)</i> , volume 97 of <i>Proceedings of Machine Learning Research</i> , pages 2151–2159. PMLR. | 1037 |
| 985 | | 1038 |
| 986 | | 1039 |
| 987 | | 1040 |
| 988 | | 1041 |
| 989 | | 1042 |
| 990 | Elizabeth Goh and colleagues. 2024. Large language model influence on diagnostic reasoning . <i>JAMA Network Open</i> . | 1043 |
| 991 | | 1044 |
| 992 | | 1045 |
| 993 | R. Golan, S. J. Ripps, R. Reddy, J. Loloi, A. P. Bernstein, Z. M. Connelly, N. S. Golan, R. Ramasamy, S. Ripps, and R. V. Reddy. 2023. Chatgpt’s ability to assess quality and readability of online medical information: evidence from a cross-sectional study . <i>Cureus</i> , 15(7). | 1046 |
| 994 | | 1047 |
| 995 | | 1048 |
| 996 | | 1049 |
| 997 | | 1050 |
| 998 | Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and Jian Guo. 2024. A survey on llm-as-a-judge . <i>arXiv preprint arXiv:2411.15594</i> . | 1051 |
| 999 | | 1052 |
| 1000 | | 1053 |
| 1001 | | 1054 |
| | Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. On calibration of modern neural networks . In <i>Proceedings of the 34th International Conference on Machine Learning (ICML 2017)</i> , volume 70 of <i>Proceedings of Machine Learning Research</i> , pages 1321–1330. PMLR. | 1055 |
| | | 1056 |
| | | 1057 |
| | | 1058 |
| | | 1059 |
| | Paul Hager, Friederike Jungmann, Robbie Holland, Kunal Bhagat, Inga Hubrecht, Manuel Knauer, Jakob Vielhauer, Marcus Makowski, Rickmer Braren, Georgios Kaissis, and Daniel Rueckert. 2024. Evaluation and mitigation of the limitations of large language models in clinical decision-making . <i>Nature Medicine</i> , 30(9):2613–2622. | 1060 |
| | | 1061 |
| | | 1062 |
| | | 1063 |
| | | 1064 |
| | | 1065 |
| | | 1066 |
| | | 1067 |
| | | 1068 |
| | | 1069 |
| | | 1070 |
| | | 1071 |
| | | 1072 |
| | | 1073 |
| | | 1074 |
| | | 1075 |
| | | 1076 |
| | | 1077 |
| | | 1078 |
| | | 1079 |
| | | 1080 |
| | | 1081 |
| | | 1082 |
| | | 1083 |
| | | 1084 |
| | | 1085 |
| | | 1086 |
| | | 1087 |
| | | 1088 |
| | | 1089 |
| | | 1090 |
| | | 1091 |
| | | 1092 |
| | | 1093 |
| | | 1094 |
| | | 1095 |
| | | 1096 |
| | | 1097 |
| | | 1098 |
| | | 1099 |
| | | 1100 |
| | | 1101 |
| | | 1102 |
| | | 1103 |
| | | 1104 |
| | | 1105 |
| | | 1106 |
| | | 1107 |
| | | 1108 |
| | | 1109 |
| | | 1110 |
| | | 1111 |
| | | 1112 |
| | | 1113 |
| | | 1114 |
| | | 1115 |
| | | 1116 |
| | | 1117 |
| | | 1118 |
| | | 1119 |
| | | 1120 |
| | | 1121 |
| | | 1122 |
| | | 1123 |
| | | 1124 |
| | | 1125 |
| | | 1126 |
| | | 1127 |
| | | 1128 |
| | | 1129 |
| | | 1130 |
| | | 1131 |
| | | 1132 |
| | | 1133 |
| | | 1134 |
| | | 1135 |
| | | 1136 |
| | | 1137 |
| | | 1138 |
| | | 1139 |
| | | 1140 |
| | | 1141 |
| | | 1142 |
| | | 1143 |
| | | 1144 |
| | | 1145 |
| | | 1146 |
| | | 1147 |
| | | 1148 |
| | | 1149 |
| | | 1150 |
| | | 1151 |
| | | 1152 |
| | | 1153 |
| | | 1154 |
| | | 1155 |
| | | 1156 |
| | | 1157 |
| | | 1158 |
| | | 1159 |
| | | 1160 |
| | | 1161 |
| | | 1162 |
| | | 1163 |
| | | 1164 |
| | | 1165 |
| | | 1166 |
| | | 1167 |
| | | 1168 |
| | | 1169 |
| | | 1170 |
| | | 1171 |
| | | 1172 |
| | | 1173 |
| | | 1174 |
| | | 1175 |
| | | 1176 |
| | | 1177 |
| | | 1178 |
| | | 1179 |
| | | 1180 |
| | | 1181 |
| | | 1182 |
| | | 1183 |
| | | 1184 |
| | | 1185 |
| | | 1186 |
| | | 1187 |
| | | 1188 |
| | | 1189 |
| | | 1190 |
| | | 1191 |
| | | 1192 |
| | | 1193 |
| | | 1194 |
| | | 1195 |
| | | 1196 |
| | | 1197 |
| | | 1198 |
| | | 1199 |
| | | 1200 |
| | | 1201 |
| | | 1202 |
| | | 1203 |
| | | 1204 |
| | | 1205 |
| | | 1206 |
| | | 1207 |
| | | 1208 |
| | | 1209 |
| | | 1210 |
| | | 1211 |
| | | 1212 |
| | | 1213 |
| | | 1214 |
| | | 1215 |
| | | 1216 |
| | | 1217 |
| | | 1218 |
| | | 1219 |
| | | 1220 |
| | | 1221 |
| | | 1222 |
| | | 1223 |
| | | 1224 |
| | | 1225 |
| | | 1226 |
| | | 1227 |
| | | 1228 |
| | | 1229 |
| | | 1230 |
| | | 1231 |
| | | 1232 |
| | | 1233 |
| | | 1234 |
| | | 1235 |
| | | 1236 |
| | | 1237 |
| | | 1238 |
| | | 1239 |
| | | 1240 |
| | | 1241 |
| | | 1242 |
| | | 1243 |
| | | 1244 |
| | | 1245 |
| | | 1246 |
| | | 1247 |
| | | 1248 |
| | | 1249 |
| | | 1250 |
| | | 1251 |
| | | 1252 |
| | | 1253 |
| | | 1254 |
| | | 1255 |
| | | 1256 |
| | | 1257 |
| | | 1258 |
| | | 1259 |
| | | 1260 |
| | | 1261 |
| | | 1262 |
| | | 1263 |
| | | 1264 |
| | | 1265 |
| | | 1266 |
| | | 1267 |
| | | 1268 |
| | | 1269 |
| | | 1270 |
| | | 1271 |
| | | 1272 |
| | | 1273 |
| | | 1274 |
| | | 1275 |
| | | 1276 |
| | | 1277 |
| | | 1278 |
| | | 1279 |
| | | 1280 |
| | | 1281 |
| | | 1282 |
| | | 1283 |
| | | 1284 |
| | | 1285 |
| | | 1286 |
| | | 1287 |
| | | 1288 |
| | | 1289 |
| | | 1290 |
| | | 1291 |
| | | 1292 |
| | | 1293 |
| | | 1294 |
| | | 1295 |
| | | 1296 |
| | | 1297 |
| | | 1298 |
| | | 1299 |
| | | 1300 |
| | | 1301 |
| | | 1302 |
| | | 1303 |
| | | 1304 |
| | | 1305 |
| | | 1306 |
| | | 1307 |
| | | 1308 |
| | | 1309 |
| | | 1310 |
| | | 1311 |
| | | 1312 |
| | | 1313 |
| | | 1314 |
| | | 1315 |
| | | 1316 |
| | | 1317 |
| | | 1318 |
| | | 1319 |
| | | 1320 |
| | | 1321 |
| | | 1322 |
| | | 1323 |
| | | 1324 |
| | | 1325 |
| | | 1326 |
| | | 1327 |
| | | 1328 |
| | | 1329 |
| | | 1330 |
| | | 1331 |
| | | 1332 |
| | | 1333 |
| | | 1334 |
| | | 1335 |
| | | 1336 |
| | | 1337 |
| | | 1338 |
| | | 1339 |
| | | 1340 |
| | | 1341 |
| | | 1342 |
| | | 1343 |
| | | 1344 |
| | | 1345 |
| | | 1346 |
| | | 1347 |
| | | 1348 |
| | | 1349 |
| | | 1350 |
| | | 1351 |
| | | 1352 |
| | | 1353 |
| | | 1354 |
| | | 1355 |
| | | 1356 |
| | | 1357 |
| | | 1358 |
| | | 1359 |
| | | 1360 |
| | | 1361 |
| | | 1362 |
| | | 1363 |
| | | 1364 |
| | | 1365 |
| | | 1366 |
| | | 1367 |
| | | 1368 |
| | | 1369 |
| | | 1370 |
| | | 1371 |
| | | 1372 |
| | | 1373 |
| | | 1374 |
| | | 1375 |
| | | 1376 |
| | | 1377 |
| | | 1378 |
| | | 1379 |
| | | 1380 |
| | | 1381 |
| | | 1382 |
| | | 1383 |
| | | 1384 |
| | | 1385 |
| | | 1386 |
| | | 1387 |
| | | 1388 |
| | | 1389 |
| | | 1390 |
| | | 1391 |
| | | 1392 |
| | | 1393 |
| | | 1394 |
| | | 1395 |
| | | 1396 |
| | | 1397 |
| | | 1398 |
| | | 1399 |
| | | 1400 |
| | | 1401 |
| | | 1402 |
| | | 1403 |
| | | 1404 |
| | | 1405 |
| | | 1406 |

| | | |
|------|--|------|
| 1055 | IPDAS Collaboration. 2024. International patient decision aid standards (ipdas) collaboration . Accessed 2025-10-05. | 1111 |
| 1056 | | 1112 |
| 1057 | | 1113 |
| 1058 | Akshay Jagatap, Srujana Merugu, and Prakash Mandayam Comar. 2025. Rxlens: Multi-agent llm-powered scan and order for pharmacy . In <i>Proceedings of the 2025 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 3: Industry Track)</i> , pages 822–832, Albuquerque, New Mexico. Association for Computational Linguistics. | 1114 |
| 1059 | | 1115 |
| 1060 | | 1116 |
| 1061 | | 1117 |
| 1062 | | 1118 |
| 1063 | | 1119 |
| 1064 | | 1120 |
| 1065 | | 1121 |
| 1066 | Prithvi Jayaraman, Jack Desman, Majid Sabounchi, Girish N. Nadkarni, and Ankit Sakhuja. 2024. A primer on reinforcement learning in medicine for clinicians . <i>npj Digital Medicine</i> , 7(1):337. | 1122 |
| 1067 | | 1123 |
| 1068 | | 1124 |
| 1069 | | 1125 |
| 1070 | Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Delong Chen, Wenliang Dai, Ho Shu Chan, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation . <i>ACM Computing Surveys</i> , 55(12). | 1126 |
| 1071 | | 1127 |
| 1072 | | 1128 |
| 1073 | | 1129 |
| 1074 | | 1130 |
| 1075 | | 1131 |
| 1076 | Yixing Jiang, Kameron C. Black, Gloria Geng, Danny Park, James Zou, Andrew Y. Ng, and Jonathan H. Chen. 2025. Medagentbench: A virtual ehr environment to benchmark medical llm agents . <i>NEJM AI</i> . | 1132 |
| 1077 | | 1133 |
| 1078 | | 1134 |
| 1079 | | 1135 |
| 1080 | | 1136 |
| 1081 | D. Jin, E. Pan, N. Oufattole, W.-H. Weng, H. Fang, and P. Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. <i>Applied Sciences</i> , 11(14). | 1137 |
| 1082 | | 1138 |
| 1083 | | 1139 |
| 1084 | | 1140 |
| 1085 | | 1141 |
| 1086 | Q. Jin, W. Kim, Q. Chen, D. C. Comeau, L. Yeganova, W. J. Wilbur, and Z. Lu. 2023. Medcpt: Contrastive pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical information retrieval . <i>Bioinformatics</i> , 39(11). | 1142 |
| 1087 | | 1143 |
| 1088 | | 1144 |
| 1089 | | 1145 |
| 1090 | | 1146 |
| 1091 | A. E. Johnson, T. J. Pollard, S. J. Berkowitz, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, R. G. Mark, and S. Horng. 2019. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports . <i>Scientific data</i> , 6(1). | 1147 |
| 1092 | | 1148 |
| 1093 | | 1149 |
| 1094 | | 1150 |
| 1095 | | 1151 |
| 1096 | Dongwon Jung, Qin Liu, Tenghao Huang, Ben Zhou, and Muhao Chen. 2024. Familiarity-aware evidence compression for retrieval-augmented generation . <i>Preprint</i> , arXiv:2409.12468. | 1152 |
| 1097 | | 1153 |
| 1098 | | 1154 |
| 1099 | | 1155 |
| 1100 | Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Samuel Bowman, Stanislav Fort, and 17 others. 2022. Language models (mostly) know what they know . <i>arXiv preprint arXiv:2207.05221</i> . | 1156 |
| 1101 | | 1157 |
| 1102 | | 1158 |
| 1103 | | 1159 |
| 1104 | | 1160 |
| 1105 | | 1161 |
| 1106 | | 1162 |
| 1107 | | 1163 |
| 1108 | | 1164 |
| 1109 | Nikhil Khandekar, Qiao Jin, Guangzhi Xiong, Soren Dunn, Serina S. Applebaum, Zain Anwar, Maame Sarfo-Gyamfi, Conrad W. Safranek, Abid A. Anwar, Andrew Zhang, Aidan Gilson, Maxwell B. Singer, Amisha Dave, Andrew Taylor, Aidong Zhang, Qingyu Chen, and Zhiyong Lu. 2024. Medcalc-bench: Evaluating large language models for medical calculations . <i>Preprint</i> , arXiv:2406.12036. | 1165 |
| 1110 | | 1166 |
| 1111 | | 1167 |
| 1112 | | 1168 |
| 1113 | | 1169 |
| 1114 | | 1170 |
| 1115 | | 1171 |
| 1116 | | 1172 |
| 1117 | | 1173 |
| 1118 | | 1174 |
| 1119 | | 1175 |
| 1120 | | 1176 |
| 1121 | | 1177 |
| 1122 | | 1178 |
| 1123 | | 1179 |
| 1124 | | 1180 |
| 1125 | | 1181 |
| 1126 | | 1182 |
| 1127 | | 1183 |
| 1128 | | 1184 |
| 1129 | | 1185 |
| 1130 | | 1186 |
| 1131 | | 1187 |
| 1132 | | 1188 |
| 1133 | | 1189 |
| 1134 | | 1190 |
| 1135 | | 1191 |
| 1136 | | 1192 |
| 1137 | | 1193 |
| 1138 | | 1194 |
| 1139 | | 1195 |
| 1140 | | 1196 |
| 1141 | | 1197 |
| 1142 | | 1198 |
| 1143 | | 1199 |
| 1144 | | 1200 |
| 1145 | | 1201 |
| 1146 | | 1202 |
| 1147 | | 1203 |
| 1148 | | 1204 |
| 1149 | | 1205 |
| 1150 | | 1206 |
| 1151 | | 1207 |
| 1152 | | 1208 |
| 1153 | | 1209 |
| 1154 | | 1210 |
| 1155 | | 1211 |
| 1156 | | 1212 |
| 1157 | | 1213 |
| 1158 | | 1214 |
| 1159 | | 1215 |
| 1160 | | 1216 |
| 1161 | | 1217 |
| 1162 | | 1218 |
| 1163 | | 1219 |
| 1164 | | 1220 |
| 1165 | | 1221 |

| | | | |
|------|--|---|------|
| 1166 | G. Liu, Y. Liao, F. Wang, B. Zhang, L. Zhang, X. Liang, | Mary L. McHugh. 2012. Interrater reliability: the kappa | 1222 |
| 1167 | X. Wan, S. Li, Z. Li, and S. Zhang. 2021. Medical- | statistic . <i>Biochemia Medica</i> , 22(3):276–282. | 1223 |
| 1168 | vlbert: Medical visual language bert for covid-19 | | |
| 1169 | ct report generation with alternate learning. <i>IEEE</i> | Nikita Mehandru, Niloufar Golchini, David Bamman, | 1224 |
| 1170 | <i>Transactions on Neural Networks and Learning Sys-</i> | Travis Zack, Melanie F. Molina, and Ahmed Alaa. | 1225 |
| 1171 | <i>tems</i> , 32(9):3786–3797. | 2025. Er-reason: A benchmark dataset for llm-based | 1226 |
| | | clinical reasoning in the emergency room . <i>Preprint</i> , | 1227 |
| | | arXiv:2505.22919. | 1228 |
| 1172 | W. Liu, J. Tang, Y. Cheng, W. Li, Y. Zheng, and | Xiangbin Meng, Xiangyu Yan, Kuo Zhang, Da Liu, Xi- | 1229 |
| 1173 | X. Liang. 2022. Meddg: an entity-centric medi- | aojuan Cui, Yaodong Yang, Muhan Zhang, Chunxia | 1230 |
| 1174 | cal consultation dataset for entity-aware medical dia- | Cao, Jingjia Wang, Xuliang Wang, Jun Gao, Ya- | 1231 |
| 1175 | logue generation. <i>CCF International Conference on</i> | Geng-Shun Wang, Jun Ji, Zhen Qiu, Min Li, | 1232 |
| 1176 | <i>Natural Language Processing and Chinese Comput-</i> | Changhong Qian, Tengda Guo, Shu Ma, Zhi Wang, | 1233 |
| 1177 | <i>ing</i> , pages 447–459. | and Yue-Dong Tang. 2024. The application of large | 1234 |
| | | language models in medicine: A scoping review . | 1235 |
| 1178 | Xiaoxuan Liu, S. C. Rivera, D. Moher, and et al. 2020. | <i>iScience</i> , 27(5):109713. | 1236 |
| 1179 | Consort-ai extension: Reporting guidelines for clinical | | |
| 1180 | trial reports for interventions involving artificial | | |
| 1181 | intelligence . <i>Nature Medicine</i> , 26:1364–1374. | | |
| | | Margaret Mitchell, Simone Wu, Andrew Zaldivar, | 1237 |
| 1182 | Juan Antonio Lossio-Ventura, Ran Sun, Sebastien Bous- | Parker Barnes, Lucy Vasserman, Ben Hutchinson, | 1238 |
| 1183 | sard, and Tina Hernandez-Boussard. 2023. Clinical | Elena Spitzer, Inioluwa Deborah Raji, and Timnit | 1239 |
| 1184 | concept recognition: Evaluation of existing sys- | Gebru. 2019. Model cards for model reporting . In | 1240 |
| 1185 | tems on ehrr . <i>Frontiers in Artificial Intelligence</i> , | <i>Proceedings of the Conference on Fairness, Account-</i> | 1241 |
| 1186 | 5:1051724. | <i>ability, and Transparency (FAT*)</i> , pages 220–229. | 1242 |
| | | | |
| 1187 | Zhong Lou, Chen Zhao, Wenhao Li, Chenyu Zhang, | Subhabrata Mukherjee, Paul Gamble, Markel | 1243 |
| 1188 | Qi Zhu, Xuehai Pan, Mohit Bansal, and Kai-Wei | Sanz Ausin, Neel Kant, Kriti Aggarwal, Neha Man- | 1244 |
| 1189 | Chang. 2024. Anchoring effects in large language | junath, Debajyoti Datta, Zhengliang Liu, Jiayuan | 1245 |
| 1190 | models . <i>arXiv preprint arXiv:2412.06593</i> . | Ding, Sophia Busacca, Cezanne Bianco, Swapnil | 1246 |
| | | Sharma, Rae Lasko, Michelle Voisard, Sanchay | 1247 |
| 1191 | Ruofan Lu, Yichen Li, and Yintong Huo. 2025. Ex- | Harneja, Darya Filippova, Gerry Meixiong, Kevin | 1248 |
| 1192 | ploring autonomous agents: A closer look at why | Cha, Amir Youssefi, and 7 others. 2024. Polaris: | 1249 |
| 1193 | they fail when completing tasks . <i>arXiv preprint</i> | A safety-focused llm constellation architecture for | 1250 |
| 1194 | <i>arXiv:2508.13143</i> . | healthcare . <i>Preprint</i> , arXiv:2403.13313. | 1251 |
| | | | |
| 1195 | Ling Luo, Po-Ting Lai, Chih-Hsuan Wei, Cecilia N. | National Institute of Standards and Technology. 2023. | 1252 |
| 1196 | Arighi, and Zhiyong Lu. 2022. Biored: A rich | Artificial intelligence risk management framework (ai | 1253 |
| 1197 | biomedical relation extraction dataset. <i>Briefings in</i> | rmf 1.0). https://nvlpubs.nist.gov/nistpubs/ | 1254 |
| 1198 | <i>Bioinformatics</i> . | ai/nist.ai.100-1.pdf . Accessed Oct 5, 2025. | 1255 |
| | | | |
| 1199 | Zizhan Ma, Wenxuan Wang, Guo Yu, Yiu-Fai Che- | Zabir Al Nazi and Wei Peng. 2024. Large language | 1256 |
| 1200 | ung, Meidan Ding, Jie Liu, Wenting Chen, and Lin- | models in healthcare and medical domain: A review . | 1257 |
| 1201 | lin Shen. 2025. Beyond the leaderboard: Rethink- | <i>Informatics</i> , 11(3):57. | 1258 |
| 1202 | ing medical benchmarks for large language models . | | |
| 1203 | <i>Preprint</i> , arXiv:2508.04325. | | |
| | | Denis Newman-Griffis, Guy Divita, Bart Desmet, Ayah | 1259 |
| 1204 | Andreas Madsen, Sarath Chandar, and Siva Reddy. 2024. | Zirikly, Carolyn P Rosé, and Eric Fosler-Lussier. | 1260 |
| 1205 | Are self-explanations from large language models | 2021. Ambiguity in medical concept normalization: | 1261 |
| 1206 | faithful? <i>arXiv preprint arXiv:2401.07927</i> . | An analysis of types and coverage in electronic health | 1262 |
| | | record datasets . <i>Journal of the American Medical In-</i> | 1263 |
| 1207 | Diwakar Mahajan, Jennifer J. Liang, Ching-Huei Tsou, | <i>formatics Association</i> , 28(3):516–532. | 1264 |
| 1208 | and Özlem Uzuner. 2023. Overview of the 2022 | | |
| 1209 | n2c2 shared task on contextualized medication event | S. Nowak, D. Biesner, Y. Layer, M. Theis, H. Schneider, | 1265 |
| 1210 | extraction in clinical notes . <i>Journal of Biomedical</i> | W. Block, B. Wulff, U. Attenberger, R. Sifa, and | 1266 |
| 1211 | <i>Informatics</i> , 144:104432. | A. Sprinkart. 2023. Transformer-based structuring | 1267 |
| | | of free-text radiology report databases . <i>European</i> | 1268 |
| 1212 | Christopher D. Manning, Prabhakar Raghavan, and Hin- | <i>Radiology</i> , 33(6):4228–4236. | 1269 |
| 1213 | rich Schütze. 2008. Introduction to Information Re- | | |
| 1214 | trieval . Cambridge University Press, Cambridge, | Juliana Damasio Oliveira, Henrique D. P. Santos, Ana | 1270 |
| 1215 | UK. | Helena D. P. S. Ulbrich, Julia Colleoni Couto, | 1271 |
| | | Marcelo Arocha, Joaquim Santos, Manuela Martins | 1272 |
| 1216 | Joshua Maynez, Shashi Narayan, Bernd Bohnet, and | Costa, Daniela Faccio, Fabio O. Tabalipa, and Ro- | 1273 |
| 1217 | Ryan McDonald. 2020. On faithfulness and factu- | drigo F. Nogueira. 2025. Development and evaluation | 1274 |
| 1218 | ality in abstractive summarization . In <i>Proceedings</i> | of a clinical note summarization system using | 1275 |
| 1219 | <i>of the 58th Annual Meeting of the Association for</i> | large language models . <i>Communications Medicine</i> , | 1276 |
| 1220 | <i>Computational Linguistics</i> , pages 1906–1919, On- | 5(1):376. | 1277 |
| 1221 | line. Association for Computational Linguistics. | | |

| | | | |
|------|--|--|------|
| 1278 | Jesutofunmi A. Omiye, Jenna C. Lester, Simon Spichak, | Kirk Roberts, Tasmee Alam, Steven Bedrick, Dina | 1332 |
| 1279 | Veronica Rotemberg, and Roxana Daneshjou. | Demner-Fushman, Kyle Lo, Ian Soboroff, Ellen | 1333 |
| 1280 | 2023. Large language models propagate race-based | Voorhees, Lucy Lu Wang, and William R. Hersh. | 1334 |
| 1281 | medicine . <i>npj Digital Medicine</i> , 6:195. | 2021. Searching for scientific evidence in a pan- | 1335 |
| | | demic: An overview of trec-covid . <i>Journal of</i> | 1336 |
| 1282 | OpenAI. 2025. Introducing GPT-5 . https:// | <i>Biomedical Informatics</i> , 121:103865. | 1337 |
| 1283 | openai.com/index/introducing-gpt-5/ . Ac- | | |
| 1284 | cessed: 2025-10-02. | Aswin RRV, Jacob Dineen, Divij Handa, Md Nayem | 1338 |
| | | Uddin, Mihir Parmar, Chitta Baral, and Ben Zhou. | 1339 |
| 1285 | L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, | 2025. Thinktuning: Instilling cognitive reflections | 1340 |
| 1286 | P. Mishkin, C. Zhang, S. Agarwal, K. Slama, and | without distillation . <i>Preprint</i> , arXiv:2508.07616. | 1341 |
| 1287 | A. Ray. 2022. Training language models to follow in- | | |
| 1288 | structions with human feedback . <i>Advances in neural</i> | Atiqer Rahman Sarkar, Yao-Shun Chuang, Noman | 1342 |
| 1289 | <i>information processing systems</i> , 35. | Mohammed, et al., and Xiaoqian Jiang. 2024. De- | 1343 |
| | | identification is not enough: a comparison between | 1344 |
| 1290 | A. Pal, L. K. Umapathi, and M. Sankarasubbu. 2022. | de-identified and synthetic clinical notes . <i>Scientific</i> | 1345 |
| 1291 | Medmcqa: A large-scale multi-subject multi-choice | <i>Reports</i> , 14(29669). | 1346 |
| 1292 | dataset for medical domain question answering. <i>Con-</i> | | |
| 1293 | <i>ference on health, inference, and learning</i> , pages | Thomas Savage, Ashwin Nayak, Robert Gallo, Ekanath | 1347 |
| 1294 | 248–260. | Rangan, and Jonathan H. Chen. 2024a. Diagnostic | 1348 |
| | | reasoning prompts reveal the potential for large lan- | 1349 |
| 1295 | C. Peng, X. Yang, A. Chen, K. E. Smith, N. PourNeja- | guage model interpretability in medicine . <i>npj Digital</i> | 1350 |
| 1296 | tian, A. B. Costa, C. Martin, M. G. Flores, Y. Zhang, | <i>Medicine</i> , 7(20). | 1351 |
| 1297 | and T. Magoc. 2023. A study of generative large | | |
| 1298 | language model for medical research and healthcare . | Thomas Savage, Ashwin Nayak, Robert Gallo, Ekanath | 1352 |
| 1299 | <i>NPJ Digital Medicine</i> , 6(1). | Rangan, and Jonathan H. Chen. 2024b. Diagnostic | 1353 |
| | | reasoning prompts reveal the potential for large lan- | 1354 |
| 1300 | Jahan C. Penny-Dimri, Magdalena Bachmann, | guage model interpretability in medicine . <i>npj Digital</i> | 1355 |
| 1301 | William R. Cooke, Sam Mathewlynn, Samuel | <i>Medicine</i> , 7(1):20. | 1356 |
| 1302 | Dockree, John Tolladay, Jannik Kossen, Lin Li, | | |
| 1303 | Yarin Gal, and Gabriel Davis Jones. 2025. Reducing | Samuel Schmidgall, Rojin Ziaei, Carl Harris, Eduardo | 1357 |
| 1304 | large language model safety risks in women’s health | Reis, Jeffrey Jopling, and Michael Moor. 2024. | 1358 |
| 1305 | using semantic entropy . <i>Preprint</i> , arXiv:2503.00269. | Agentclinic: a multimodal agent benchmark to evalu- | 1359 |
| | | ate ai in simulated clinical environments . <i>Preprint</i> , | 1360 |
| 1306 | Pouya Pezeshkpour and Estevam Hruschka. 2023. | arXiv:2405.07960. | 1361 |
| 1307 | Large language models sensitivity to the order of | Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane | 1362 |
| 1308 | options in multiple-choice questions . <i>arXiv preprint</i> | Suhr. 2023. Quantifying language models’ sensitiv- | 1363 |
| 1309 | <i>arXiv:2308.11483</i> . | ity to spurious features in prompt design or: How | 1364 |
| | | i learned to start worrying about prompt format- | 1365 |
| 1310 | Daniil Polykovskiy, Alexander Zhebrak, Benjamin | ting . <i>arXiv preprint arXiv:2310.11324</i> . ICLR 2024 | 1366 |
| 1311 | Sanchez-Lengeling, Sergey Golovanov, Oktai | camera-ready version. | 1367 |
| 1312 | Tatanov, Stanislav Belyaev, Rauf Kurbanov, Alek- | | |
| 1313 | sey Artamonov, Vladimir Aladinskiy, Mark Veselov, | Isabel Segura-Bedmar, Paloma Martínez, and David | 1368 |
| 1314 | Artur Kadurin, Simon Johansson, Hongming Chen, | Herrero-Zazo. 2013. Semeval-2013 task 9: Extrac- | 1369 |
| 1315 | Sergey Nikolenko, Alan Aspuru-Guzik, and Alex | tion of drug–drug interactions from biomedical texts | 1370 |
| 1316 | Zhavoronkov. 2018. Molecular sets (moses): A | (ddiextraction 2013) . In <i>Proceedings of SemEval</i> | 1371 |
| 1317 | benchmarking platform for molecular generation | <i>2013</i> . | 1372 |
| 1318 | models . <i>Preprint</i> , arXiv:1811.12823. | Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. | 1373 |
| | | BLEURT: Learning robust metrics for text genera- | 1374 |
| 1319 | Maja Popović. 2015. chrF: character n-gram f-score | tion . In <i>Proceedings of the 58th Annual Meeting of</i> | 1375 |
| 1320 | for automatic MT evaluation . In <i>Proceedings of the</i> | <i>the Association for Computational Linguistics</i> , pages | 1376 |
| 1321 | <i>Tenth Workshop on Statistical Machine Translation</i> , | 7881–7892. Association for Computational Linguis- | 1377 |
| 1322 | pages 392–395, Lisbon, Portugal. Association for | <i>tics</i> . | 1378 |
| 1323 | Computational Linguistics. | E. Sezgin, F. Chekeni, J. Lee, and S. Keim. 2023. Clin- | 1379 |
| | | ical accuracy of large language models and google | 1380 |
| 1324 | Sushant Prakash, Katherine Heller, Alan Karthike- | search responses to postpartum depression questions: | 1381 |
| 1325 | salingam, Christopher Semturs, Joelle Barral, Greg | cross-sectional study . <i>Journal of Medical Internet</i> | 1382 |
| 1326 | Corrado, Yossi Matias, Jamila Smith-Loud, Ivor | <i>Research</i> , 25. | 1383 |
| 1327 | Horn, and Karan Singhal. 2024. A toolbox for surfac- | | |
| 1328 | ing health equity harms and biases in large language | Mrinank Sharma, Meg Tong, Tomasz Korbak, David | 1384 |
| 1329 | models . <i>Nature Medicine</i> , 30:3590–3600. | Duvenaud, Amanda Askill, Samuel R. Bowman, | 1385 |
| | | Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, | 1386 |
| 1330 | Qwen Team. 2025. Qwen3 technical report . <i>arXiv</i> | Scott R. Johnston, Shauna Kravec, Timothy Maxwell, | 1387 |
| 1331 | <i>preprint arXiv:2505.09388</i> . | | |

| | | | |
|------|--|--|------|
| 1388 | Sam McCandlish, Kamal Ndousse, Oliver Rausch, | Ivan Sviridov, Amina Miftakhova, Artemiy | 1445 |
| 1389 | Nicholas Schiefer, Da Yan, Miranda Zhang, and | Tereshchenko, Galina Zubkova, Pavel Blinov, | 1446 |
| 1390 | Ethan Perez. 2023. Towards understanding syco- | and Andrey Savchenko. 2025. 3mdbench: Medical | 1447 |
| 1391 | phancy in language models. <i>arXiv preprint</i> | multimodal multi-agent dialogue benchmark. <i>arXiv</i> | 1448 |
| 1392 | arXiv:2310.13548. | preprint arXiv:2504.13861. | 1449 |
| 1393 | Ming Shen, Zhikun Xu, Jacob Dineen, Xiao Ye, and | Thomas Yu Chow Tam, Sonish Sivarajkumar, Sumit | 1450 |
| 1394 | Ben Zhou. 2025. Bow: Reinforcement learning | Kapoor, Alisa V. Stolyar, Katelyn Polanska, Kar- | 1451 |
| 1395 | for bottlenecked next word prediction. <i>Preprint,</i> | leigh R. McCarthy, Hunter Osterhoudt, Xizhi Wu, | 1452 |
| 1396 | arXiv:2506.13502. | Shyam Visweswaran, Sunyang Fu, Piyush Mathur, | 1453 |
| 1397 | Sarah J. Shoemaker, Michael S. Wolf, and Cindy | Giovanni E. Cacciamani, Cong Sun, Yifan Peng, and | 1454 |
| 1398 | Brach. 2013. The Patient Education Materials | Yanshan Wang. 2024. A framework for human eval- | 1455 |
| 1399 | Assessment Tool (PEMAT) and User's Guide (Version | uation of large language models in healthcare de- | 1456 |
| 1400 | 1.0). Agency for Healthcare Research and Quality, | rived from literature review. <i>npj Digital Medicine,</i> | 1457 |
| 1401 | Rockville, MD. AHRQ Publication No. 14-0002-EF. | 7(1):258. | 1458 |
| 1402 | Updated August 2014. | Sijun Tan, Siyuan Zhuang, Kyle Montgomery, | 1459 |
| 1403 | Sina Shool, Sara Adimi, Reza Saboori Amleshi, Ehsan | William Y. Tang, Alejandro Cuadron, Chenguang | 1460 |
| 1404 | Bitaraf, Reza Golpira, and Mahmood Tara. 2025. A | Wang, Raluca Ada Popa, and Ion Stoica. 2025. | 1461 |
| 1405 | systematic review of large language model (llm) | Judgebench: A benchmark for evaluating llm-based | 1462 |
| 1406 | evaluations in clinical medicine. <i>BMC Medical Informat-</i> | judges. In <i>Proceedings of the International Confer-</i> | 1463 |
| 1407 | ics and Decision Making, 25(117). | ence on Learning Representations (ICLR). ICLR | 1464 |
| 1408 | Som Singh, Aleena Jamal, and Fawad Qureshi. 2024. | 2025 (published version); arXiv:2410.12784. | 1465 |
| 1409 | Readability metrics in patient education: Where do | Ting Fang Tan, Kabilan Elangovan, Jasmine Ong, | 1466 |
| 1410 | we innovate? <i>Clinics and Practice,</i> 14(6):2341– | Nigam Shah, Joseph Sung, Tien Yin Wong, Lan | 1467 |
| 1411 | 2349. | Xue, Nan Liu, Haibo Wang, Chang Fu Kuo, Simon | 1468 |
| 1412 | K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. | Chesterman, Zee Kin Yeong, and Daniel S. W. Ting. | 1469 |
| 1413 | Chung, N. Scales, A. Tanwani, H. Cole-Lewis, and | 2024. A proposed s.c.o.r.e. evaluation framework for | 1470 |
| 1414 | S. Pfohl. 2023. Large language models encode clinical | large language models: Safety, consensus, objectiv- | 1471 |
| 1415 | knowledge. <i>Nature,</i> 620(7972):172–180. | ity, reproducibility and explainability. <i>arXiv preprint</i> | 1472 |
| 1416 | Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, | arXiv:2407.07666. | 1473 |
| 1417 | Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin | Liyan Tang, Zhaoyi Sun, Betina Idnay, Jordan G. Nestor, | 1474 |
| 1418 | Clark, Stephen R. Pfohl, Heather Cole-Lewis, Dar- | Ali Soroush, Pierre A. Elias, Ziyang Xu, Ying Ding, | 1475 |
| 1419 | lene Neal, Qazi Mamunur Rashid, Mike Schaecker- | Greg Durrett, Justin F. Rousseau, Chunhua Weng, | 1476 |
| 1420 | mann, Amy Wang, Dev Dash, Jonathan H. Chen, | and Yifan Peng. 2023. Evaluating large language | 1477 |
| 1421 | Nigam H. Shah, Sami Lachgar, Philip Andrew Mans- | models on medical evidence summarization. <i>npj</i> | 1478 |
| 1422 | field, and 16 others. 2025. Toward expert-level medi- | Digital Medicine, 6(1):158. | 1479 |
| 1423 | cal question answering with large language models. | G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhu- | 1480 |
| 1424 | <i>Nature Medicine,</i> 31(3):943–950. | patiraju, S. Pathak, L. Sifre, M. Rivi'ere, M. S. Kale, | 1481 |
| 1425 | Ali Soroush, Benjamin S. Glicksberg, Eyal Zimlich- | and J. Love. 2024. Gemma: Open models based | 1482 |
| 1426 | man, Yiftach Barash, Robert Freeman, Alexander W. | on gemini research and technology. <i>arXiv preprint</i> | 1483 |
| 1427 | Charney, Girish N. Nadkarni, and Eyal Klang. 2024. | arXiv:2403.08295. | 1484 |
| 1428 | Large language models are poor medical coders — | N. Thakur, N. Reimers, A. R'uckl'e, A. Srivastava, | 1485 |
| 1429 | benchmarking of medical code querying. <i>NEJM AI,</i> | and I. Gurevych. 2021. Beir: A heterogenous bench- | 1486 |
| 1430 | 1(5). | mark for zero-shot evaluation of information retrieval | 1487 |
| 1431 | Adarsh Srinivasan, Jacob Dineen, Muhammad Umar | models. <i>arXiv preprint arXiv:2104.08663.</i> | 1488 |
| 1432 | Afzal, Muhammad Uzair Sarfraz, Irbaz B. Riaz, and | Aravinthan J. Thirunavukarasu, Daniel Shu Wei Ting, | 1489 |
| 1433 | Ben Zhou. 2025. Recap: Transparent inference-time | Karthik Elangovan, Luis Gutierrez, Thelma F. Tan, | 1490 |
| 1434 | emotion alignment for medical dialogue systems. | and Daniel Shu Jin Ting. 2023. Large language mod- | 1491 |
| 1435 | <i>Preprint,</i> arXiv:2509.10746. | els in medicine. <i>Nature Medicine,</i> 29:1930–1940. | 1492 |
| 1436 | Dawn Stacey, Robert J. Volk, and IPDAS Evidence Up- | Jeremias Traub, Till J. Bungert, Carsten T. Lüth, | 1493 |
| 1437 | date Leads. 2021. The international patient decision | Michael Baumgartner, Klaus H. Maier-Hein, Lena | 1494 |
| 1438 | aid standards (ipdas) collaboration: Evidence update | Maier-Hein, and Paul F. Jaeger. 2024. Overcoming | 1495 |
| 1439 | 2.0. <i>Medical Decision Making,</i> 41(7):729–733. | common flaws in the evaluation of selective classifi- | 1496 |
| 1440 | Mark Steyvers, Heliodoro Tejada, Aakriti Kumar, Cata- | cation systems. In <i>Advances in Neural Information</i> | 1497 |
| 1441 | rina Belem, Sheer Karny, Xinyue Hu, Lukas W. | Processing Systems. NeurIPS 2024 (spotlight). | 1498 |
| 1442 | Mayer, and Padhraic Smyth. 2025. What large lan- | Miles Turpin, Julian Michael, Ethan Perez, and | 1499 |
| 1443 | guage models know and what people think they know. | Samuel R. Bowman. 2023. Language models don't | 1500 |
| 1444 | <i>Nature Machine Intelligence,</i> 7:221–231. | always say what they think: Unfaithful explanations | 1501 |

| | | | |
|------|--|---|------|
| 1502 | in chain-of-thought prompting. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> . | Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. Benchmarking retrieval-augmented generation for medicine. <i>Preprint</i> , arXiv:2402.13178. | 1555 |
| 1503 | | | 1556 |
| 1504 | Özlem Uzuner, Brett R. South, Shuying Shen, and Scott L. DuVall. 2011. 2010 i2b2/va challenge on concepts, assertions, and relations in clinical text. <i>Journal of the American Medical Informatics Association</i> , 18(5):552–556. | Ran Xu, Yuchen Zhuang, Yishan Zhong, Yue Yu, Zifeng Wang, Xiangru Tang, Hang Wu, May D. Wang, Peifeng Ruan, Donghan Yang, Tao Wang, Guanghua Xiao, Xin Liu, Carl Yang, Yang Xie, and Wenqi Shi. 2025a. Medagentgym: A scalable agentic training environment for code-centric reasoning in biomedical data science. <i>Preprint</i> , arXiv:2506.04405. | 1557 |
| 1505 | | | 1558 |
| 1506 | | | 1559 |
| 1507 | | | 1560 |
| 1508 | | | 1561 |
| 1509 | Phuc Phan Van, Dat Nguyen Minh, An Dinh Ngoc, and Huy Phan Thanh. 2024. Rx strategist: Prescription verification using llm agents system. <i>Preprint</i> , arXiv:2409.03440. | | 1562 |
| 1510 | | | 1563 |
| 1511 | | | 1564 |
| 1512 | | | 1565 |
| 1513 | Li Wang, Xi Chen, Xiangwen Deng, Hao Wen, Mingke You, Weizhi Liu, Qi Li, and Jian Li. 2024. Prompt engineering in consistency and reliability with the evidence-based guideline for LLMs. <i>npj Digital Medicine</i> , 7(1):41. | Zhikun Xu, Ming Shen, Jacob Dineen, Zhaonan Li, Xiao Ye, Shijie Lu, Aswin RRV, Chitta Baral, and Ben Zhou. 2025b. Tow: Thoughts of words improve reasoning in large language models. <i>Preprint</i> , arXiv:2410.16235. | 1566 |
| 1514 | | | 1567 |
| 1515 | | | 1568 |
| 1516 | | | 1569 |
| 1517 | | | 1570 |
| 1518 | Lucy Lu Wang, Jay DeYoung, and Byron Wallace. 2022. Overview of MSLR2022: A shared task on multi-document summarization for literature reviews. In <i>Proceedings of the Third Workshop on Scholarly Document Processing</i> , pages 175–180, Gyeongju, Republic of Korea. Association for Computational Linguistics. | Weixiang Yan, Haitian Liu, Tengxiao Wu, Qian Chen, Wen Wang, Haoyuan Chai, Jiayi Wang, Weishan Zhao, Yixin Zhang, Renjun Zhang, Li Zhu, and Xuan-dong Zhao. 2024. Clinicallab: Aligning agents for multi-departmental clinical diagnostics in the real world. <i>Preprint</i> , arXiv:2406.13890. | 1571 |
| 1519 | | | 1572 |
| 1520 | | | 1573 |
| 1521 | | | 1574 |
| 1522 | | | 1575 |
| 1523 | | | 1576 |
| 1524 | | | 1577 |
| 1525 | Xiaosong Wang, Yifan Peng, Le Lu, and 1 others. 2017. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 3462–3471. | Rui Yang, Yilin Ning, Emilia Keppo, Mingxuan Liu, Chuan Hong, Danielle S. Bitterman, Jasmine C. L. Ong, Daniel S. W. Ting, and Nan Liu. 2025. Retrieval-augmented generation for generative artificial intelligence in health care. <i>npj Health Systems</i> , 2(2):2. | 1578 |
| 1526 | | | 1579 |
| 1527 | | | 1580 |
| 1528 | | | 1581 |
| 1529 | | | 1582 |
| 1530 | | | 1583 |
| 1531 | Koki Wataoka, Tsubasa Takahashi, and Ryokan Ri. 2024. Self-preference bias in llm-as-a-judge. <i>arXiv preprint arXiv:2410.21819</i> . | Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, Nitesh V. Chawla, and Xi-angliang Zhang. 2025a. Justice or prejudice? quantifying biases in llm-as-a-judge. In <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> . ICLR 2025 Poster; arXiv:2410.02736. | 1584 |
| 1532 | | | 1585 |
| 1533 | | | 1586 |
| 1534 | Christopher Y. K. Williams, Jaskaran Bains, Tianyu Tang, Kishan Patel, Alexa N. Lucas, Fiona Chen, Brenda Y. Miao, Atul J. Butte, and Aaron E. Kornblith. 2025. Evaluating large language models for drafting emergency department encounter summaries. <i>PLOS Digital Health</i> , 4(6):e0000899. | Xiao Ye, Shaswat Shrivastava, Zhaonan Li, Jacob Dineen, Shijie Lu, Avneet Ahuja, Ming Shen, Zhikun Xu, and Ben Zhou. 2025b. Cc-learn: Cohort-based consistency learning. <i>Preprint</i> , arXiv:2506.15662. | 1587 |
| 1535 | | | 1588 |
| 1536 | | | 1589 |
| 1537 | | | 1590 |
| 1538 | | | 1591 |
| 1539 | | | 1592 |
| 1540 | World Health Organization. 2025. Ethics and governance of artificial intelligence for health: Guidance on large multi-modal models. https://www.who.int/publications/i/item/9789240084759 . Accessed Oct 5, 2025. | Wen-wai Yim, Yujuan Fu, Asma Ben Abacha, Neal Snider, Thomas Lin, and Meliha Yetisgen. 2023. Acibench: A novel ambient clinical intelligence dataset for benchmarking automatic visit note generation. <i>Scientific Data</i> , 10. | 1593 |
| 1541 | | | 1594 |
| 1542 | | | 1595 |
| 1543 | | | 1596 |
| 1544 | | | 1597 |
| 1545 | K. Wu, R. Wang, H. Quan, C. Lin, E. Lo, H. Tu, F. Yu, and H. T. Lin. 2025. An automated framework for assessing how well large language models cite relevant medical references. <i>Nature Communications</i> , 16(1):427. | Erjian Yu, Xiaofeng Chu, Wei Zhang, Xiangbin Meng, Yaodong Yang, Xue Ji, and Chao Wu. 2025. Large language models in medicine: Applications, challenges, and future directions. <i>International Journal of Medical Sciences</i> , 22(11):2792–2801. | 1598 |
| 1546 | | | 1599 |
| 1547 | | | 1600 |
| 1548 | | | 1601 |
| 1549 | | | 1602 |
| 1550 | Hanguang Xiao, Feizhong Zhou, Xingyue Liu, Tianqi Liu, Zhipeng Li, Xin Liu, and Xiaoxuan Huang. 2025. A comprehensive survey of large language models and multimodal large language models in medicine. <i>Information Fusion</i> , 117:102888. | Matthew Yue, Zhikun Xu, Vivek Gupta, Thao Ha, Liesal Sharabi, and Ben Zhou. 2025. Relate-sim: Leveraging turning point theory and llm agents to predict and understand long-term relationship dynamics through interactive narrative simulations. <i>Preprint</i> , arXiv:2510.00414. | 1603 |
| 1551 | | | 1604 |
| 1552 | | | 1605 |
| 1553 | | | 1606 |
| 1554 | | | 1607 |
| | | | 1608 |
| | | | 1609 |

| | | | |
|------|---|--|------|
| 1610 | T. Zack, E. Lehman, M. Suzgun, J. A. Rodriguez, L. A. Celi, J. Gichoya, D. Jurafsky, P. Szolovits, D. W. Bates, and R.-E. E. Abdunour. 2024. Assessing the potential of gpt-4 to perpetuate racial and gender biases in health care: a model evaluation study. <i>The Lancet Digital Health</i> , 6(1). | are scalable judges. In <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> . ICLR 2025 Spotlight; arXiv:2310.17631. | 1666 |
| 1611 | | | 1667 |
| 1612 | | | 1668 |
| 1613 | | | |
| 1614 | | A Appendix | 1669 |
| 1615 | | | |
| 1616 | Cyril Zakka, Joseph Cho, Gracia Fahed, Rohan Shad, Michael Moor, Robyn Fong, Dhamanpreet Kaur, Vishnu Ravi, Oliver Aalami, Roxana Daneshjou, Akshay Chaudhari, and William Hiesinger. 2024. Almanac copilot: Towards autonomous electronic health record navigation . Preprint, arXiv:2405.07896. | A.1 Full Table | 1670 |
| 1617 | | A.2 LLM Usage | 1671 |
| 1618 | | We used LLMs for language polishing and organization only; all technical claims and citations were authored and verified by the authors | 1672 |
| 1619 | | | 1673 |
| 1620 | | | 1674 |
| 1621 | | | |
| 1622 | | | |
| 1623 | Guangtao Zeng, Wenmian Yang, Zeqian Ju, Yue Yang, Sicheng Wang, Ruisi Zhang, Meng Zhou, Jiaqi Zeng, Xiangyu Dong, Ruoyu Zhang, Hongchao Fang, Penghui Zhu, Shu Chen, and Pengtao Xie. 2020. Meddialog: Large-scale medical dialogue datasets . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9241–9250, Online. Association for Computational Linguistics. | | |
| 1624 | | | |
| 1625 | | | |
| 1626 | | | |
| 1627 | | | |
| 1628 | | | |
| 1629 | | | |
| 1630 | | | |
| 1631 | | | |
| 1632 | N. Zhang, M. Chen, Z. Bi, X. Liang, L. Li, X. Shang, K. Yin, C. Tan, J. Xu, and F. Huang. 2021. Cblue: A chinese biomedical language understanding evaluation benchmark. <i>arXiv preprint arXiv:2106.08087</i> . | | |
| 1633 | | | |
| 1634 | | | |
| 1635 | | | |
| 1636 | Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating text generation with BERT . In <i>International Conference on Learning Representations (ICLR 2020)</i> . | | |
| 1637 | | | |
| 1638 | | | |
| 1639 | | | |
| 1640 | | | |
| 1641 | Y. Zhang, M. Lang, J. Jiang, Z. Gao, F. Xu, T. Litfin, K. Chen, J. Singh, X. Huang, and G. Song. 2024a. Multiple sequence alignment-based rna language model and its application to structural inference. <i>Nucleic Acids Research</i> , 52(1). | | |
| 1642 | | | |
| 1643 | | | |
| 1644 | | | |
| 1645 | | | |
| 1646 | Yubo Zhang, Shudi Hou, Mingyu Derek Ma, Wei Wang, Muhao Chen, and Jieyu Zhao. 2024b. CLIMB: A benchmark of clinical bias in large language models . Preprint, arXiv:2407.05250. | | |
| 1647 | | | |
| 1648 | | | |
| 1649 | | | |
| 1650 | Zhengyun Zhao, Qiao Jin, Fangyuan Chen, Tuorui Peng, and Sheng Yu. 2023. A large-scale dataset of patient summaries for retrieval-based clinical decision support systems . <i>Scientific Data</i> , 10(1):909. | | |
| 1651 | | | |
| 1652 | | | |
| 1653 | | | |
| 1654 | Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena . In <i>NeurIPS 2023 Datasets and Benchmarks Track</i> . | | |
| 1655 | | | |
| 1656 | | | |
| 1657 | | | |
| 1658 | | | |
| 1659 | | | |
| 1660 | Ben Zhou, Hongming Zhang, Sihao Chen, Dian Yu, Hongwei Wang, Baolin Peng, Dan Roth, and Dong Yu. 2024. Conceptual and unbiased reasoning in language models . Preprint, arXiv:2404.00205. | | |
| 1661 | | | |
| 1662 | | | |
| 1663 | | | |
| 1664 | Lianghui Zhu, Xinggang Wang, and Xinlong Wang. 2025. Judgelm: Fine-tuned large language models | | |
| 1665 | | | |

Table 3: Expanded benchmarks by task class. Each row lists what the task can establish for clinical readiness, common evaluation metrics, the target level (L0–L3), and representative benchmarks.

| Task | Usefulness | Metric | Level | Benchmark |
|---------------------------------|---|--|-------|---|
| Knowledge recall / exam Q&A | Breadth of factual recall and clinical reasoning on multiple-choice or open-ended exam questions across specialties and difficulty levels. Evaluates recall under constrained formats and coverage gates, and whether models hallucinate when unsure (Arora et al., 2025). | Exact match/F1; accuracy; subject/difficulty slices; contamination checks; calibration on unanswerable or uncertain questions. | L0 | <i>MedMCQA</i> (Pal et al., 2022); <i>MultiMedQA</i> (MedQA, MedMCQA, PubMedQA, MMLU) (Singhal et al., 2023); <i>Mirage RAG suite</i> (MMLU-Med, MedQA-US, PubMedQA, BioASQ Y/N) (Xiong et al., 2024) |
| Summarization / transformation | Fidelity and completeness of summaries or transformations of clinical notes, conversations, or research literature. Tests whether models can produce coherent, structurally complete summaries without omissions or hallucinations (DeYoung et al., 2021; Aali et al., 2025). | Hallucination/omission rate; ROUGE/BERTScore/chrF; section completeness; clinical correctness; expert ratings. | L1 | <i>MS²</i> (DeYoung et al., 2021); <i>MIMIC-IV-BHC</i> (Aali et al., 2025); <i>MTS-Dialog</i> (Ben Abacha et al., 2023); <i>ACI-Bench</i> (Yim et al., 2023) |
| Retrieval-augmented QA | Attribution and faithfulness of answers to retrieved sources and freshness/recency of information. Evaluates how well models retrieve and ground answers in relevant documents (Xiong et al., 2024). | Faithfulness/attribution; source contradiction rate; Recall@k, nDCG, MRR; answer correctness; freshness. | L0,L1 | <i>Mirage RAG benchmark</i> (MMLU-Med, MedMCQA, PubMedQA, BioASQ) (Xiong et al., 2024); <i>HealthSearchQA</i> (part of <i>MultiMedQA</i>) (Singhal et al., 2023) |
| Evidence based fact checking | Reliability of claims and ability to verify or refute medical statements using evidence. Useful for ensuring LLM outputs do not propagate misinformation. | Claim classification accuracy; evidence recall/precision; F1 for true/false/unfounded labels; citation quality. | L0,L1 | <i>MedFact</i> (Chen et al., 2025b) |
| Information extraction / coding | Structured accuracy on entity recognition, relation extraction, coding, and normalization tasks. Establishes ability to extract structured data from unstructured texts (Luo et al., 2022; Li et al., 2016; Uzuner et al., 2011; Zhang et al., 2021). | Mention/cluster F1; relation F1; coding/normalization accuracy; entity linking accuracy. | L1 | <i>BioRED</i> (Luo et al., 2022); <i>BC5CDR</i> (Li et al., 2016); <i>n2c2 2010 (i2b2)</i> (Uzuner et al., 2011); <i>CBLUE</i> (Zhang et al., 2021) |

Table 3: Expanded benchmarks by task class (continued).

| Task | Usefulness | Metric | Level | Benchmark |
|--|--|---|-------|--|
| Decision support / triage (simulation) | Selective reliability for clinical decision making: calibration at deployable thresholds, risk coverage trade-offs, harm proxies, and quantitative reasoning. Includes simulation of triage, diagnosis, personalized diabetes management, and medical calculations (Arora et al., 2025; Bedi et al., 2025a; Cardei et al., 2025; Mehandru et al., 2025; Khandekar et al., 2024). | ECE; Brier; NLL; risk coverage curves; contraindication/near miss rates; accuracy, groundedness, safety, clarity, actionability; MAE for calculations. | L2 | <i>HealthBench</i> (Arora et al., 2025); <i>MedHELM</i> (Bedi et al., 2025a); <i>DexBench</i> (Cardei et al., 2025); <i>ER-Reason</i> (Mehandru et al., 2025); <i>MedCalc-Bench</i> (Khandekar et al., 2024) |
| Clinical dialogue | Communication quality and human factors in multi-turn doctor patient conversations or simulated OSCE interviews. Measures goal completion, uncertainty marking, empathy, and adherence to safety rails; also covers note generation from visit dialogues (Zeng et al., 2020; Yim et al., 2023; Fareez et al., 2022). | Goal completion; uncertainty/hedging tags; rubric-based ratings; guideline-contradiction flags; empathy/communication scores; note-generation quality; ROUGE/BERTScore. | L2,L3 | <i>MedDialog</i> (Chinese/English) (Zeng et al., 2020); <i>MTS-Dialog</i> (Ben Abacha et al., 2023); <i>ACI-Bench</i> (Yim et al., 2023); <i>OSCE simulated interview dataset</i> (Fareez et al., 2022) |
| Multimodal (imaging + text) | Linkage between radiology images and free text reports or classification labels; evaluates image understanding, report generation, and cross-modal retrieval (Johnson et al., 2019; Wang et al., 2017). | Report correctness; finding detection/linking; classification accuracy; precision/recall/F1; bounding box/segmentation metrics. | L2,L3 | <i>MIMIC-CXR</i> , <i>MIMIC-CXR-JPG</i> (Johnson et al., 2019); <i>NIH ChestX-ray</i> (Wang et al., 2017) |
| Patient retrieval | Ability to retrieve relevant literature or similar patient summaries to support clinicians. Tests retrieval quality and ranking of semantically similar patients or articles (Zhao et al., 2023). | Recall@k; nDCG; MRR; patient-similarity accuracy; retrieval precision. | L1 | <i>PMC-Patients</i> (Zhao et al., 2023) |
| Molecular / drug discovery | Validity and diversity of generated molecules and optimization of proxy properties for research settings. | Validity/diversity/novelty proxies; property-optimization success; synthetic accessibility; logP and QED scores. | L3 | <i>MOSES</i> (Polykovskiy et al., 2018); <i>GuacaMol</i> (Brown et al., 2019) |

A.3 Autonomy Escalation Criteria

| Current level | Next level | Example evidence required before escalation |
|---|---|---|
| L0 – Inform | L1 – Information Transformation & Aggregation | Stable performance on medical Q&A and patient-facing explanation benchmarks in the intended specialties; hallucination and unsafe-content rates below predefined thresholds on representative L0 queries; readability and lay appropriateness confirmed by human raters; and risk-controlled selectivity (ability to abstain or defer) for out-of-scope informational queries. |
| L1 – Information Transformation & Aggregation | L2 – Decision Support | High and stable extraction / normalization performance for clinically critical fields (e.g., problems, medications, laboratory results, temporality) on internal and public test sets; coverage of mandatory fields with low rates of clinically important omissions; RAG evaluations showing that responses are grounded in retrieved evidence with low contradiction to source rates; privacy and provenance checks demonstrating acceptable de-identification risk and machine readable lineage; and, for target L2 tasks, pilot studies showing that clinicians can reliably base case workups on the transformed EHR substrate. |
| L2 – Decision Support | L3 – Agents under Human Supervision | All L2 requirements satisfied for the tasks in scope, including calibrated risk coverage with a tuned abstention policy; simulation or shadow deployments showing high end-to-end task success under human oversight for representative workflows; tool-use correctness audits (API choice, parameters, sequencing) with low failure and silent-error rates; clinician-in-the-loop studies indicating acceptable edit / override rates and no strong evidence of harmful automation bias; and complete, queryable audit trails (e.g., FHIR Provenance / AuditEvent) for data access, prompts, model invocations, tool calls, and approvals. |

Table 4: Example autonomy escalation criteria: illustrative evidence requirements for moving a medical LLM system from one autonomy level (L0–L2) to the next (L1–L3). We intentionally do not prescribe specific numerical thresholds here, because acceptable error rates, coverage targets, and abstention rates depend on deployment context, task criticality, baseline performance of existing workflows, and local regulatory or institutional standards. Concrete numerical criteria should therefore be instantiated by deployers in collaboration with clinical, safety, and governance stakeholders and updated over time as evidence, norms, and models evolve.

A.4 Survey Comparison

Table 5: Comparison between our survey and prior medical LLM evaluation surveys focusing on specific aspects of evaluation.

| Survey | Focus of prior survey | Difference from our autonomy-aware survey |
|--|--|---|
| Evaluating large language models and agents in healthcare: Key challenges in clinical applications (Chen et al., 2025c) | Examines the evaluation of large language models and agents in healthcare, emphasizing agentic, tool-using scenarios and challenges in clinical applications. | Provides a detailed discussion of agentic and workflow-oriented use cases, but does not organize evaluations by levels of autonomy or explicitly link metrics to escalation decisions across levels. Our survey structures agentic evaluations along an autonomy ladder and specifies the required evidence for each level. |
| A systematic review of large language model (LLM) evaluations in clinical medicine (Shool et al., 2025) | Systematically reviews evaluations of LLMs in clinical medicine, with a strong emphasis on benchmark and task performance, often centered on accuracy-heavy metrics. | Covers a broad range of clinical evaluations but primarily treats metrics as task-level accuracy and performance measures, without structuring them by autonomy level or mapping them to deployment decisions. We instead embed such metrics within an autonomy-aware evaluation framework. |
| A framework for human evaluation of large language models in healthcare derived from literature review (QUEST) (Tam et al., 2024) | Proposes the QUEST framework for human evaluation of LLMs in healthcare, focusing on dimensions such as quality, utility, and safety of model outputs as judged by human raters. | Provides rich guidance on human evaluation metrics but does not systematically connect these metrics to levels of autonomy or to the escalation of model responsibilities. Our survey incorporates human evaluation as one component of the evidence required at each autonomy level. |
| A Proposed S.C.O.R.E. Evaluation Framework for Large Language Models: Safety, Consensus, Objectivity, Reproducibility and Explainability (Tan et al., 2024) | Introduces a multidimensional framework (S.C.O.R.E.) for evaluating LLMs along safety, consensus, objectivity, reproducibility, and explainability. | Offers a multidimensional view of evaluation but does not explicitly map these dimensions to an autonomy axis or specify what level of S.C.O.R.E. evidence is required before a model can be escalated to higher-impact roles. Our survey directly ties such multidimensional metrics to autonomy levels and escalation criteria. |
| Foundation metrics for evaluating effectiveness of healthcare conversations powered by generative AI (Abbasian et al., 2024) | Defines foundation metrics for evaluating the effectiveness of healthcare conversations powered by generative AI, focusing on conversation quality and interaction outcomes. | Concentrates on conversational metrics rather than autonomy-aware evaluation, and does not map conversational evidence to what the system is allowed to do in clinical workflows. Our framework relates conversational evaluation to specific autonomy levels and associated permissions. |

Table 6: Comparison between our survey and additional related surveys identified in our extended search.

| Survey | Focus of prior survey | Difference from our autonomy-aware survey |
|---|---|--|
| The future landscape of large language models in medicine (Clusmann et al., 2023) | Discusses the broader landscape of how large language models could fit into medicine, including opportunities, risks, and potential future directions for clinical integration. | Provides a broad, forward-looking view of LLMs in medicine but does not map evaluations to specific levels of autonomy, making it difficult to decide when to escalate an LLM from a chatbot to higher-impact roles in real-world applications. Our survey explicitly aligns evaluations and metrics with autonomy levels. |
| A scoping review of large language models for generative tasks in mental health care (Hua et al., 2025) | Reviews the use of LLMs for generative tasks in mental health care, highlighting applications, challenges, and issues such as non-standard and ad-hoc evaluation scales. | Identifies important evaluation problems (e.g., heterogeneous and ad-hoc scales) but does not specify what concrete evidence is sufficient to deploy a system in real mental health settings. We instead specify the necessary evaluations for each task within each autonomy level. |
| The application of large language models in medicine: A scoping review (Meng et al., 2024) | Synthesizes 550 studies on LLMs in medicine, characterizing the gap between benchmark-focused evaluations and real-world deployment across many clinical and research settings. | Clearly documents the benchmark–deployment gap but does not provide a concrete plan to close it. In contrast, we pair evaluation metrics with corresponding checks for each autonomy level, turning this gap into a deployable test plan. |
| Large language models in medicine (Thirunavukarasu et al., 2023) | Surveys opportunities, risks, and potential applications of LLMs in medicine, with emphasis on clinical use cases and safety concerns. | Discusses opportunities and risks across medical research and practice but does not specify evaluations tied to what the model is allowed to do in clinical workflows. Our framework links required evidence directly to permitted actions at each autonomy level. |
| A Systematic Review of Large Language Models in Medical Specialties: Applications, Challenges and Future Directions (Alkalbani et al., 2025) | Systematically reviews LLM applications across medical specialties, reporting heterogeneous results and emphasizing the need for standards in evaluation and deployment. | Catalogs applications and argues that standards are necessary, but does not describe which additional evidence is required at each level of autonomy. We identify the evidence and checks required before escalation between levels. |
| A comprehensive survey of large language models and multimodal large language models in medicine (Xiao et al., 2025) | Offers general guidance on medical LLMs and MLLMs, including architectures, datasets, applications, and evaluation, with emphasis on multimodal capabilities. | Provides broad evaluation guidance for LLMs and MLLMs but is not autonomy-aware and does not indicate what evidence must be demonstrated before escalation. Our survey is explicitly structured around autonomy levels and escalation criteria. |
| Large Language Models in Medicine: Applications, Challenges, and Future Directions (Yu et al., 2025) | Reviews applications, challenges, and future directions of LLMs in medicine, and calls for standardized benchmarks and multimodal evaluations. | Calls for standardized benchmarks and multimodal evaluations, but does not provide a schema that specifies what additional evidence is needed for each level of autonomy. In contrast, we specify which checks are necessary at each level and how they relate to real-world deployment decisions. |

| | | | |
|------|---|------|--|
| 1677 | A.5 Case Study & Worked example: | | |
| 1678 | Evaluating an EHR copilot across | | |
| 1679 | autonomy levels | | |
| 1680 | In this worked example, we show how the same | | |
| 1681 | high-level application—an EHR copilot for outpa- | | |
| 1682 | tient anticoagulant management—can be instan- | | |
| 1683 | tiated and evaluated at different autonomy levels | | |
| 1684 | (L0–L3). The goal is to make concrete how our | | |
| 1685 | autonomy lens guides (i) which actions the system | | |
| 1686 | is allowed to take, and (ii) which evidence and met- | | |
| 1687 | rics are required at each level, building directly on | | |
| 1688 | the definitions in Section 4. | | |
| 1689 | Clinical scenario. A primary care clinician is | | |
| 1690 | starting a new oral anticoagulant for a patient with | | |
| 1691 | non-valvular atrial fibrillation. The overall work- | | |
| 1692 | flow includes (1) educating the patient about the | | |
| 1693 | medication, (2) reviewing the chart for contraindi- | | |
| 1694 | cations and drug to drug interactions, (3) proposing | | |
| 1695 | a dosing and monitoring plan, and (4) drafting or- | | |
| 1696 | ders and patient-facing instructions in the EHR for | | |
| 1697 | clinician sign-off. We consider four configurations | | |
| 1698 | of an EHR copilot for this scenario, one per auton- | | |
| 1699 | omy level. | | |
| 1700 | L0: Inform-only patient education | | |
| 1701 | At autonomy level L0, the copilot acts purely as | | |
| 1702 | an informational tool. It does <i>not</i> access patient- | | |
| 1703 | specific EHR data and does <i>not</i> give individualized | | |
| 1704 | recommendations. Instead, it: | | |
| 1705 | • answers general questions such as “What is | | |
| 1706 | an anticoagulant?”, “Why is this medication | | |
| 1707 | prescribed?”, or “What are common side ef- | | |
| 1708 | fects?”; | | |
| 1709 | • rewrites clinician-authored discharge instruc- | | |
| 1710 | tions into plain language; | | |
| 1711 | • explains standard lab tests (e.g., “What does a | | |
| 1712 | kidney function test measure?”). | | |
| 1713 | The outputs are generic and include explicit dis- | | |
| 1714 | claimers (e.g., “This is general information and not | | |
| 1715 | personal medical advice.”). | | |
| 1716 | Evaluation at L0. For this configuration, evalua- | | |
| 1717 | tion targets match our L0 focus in Section 4.1: | | |
| 1718 | • Factual accuracy and completeness: ques- | | |
| 1719 | tion answering accuracy on consumer-health | | |
| 1720 | datasets (e.g., HealthSearchQA-style corpora), | | |
| 1721 | coverage of key points in patient education | | |
| 1722 | materials. | | |
| | • Readability and organization: readability | 1723 | |
| | indices and human ratings of clarity, struc- | 1724 | |
| | ture, and appropriateness of language for lay | 1725 | |
| | readers. | 1726 | |
| | • Safety and tone: rubric-based assessment | 1727 | |
| | of whether answers avoid giving individu- | 1728 | |
| | alized advice, avoid discouraging guideline- | 1729 | |
| | concordant care, and encourage seeking pro- | 1730 | |
| | fessional help when appropriate. | 1731 | |
| | No calibration, personalization, or tool-use met- | 1732 | |
| | rics are required, because the system does not see | 1733 | |
| | patient data and cannot act on the record. | 1734 | |
| | L1: Information transformation and | 1735 | |
| | aggregation over the EHR | 1736 | |
| | At autonomy level L1, the copilot is allowed to | 1737 | |
| | read the patient’s EHR but remains descriptive and | 1738 | |
| | non-prescriptive. It can: | 1739 | |
| | • extract and normalize key facts from the chart | 1740 | |
| | (problems, active medications, allergies, re- | 1741 | |
| | cent labs, relevant imaging reports) into a stan- | 1742 | |
| | dardized schema (e.g., FHIR resources or an | 1743 | |
| | OMOP-style common data model); | 1744 | |
| | • generate a concise clinical summary highlight- | 1745 | |
| | ing bleeding and thrombotic risk factors (e.g., | 1746 | |
| | prior gastrointestinal bleeding, renal function, | 1747 | |
| | age, concomitant antiplatelets); | 1748 | |
| | • retrieve and surface relevant guideline snip- | 1749 | |
| | pets or drug monographs via retrieval- | 1750 | |
| | augmented generation, with inline citations. | 1751 | |
| | The copilot presents <i>what is known</i> and <i>what the</i> | 1752 | |
| | <i>literature says</i> but does not propose a concrete plan | 1753 | |
| | for this patient. | 1754 | |
| | Evaluation at L1. Relative to L0, evaluation now | 1755 | |
| | needs to cover transformation quality, aggregation, | 1756 | |
| | and provenance, as outlined in Section 4.2: | 1757 | |
| | • Extraction and normalization: span-level | 1758 | |
| | precision/recall/ F_1 for entities (problems, | 1759 | |
| | medications, allergies, labs) on de-identified | 1760 | |
| | note benchmarks; concept mapping accuracy | 1761 | |
| | to standard vocabularies. | 1762 | |
| | • Retrieval and aggregation: retrieval metrics | 1763 | |
| | (e.g., Recall@k, MRR, nDCG) for guideline | 1764 | |
| | and evidence search, plus grounded-answer | 1765 | |
| | and contradiction-to-source rates for gener- | 1766 | |
| | ated summaries. | 1767 | |

- logging provenance (which data, which model, which tools) for each drafted action.

Crucially, all drafted actions are *queued for explicit clinician review and sign-off*. The agent is not allowed to execute orders or send messages autonomously; it operates “under supervision” as in our L3 definition.

Evaluation at L3. For L3, the primary target shifts to supervised action quality and human–AI interaction, as discussed in Section 4.4:

- **End-to-end task success:** in simulation or shadow mode, the proportion of encounters where the agent-produced plan and drafted actions are judged clinically acceptable, and the degree to which they reduce manual workload (e.g., reduction in time spent on order entry and documentation).
- **Tool/API correctness:** validity of API calls (correct endpoints, parameters, and patient identifiers), syntactic correctness (e.g., FHIR-conformant requests), and semantic correctness (the data retrieved or orders drafted match the intended patient and plan).
- **Oversight effectiveness:** rates at which clinicians accept, modify, or reject drafted actions; qualitative or quantitative signals of over-reliance (rubber-stamping) or alert fatigue; changes in error rates with versus without the copilot.
- **Auditability and governance:** completeness of machine-readable provenance (e.g., FHIR AuditEvent and Provenance resources) that tie each action to its inputs, prompts, model versions, tools, and human approvals, enabling post-hoc investigation of failures or near-misses.
- **Operational safety:** monitoring for rare but high-impact failure modes in pilot deployments (e.g., mis-addressed messages, wrong-patient orders) and documenting mitigation strategies (access controls, guardrails, kill-switches).

These measurements go beyond traditional benchmark scores. They are required to justify even supervised L3 use, because errors now manifest as mis-drafted actions rather than just incorrect text.

Summary of the worked example

This worked example illustrates how autonomy levels structure both system design and evaluation:

- L0 constrains the copilot to generic education, evaluated on accuracy, readability, and safe communication.
- L1 adds EHR-based transformation and aggregation, evaluated on extraction quality, completeness, retrieval, and provenance.
- L2 introduces patient-specific recommendations, evaluated on correctness, reasoning, calibration, safety, and fairness.
- L3 adds tool use under supervision, evaluated on supervised action quality, tool correctness, oversight, auditability, and operational safety.

Framing the same anticoagulant-management copilot across L0–L3 shows how our autonomy taxonomy operates as a practical blueprint: it makes the system’s allowed actions explicit, and aligns those actions with level-appropriate evidence and metrics, turning benchmark performance into credible, risk-aware claims about clinical utility.