

From Guidelines to Guarantees: A Graph-Based Evaluation Harness for Domain-Specific Evaluation of LLMs

Anonymous ACL submission

Abstract

Rigorous evaluation of domain-specific language models requires benchmarks that are comprehensive, contamination-resistant, and maintainable. Static, manually curated datasets do not satisfy these properties. We present a graph-based evaluation harness that transforms structured clinical guidelines into a queryable knowledge graph and dynamically instantiates evaluation queries via graph traversal. The framework provides three guarantees: (1) complete coverage of guideline relationships; (2) surface-form contamination resistance through combinatorial variation; and (3) validity inherited from expert-authored graph structure. Applied to the WHO IMCI guidelines, the harness generates clinically grounded multiple-choice questions spanning symptom recognition, treatment, severity classification, and follow-up care. Evaluation across five language models reveals systematic capability gaps. Models perform well on symptom recognition but show lower accuracy on treatment protocols and clinical management decisions. The framework supports continuous regeneration of evaluation data as guidelines evolve and generalizes to domains with structured decision logic. This provides a scalable foundation for evaluation infrastructure.

Data and Code Availability The WHO IMCI handbook is publicly available (WHO, 2014). Our graph construction, question generation code, and generated question dataset are available at [ANONYMIZED GITHUB LINK] for review. Upon acceptance, we will release the code and dataset publicly under an open-source license to enable reproducibility and extension to other clinical guidelines.

1 Introduction

1.1 The Evaluation Coverage Problem

Rigorous evaluation of language models faces a critical challenge: the distribution gap between

application-specific text and existing benchmark datasets. This gap encompasses both context (domain, localization, complexity) and coverage (tasks, content). Current medical benchmarks rely on human curation, which is resource-intensive and results in incomplete coverage of specific medical guidelines.

MCQA benchmark datasets serve dual purposes: training new models and evaluating across models. The test split has widespread utility as a yardstick for comparison across models. While vignettes and multi-turn conversation with evaluation rubrics (Tu et al., 2024; Nori et al., 2025; Arora et al., 2025) more closely resemble real-world scenarios, MCQA remains an important evaluation format because it is less ambiguous, easy to grade, and scalable.

Despite advances in model architectures and training paradigms, MCQA benchmarks remain central for both evaluation and post-training. In health-domain models, supervised finetuning continues to be useful. Within alignment, MCQA also provides naturally ranked outputs for methods such as GRPO, where correct answers serve as high-reward samples and incorrect options serve as progressively lower-reward samples without requiring expensive human ranking.

WHO guidelines are an appropriate use case for this setting because there is substantial need for AI systems that support scarce healthcare workers in low- and middle-income countries (LMICs). These guidelines are often country-specific, which makes custom evaluation necessary for accurate measurement of model performance.

1.2 Limitations of Existing Medical Benchmarks

Medical benchmarks exist in multiple languages, and rely on questions from licensing exams, textbooks, journals, and crowdsourcing (Jin et al., 2021; Pal et al., 2022; Vilares and Gómez-

083	Rodríguez, 2019; Labrak et al., 2022; Kasai et al.,	• FollowUp (15 nodes): Monitoring schedules	127
084	2023; Jin et al., 2019; Zhang et al., 2017; Olatunji	(e.g., “3 days”, “7 days”)	128
085	et al., 2024; Hendrycks et al., 2021; Alonso et al.,	• Severity (4 nodes): Triage classifications (se-	129
086	2024). Synthetic medical QA datasets employ	vere, moderate, mild, none)	130
087	diverse generation strategies: template-based ap-	Four edge types connect these nodes:	131
088	proaches as in emrQA (Pampari et al., 2018) and	• INDICATES : Symptom → Condition	132
089	RadQA (Soni et al., 2022), generation using ontol-	• TREAT : Condition → Treatment	133
090	ogy concepts (Dong et al., 2023), and LLM-based	• FOLLOW : Condition → FollowUp	134
091	generation for hallucination detection (Pal et al.,	• TRIAGE : Condition → Severity	135
092	2023).	Automated extraction via PDF parsers and	136
093	1.3 Contributions	LLMs failed to reliably capture the conditional	137
094	Our main contributions are as follows:	logic embedded in IMCI flowcharts. Relationships	138
095	1. We introduce a graph-based evaluation harn-	expressed visually through color-coded triage paths	139
096	ness that provides explicit guarantees of cov-	and nested decision branches cannot be faithfully	140
097	erage, contamination resistance, and validity.	reconstructed as directed edges by current PDF and	141
098	2. We present a method for transforming struc-	LLM pipelines. The knowledge graph was there-	142
099	tured clinical guidelines into a knowledge	fore manually curated by a co-author with over 15	143
100	graph that supports systematic evaluation.	years of clinical practice, specialized pediatric train-	144
101	3. We demonstrate dynamic evaluation through	ing, and extensive experience implementing WHO	145
102	on-demand query instantiation rather than	IMCI guidelines in sub-Saharan Africa. This clini-	146
103	static datasets.	cal authorship of the graph establishes validity at	147
104	4. We empirically show that this framework re-	the source: all generated questions inherit their ac-	148
105	veals systematic weaknesses in clinical reason-	curacy from expert-constructed relationships rather	149
106	ing that are not captured by aggregate bench-	than requiring post-hoc review of generated out-	150
107	marks.	puts.	151
108	2 Method	2.2 Evaluation Query Instantiation	152
109	2.1 Graph Construction from Clinical	We employ graph traversal to automatically instan-	153
110	Guidelines	tiate MCQA evaluation queries that ensure com-	154
111	We transform the WHO IMCI handbook (WHO,	plete coverage of medical relationships. For each	155
112	2014) into a directed graph structure. The hand-	condition node, we traverse its connected nodes to	156
113	book, an 80-page document containing flowcharts	instantiate the five question types shown in Table 1.	157
114	and checklists for childhood illness management,	The framework dynamically instantiates evalu-	158
115	is parsed to extract medical entities and their rela-	ation queries using four templates for each of five	159
116	tionships. The resulting graph contains 200+ nodes	question types while maintaining clinical relevance	160
117	and 300+ edges spanning respiratory, gastrointesti-	and variability. Random age generation is con-	161
118	nal, and infectious diseases.	strained to the condition’s valid range (e.g., 0–8	162
119	The graph schema consists of five node types:	weeks for young infants, 2–60 months for chil-	163
120	• Condition (31 nodes): Medical conditions	dren).	164
121	with age range attributes (0–2 months for	The distractor sampling algorithm prioritizes	165
122	young infants, 2–60 months for children)	clinical validity through age-stratified selection.	166
123	• Symptom (79 nodes): Observable clinical in-	For each question requiring $k = 3$ distractors, the	167
124	dicators (e.g., “fast breathing”, “convulsions”)	system first identifies all conditions sharing the	168
125	• Treatment (84 nodes): Medical interventions	same age range as the target condition, creating an	169
126	(e.g., “give oral Amoxicillin for 5 days”)	age-appropriate candidate pool.	170
		For a question with correct answer v_{corr} of type	171
		τ and target condition with age range α , we con-	172
		struct an age-appropriate distractor pool by select-	173
		ing candidate nodes that (i) match the required type	174

Table 1: Examples of auto-generated questions by relationship type.

Type	Example
Condition → Symptom	Q: A 2 year old child with Very Severe Disease would most likely present with which symptom? Options: A: convulsions, B: chest indrawing, C: pus draining from the eye, D: WFH/L 2 z-scores or more Answer: A
Symptom → Condition	Q: A 21 month old child presenting with convulsions is most likely to have: Options: A: Cough or Cold, B: Very Severe Disease, C: Severe Pneumonia or Very Severe Disease, D: Very Severe Febrile Disease with no Malaria Risk Answer: B
Condition → Treatment	Q: Which treatment is recommended for a 21 month old child with Very Severe Disease? Options: A: assess or refer for TB assessment and INH preventive therapy, B: if mouth ulcers treat with gentian violet, C: do virological test at age 4–6 weeks or repeat 6 weeks after the child stops breastfeeding, D: give first dose of intramuscular antibiotics Answer: D
Condition → FollowUp	Q: What is the appropriate follow-up schedule for a 3 year old child with Some Dehydration? Options: A: follow-up in 14 days, B: follow-up in 5 days, C: follow-up in 2 days if not improving, D: follow-up in 7 days Answer: C
Condition → Severity	Q: A 13 month old child with Very Severe Disease should be classified as: Options: A: moderate, B: mild, C: none, D: severe Answer: D

and (ii) are compatible with the target age range. Distractors are then sampled uniformly without replacement from this pool.

This construction ensures that all distractors are clinically plausible within the relevant age group while maintaining variability across generated questions. A formal specification of the distractor construction is provided in Appendix A.

The dynamic generation process creates novel evaluation instances through variation in templates, ages, and distractors while maintaining consistent difficulty and clinical relevance. This mitigates a key limitation of static benchmarks, in which models may have seen evaluation questions during training, while enabling substantial variation for robust statistical analysis.

2.3 Contamination Resistance

The harness addresses two distinct contamination risks that static benchmarks cannot mitigate.

Surface-form contamination occurs when evaluation questions appear verbatim in training data. By generating questions at evaluation time with randomized ages, distractor sampling, and template selection drawn from a large combinatorial space of possible instances, the probability of repeated surface forms is reduced relative to static benchmarks.

Relationship-level contamination occurs when a model has learned the underlying clinical relation-

ships from source documents, such that it can answer questions correctly regardless of surface form. Unlike surface-form contamination, this cannot be mitigated through variation in phrasing alone.

Rather than attempting to eliminate this form of contamination, the proposed harness enables a complementary evaluation strategy. Because evaluation queries are generated dynamically from a structured representation of the guidelines, the same framework can be applied to updated or modified guidelines that postdate model training. This allows evaluation to probe whether models have genuinely acquired generalizable clinical reasoning or are relying on memorized relationships from specific guideline versions.

In this sense, the harness supports temporal and versioned evaluation, making it possible to identify knowledge gaps as clinical guidelines evolve. This shifts evaluation from static benchmarking to continuously refreshable assessment aligned with evolving domain knowledge.

Graph-level errors represent a third risk, where inaccuracies in the knowledge graph propagate to all generated questions. Expert authorship of the graph (Section 3.1) directly addresses this by establishing the graph as a clinically verified source of evaluation truth.

231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279

3 Case Study: WHO IMCI

3.1 Clinical Expert Authorship and Validation

The knowledge graph underlying all generated questions was manually curated by a co-author who is a board-certified physician with over 15 years of clinical practice, specialized pediatric training, and extensive experience implementing WHO IMCI guidelines in clinical settings in sub-Saharan Africa. This authorship model, where domain expertise is embedded at the graph construction stage rather than applied as post-hoc review, provides stronger validity guarantees than question-level annotation alone: every generated question inherits its clinical accuracy from expert-constructed graph relationships.

To further validate the generated question set, the same expert reviewed the 432 auto-generated questions across the five relationship types: Condition \rightarrow Treatment (130), Symptom \rightarrow Condition (118), Condition \rightarrow Symptom (118), Condition \rightarrow Severity (37), and Condition \rightarrow FollowUp (29). For each question, the review assessed: (1) clinical accuracy of the correct answer, (2) appropriateness of distractors for the specified age range, and (3) clarity and unambiguity of question phrasing. Given that questions are derived from an expert-curated graph, this review serves primarily to verify that the generation pipeline correctly traverses and formats the underlying relationships rather than to establish clinical accuracy de novo.

The graph was curated by a single clinical expert, which precludes inter-rater reliability assessment. The underlying guidelines provide deterministic decision rules, which partially mitigates subjectivity in annotation. Independent validation by additional clinicians with IMCI expertise remains important future work for establishing the rigor required of a production evaluation instrument.

3.2 LLM Inference Results

We conduct baseline inference evaluation to assess out-of-the-box model performance for the closed-source models Claude Sonnet 4.6, o4-mini, and GPT-5.2, the open-weights model GPT-OSS-20B, and the domain fine-tuned model MedGemma-4B. Models receive questions in a standardized format with explicit instructions to respond with only the letter (A, B, C, or D) corresponding to the correct answer. We measure accuracy per question type with uncertainty over the template variations.

Figure 1 and Table 2 present model performance across question types.

Figure 2 presents model performance variations across clinical question types, measured as the delta between question-specific accuracy and overall model accuracy.

3.3 Key Findings

1. The three frontier closed-source models, Claude Sonnet 4.6, o4-mini, and GPT-5.2, achieve similar overall accuracy (approximately 66–68%), outperforming GPT-OSS-20B (approximately 57%) and MedGemma-4B (approximately 50%).
2. Symptom \rightarrow Condition questions show the highest performance across all models (64–82%), indicating that models better recognize symptoms than prescribe treatments or protocols.
3. Within-model performance varies substantially across question types, underscoring that aggregate accuracy obscures meaningful capability differences.
4. MedGemma-4B underperforms larger models across all question types, indicating that model scale and general reasoning capacity may dominate performance in this setting.

Unlike human-curated benchmarks, our dynamic graph-based method ensures complete coverage of all guideline relationships, consistent terminology from source documents, reduced data contamination through automated generation, and scalability to other medical guidelines.

3.4 Template Ablation Study

Figure 3 reveals substantial within-type variance across question templates, demonstrating that phrasing significantly affects model performance independently of the underlying clinical relationship being tested. The `cond_followup_t1` template (“When should a {age} old child with {cond} return for follow-up?”) consistently produces the lowest accuracy across all models (14–57%), while `cond_symp_t3` produces some of the highest (50–90%). This variance has direct implications for evaluation harness design: using multiple templates per question type, as our harness does, provides more robust estimates of model capability than single-template approaches, and averaging over

280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326

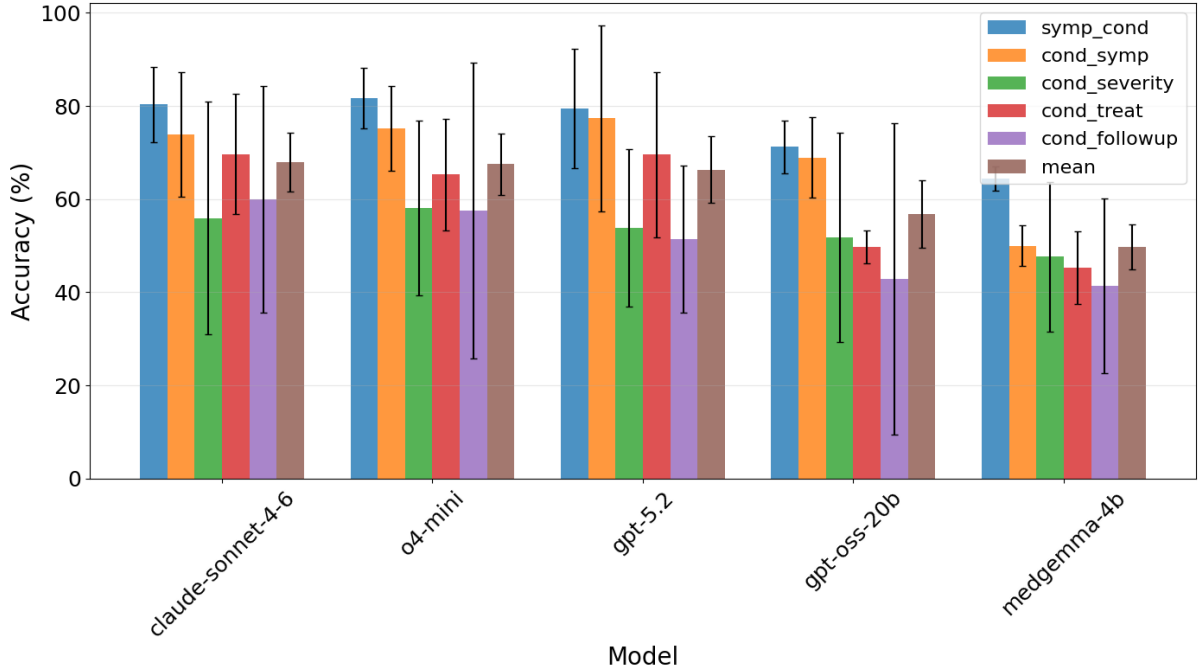


Figure 1: Model accuracy across different question types with 95% confidence intervals across five clinical question categories: symptom-condition (symp_cond), condition-symptom (cond_symp), condition-severity (cond_severity), condition-treatment (cond_treat), and condition-followup (cond_followup). Error bars represent 95% confidence intervals.

Table 2: Model performance comparison across question types. Values shown as accuracy \pm standard deviation (%).

Model	Overall	C→S	S→C	C→T	C→Sv	C→F
Claude Sonnet 4.6	68.0±13.6	73.9±8.4	80.3±5.1	69.7±8.1	56.0±15.7	60.0±15.3
GPT-5.2	66.3±15.2	77.3±12.5	79.4±8.1	69.6±11.1	53.9±10.6	51.4±9.9
o4-mini	67.5±14.0	75.1±5.7	81.6±4.1	65.3±7.5	58.0±11.8	57.5±19.9
GPT-OSS-20B	56.9±15.5	68.9±5.5	71.2±3.6	49.7±2.2	51.8±14.1	42.9±21.0
MedGemma-4B	49.8±10.4	50.0±2.7	64.4±1.6	45.4±4.9	47.6±10.1	41.4±11.8

template variants reduces the influence of phrasing artifacts on reported accuracy.

4 Evaluation Considerations

4.1 Cost and Scalability Relative to Manual Curation

Manual benchmark curation requires domain experts to author, review, and validate each question individually, a process that does not scale and produces static artifacts vulnerable to contamination. Our harness shifts the labor from question authorship to graph construction: a one-time cost that yields a large and refreshable space of evaluation instances for practical evaluation. For IMCI, manual curation of the knowledge graph by a domain clinical expert required significant upfront investment, after which 432 questions were generated automatically with validity inherited from the graph struc-

ture. Expanding across the combinatorial space induced by templates, ages, and distractors requires no additional expert labor beyond graph maintenance as guidelines are updated.

The primary scaling bottleneck is graph construction itself. Automated extraction via PDF parsers and LLMs missed critical relationships because the conditional logic in IMCI flowcharts is expressed visually through color-coded triage paths and nested decision branches that current pipelines cannot faithfully reconstruct as directed edges. Future work could reduce this bottleneck through semi-automated graph construction with expert review, particularly for guidelines with consistent structure such as WHO protocols.

4.2 Stakeholder Roles

The harness separates evaluation into three distinct stakeholder roles with different expertise require-

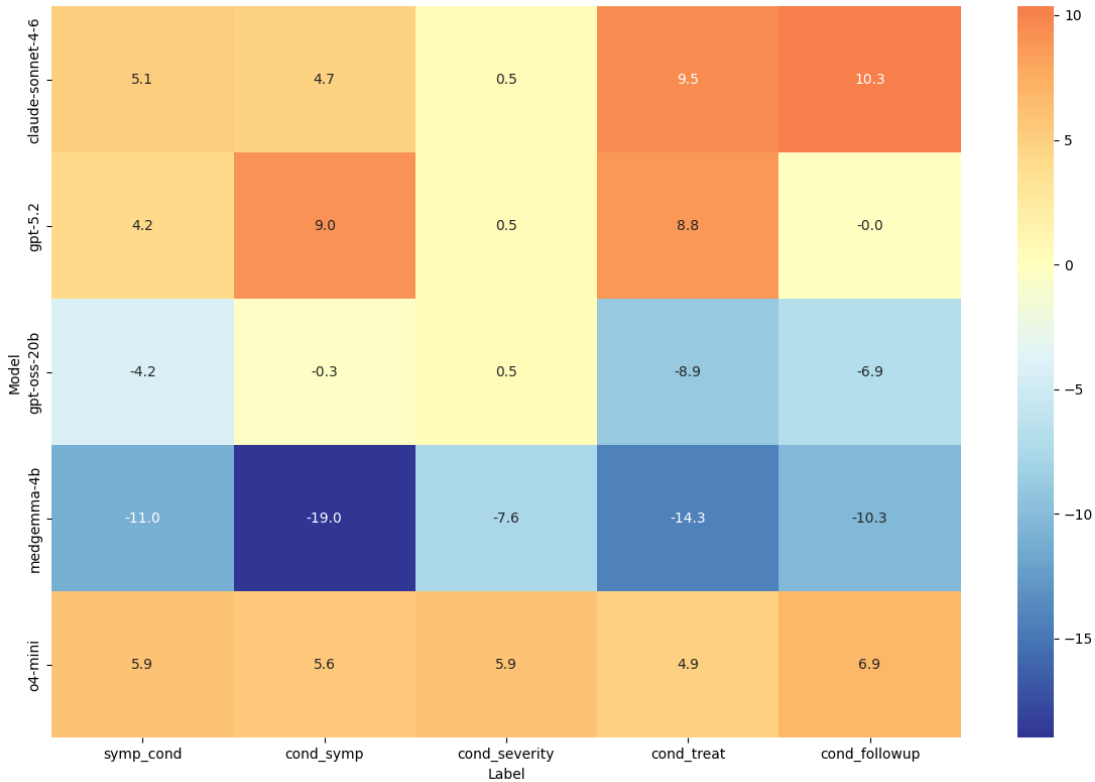


Figure 2: Accuracy delta heatmap showing the difference between question-type-specific accuracy and overall model accuracy for each model. Positive values (red/orange) indicate above-average performance for that question type, while negative values (blue) indicate below-average performance. Values are expressed as percentage points.

362 ments. *Graph constructors* require deep domain ex- 363
 364 pertise to accurately encode guideline relationships; 365
 366 in our case, a pediatrician with IMCI implemen- 367
 368 tation experience in sub-Saharan Africa. *Harness* 369
 370 *operators* require technical expertise to run gen- 371
 372 eration and evaluation pipelines but not medical 373
 374 knowledge. *Model developers* can consume evalua- 375
 376 tion results without access to the underlying graph, 377
 378 enabling third-party evaluation with separation be- 379
 380 tween evaluators and developers, a property the 381
 382 EvalEval community has identified as important 383
 384 for accountability (Reuel et al., 2025).

This separation also clarifies accountability: er- 374
 375 rors in evaluation results can be traced to graph 376
 377 inaccuracies (domain expert responsibility), gen- 378
 379 eration bugs (harness operator responsibility), or 380
 381 model failures (developer responsibility).

379 4.3 Extensibility to Other Guidelines

380 The graph schema, including conditions, symp- 381
 382 toms, treatments, follow-ups, severities, and their 383
 384 directed relationships, is not specific to IMCI. Any 385
 clinical guideline with structured decision logic 386
 is a candidate. WHO produces guidelines across 387
 malaria, tuberculosis, HIV, and maternal health that 388
 share the same flowchart structure as IMCI. Beyond 389
 healthcare, structured regulatory guidelines, legal 390
 compliance frameworks, and technical standards 391
 with explicit relationship structures could support 392
 the same approach. The primary requirement is 393
 that the source document encodes relationships ex- 394
 plicitly enough to support graph construction, a 395
 property common to clinical and regulatory guide- 396
 lines by design. 397

malaria, tuberculosis, HIV, and maternal health that 385
 share the same flowchart structure as IMCI. Beyond 386
 healthcare, structured regulatory guidelines, legal 387
 compliance frameworks, and technical standards 388
 with explicit relationship structures could support 389
 the same approach. The primary requirement is 390
 that the source document encodes relationships ex- 391
 plicitly enough to support graph construction, a 392
 property common to clinical and regulatory guide- 393
 lines by design. 394

395 4.4 Limitations

396 Question quality depends entirely on graph accu- 397
 398 racy: any errors in manual annotation propagate to 399
 400 all generated questions. The graph was curated by 401
 402 a single clinical expert, which precludes inter-rater 403
 404 reliability assessment; independent validation by 405
 406 additional clinicians remains important future work. 407
 We evaluate only MCQA format, which cannot 408
 capture the complexity of real clinical reasoning 409
 involving differential diagnosis and incomplete in- 410
 formation. Our text-only approach excludes visual 411
 diagnostic elements present in the original IMCI 412
 handbook. While question generation is automated, 413

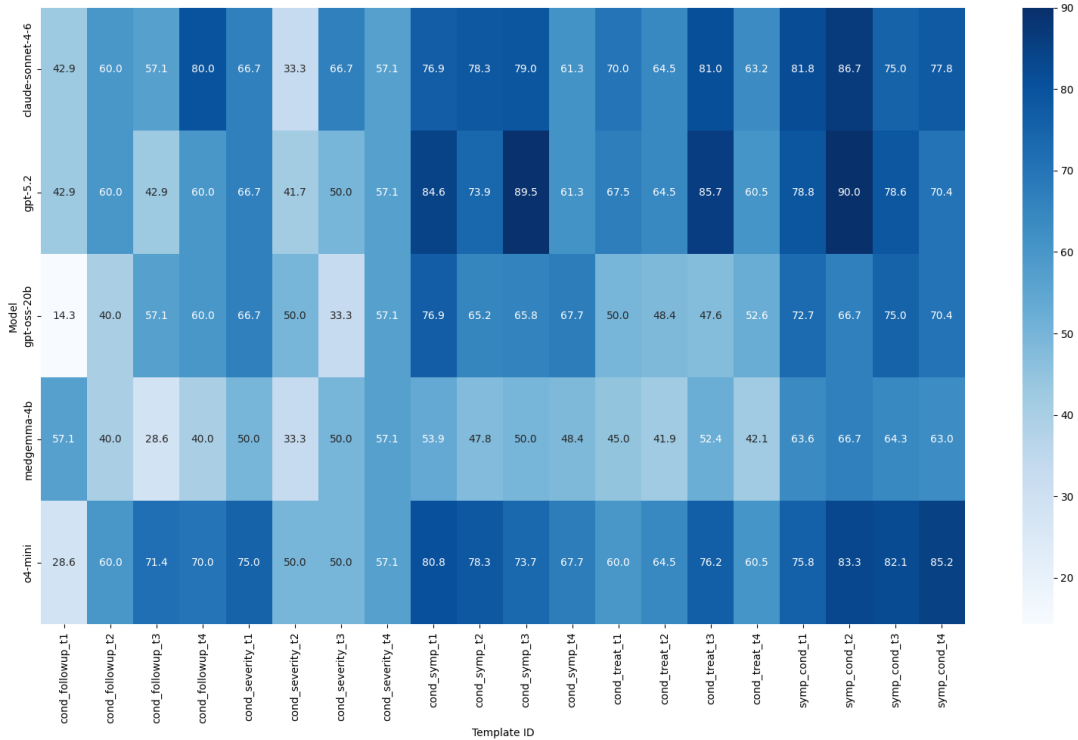


Figure 3: Accuracy by template and model. Each cell shows the accuracy (%) for a given model and question template. Darker blue indicates higher accuracy. Substantial within-type variance across templates demonstrates that question phrasing affects model performance independently of the underlying clinical relationship.

initial graph construction remains manual, limiting scalability. Our evaluation on IMCI guidelines may not generalize to other medical domains. Although the framework admits a large combinatorial space of possible instances, the practically valid subset is smaller because clinical constraints introduce dependencies among age, condition, and distractor choices, and we have not exhaustively verified all such variants.

4.5 Potential Risks

This work presents evaluation tools for medical AI systems. Models performing well on MCQA may still fail in actual clinical scenarios requiring differential diagnosis and incomplete information. Any errors in manual graph annotation propagate to evaluation, potentially validating incorrect medical knowledge. Our focus on WHO IMCI guidelines may not generalize to other healthcare contexts. This evaluation harness is intended for research purposes only and is not suitable for clinical decision-making.

5 Conclusion

This work introduces a graph-based evaluation harness for systematically instantiating evaluation

queries from clinical guidelines, demonstrated on the WHO IMCI handbook. By transforming medical guidelines into queryable graphs, the framework achieves complete coverage of encoded relationships, which is not feasible through manual curation alone. Its dynamic design allows new evaluation instances with different ages and distractors to be sampled continuously, including as guidelines are updated. While baseline inference provides initial scores, the main value lies in granular performance across relationship types, which reveals systematic strengths and weaknesses in clinical protocol understanding.

The clinical validity of the generated questions rests on expert authorship of the underlying graph rather than post-hoc sampling, a design choice that both strengthens the validity claim and clarifies the role of domain expertise in evaluation infrastructure. The graph-based approach is extensible beyond IMCI, addressing the gap between general-purpose benchmarks and real-world domain-specific applications.

References

- 455 Iñigo Alonso, Maite Oronoz, and Rodrigo Agerri. 2024. [Medexpqa: Multilingual benchmarking of large language models for medical question answering](#). *Artificial Intelligence in Medicine*, 155:102938.
- 456
- 457
- 458
- 459 Rahul K. Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, Johannes Heidecke, and Karan Singhal. 2025. [Healthbench: Evaluating large language models towards improved human health](#). *Preprint*, arXiv:2505.08775.
- 460
- 461
- 462
- 463
- 464
- 465
- 466 Hang Dong, Jiaoyan Chen, Yuan He, and Ian Horrocks. 2023. [Ontology enrichment from texts: A biomedical dataset for concept discovery and placement](#). In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM '23*, page 5316–5320, New York, NY, USA. Association for Computing Machinery.
- 467
- 468
- 469
- 470
- 471
- 472
- 473 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Proceedings of ICLR*.
- 474
- 475
- 476
- 477 Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.
- 478
- 479
- 480
- 481
- 482 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. In *Proceedings of EMNLP-IJCNLP*, pages 2567–2577.
- 483
- 484
- 485
- 486 Jungo Kasai, Yuhei Kasai, Keisuke Sakaguchi, Yutaro Yamada, and Dragomir Radev. 2023. [Evaluating gpt-4 and chatgpt on japanese medical licensing examinations](#). *Preprint*, arXiv:2303.18027.
- 487
- 488
- 489
- 490 Yanis Labrak, Adrien Bazoge, Richard Dufour, Mickael Rouvier, Emmanuel Morin, Béatrice Daille, and Pierre-Antoine Gourraud. 2022. Frenchmedmcqa: A french multiple-choice question answering dataset for medical domain. In *Proceedings of the 13th International Workshop on Health Text Mining and Information Analysis (LOUHI)*, pages 41–46, Abu Dhabi, United Arab Emirates.
- 491
- 492
- 493
- 494
- 495
- 496
- 497
- 498 Harsha Nori, Mayank Daswani, Christopher Kelly, Scott Lundberg, Marco Tulio Ribeiro, Marc Wilson, Xiaoxuan Liu, Viknesh Sounderajah, Jonathan Carlson, Matthew P Lungren, Bay Gross, Peter Hames, Mustafa Suleyman, Dominic King, and Eric Horvitz. 2025. [Sequential diagnosis with language models](#). *Preprint*, arXiv:2506.22405.
- 499
- 500
- 501
- 502
- 503
- 504
- 505 Tobi Olatunji, Abraham Owodunni, Tassallah Abdullahi, Ayokunmi Ilesanmi, Olalekan Obadun, Aimérou Ndiaye Etori, Ifeoma Okoh, Evans Doe Ocansey, Wendy Kinara, Michael Best, and 1 others. 2024. Afrimedqa: A pan-african, multi-specialty, medical question-answering benchmark dataset. *arXiv preprint arXiv:2411.15640*.
- 506
- 507
- 508
- 509
- 510
- 511
- 512 Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Conference on Health, Inference, and Learning*, pages 248–260. PMLR.
- 513
- 514
- 515
- 516
- 517 Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2023. [Med-halt: Medical domain hallucination test for large language models](#). *Preprint*, arXiv:2307.15343.
- 518
- 519
- 520
- 521 Anusri Pampari, Preethi Raghavan, Jennifer Liang, and Jian Peng. 2018. emrqa: A large corpus for question answering on electronic medical records. In *Proceedings of EMNLP*, pages 2357–2368.
- 522
- 523
- 524
- 525 Anka Reuel, Avijit Ghosh, Jenny Chim, Andrew Tran, Yanan Long, Jennifer Mickel, Usman Gohar, Srishti Yadav, Pawan Sasanka Ammanamanchi, Mowafak Allaham, Hossein A. Rahmani, Mubashara Akhtar, Felix Friedrich, Robert Scholz, Michael Alexander Riegler, Jan Batzner, and 1 others. 2025. Who evaluates AI’s social impacts? Mapping coverage and gaps in first and third party evaluations. *arXiv preprint arXiv:2511.05613*.
- 526
- 527
- 528
- 529
- 530
- 531
- 532
- 533
- 534 Sarvesh Soni, Meghana Gudala, Atieh Pajouhi, and Kirk Roberts. 2022. Radqa: A question answering dataset to improve comprehension of radiology reports. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6250–6259, Marseille, France. European Language Resources Association.
- 535
- 536
- 537
- 538
- 539
- 540
- 541 Tao Tu, Anil Palepu, Mike Schaeckermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Nenad Tomasev, Shekoofeh Azizi, Karan Singhal, Yong Cheng, Le Hou, Albert Webson, Kavita Kulkarni, S Sara Mahdavi, Christopher Semturs, Juraj Gottweis, and 6 others. 2024. Towards conversational diagnostic ai. *Nature*, 629(8010):331–338.
- 542
- 543
- 544
- 545
- 546
- 547
- 548
- 549 David Vilares and Carlos Gómez-Rodríguez. 2019. [HEAD-QA: A healthcare dataset for complex reasoning](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 960–966, Florence, Italy. Association for Computational Linguistics.
- 550
- 551
- 552
- 553
- 554
- 555 WHO. 2014. Integrated management of childhood illness - chart booklet. Technical report, World Health Organization. Technical document.
- 556
- 557
- 558 Sheng Zhang, Xin Zhang, Hui Wang, Jiajun Cheng, Pei Li, and Zhaoyun Ding. 2017. [Chinese medical question answer matching using end-to-end character-level multi-scale cnns](#). *Applied Sciences*, 7(8).
- 559
- 560
- 561

562 A Distractor Pool Construction

563 We formalize distractor construction for complete-
 564 ness. Let $G = (V, E)$ denote the IMCI knowledge
 565 graph.

566 For a question with correct answer v_{corr} of type
 567 τ and age range α , the distractor pool is defined as

$$568 \quad P_{\tau, \alpha} = \begin{cases} C_{\alpha} \setminus \{v_{\text{corr}}\}, & \tau = \text{Cond}, \\ \mathcal{N}_{\tau, \alpha} \setminus \{v_{\text{corr}}\}, & \tau \in \mathcal{T}, \\ S \setminus \{v_{\text{corr}}\}, & \tau = \text{Sev}. \end{cases} \quad (1)$$

569 where $\mathcal{T} = \{\text{Sym}, \text{Treat}, \text{FollowUp}\}$.

570 The condition set is

$$571 \quad C_{\alpha} = \{c \in V : \text{type}(c) = \text{Condition}, \text{age_range}(c) = \alpha\}, \quad (2)$$

572 and the aggregated neighborhood is

$$573 \quad \mathcal{N}_{\tau, \alpha} = \bigcup_{c \in C_{\alpha}} N_{\tau}(c). \quad (3)$$

574 The neighborhood function is

$$575 \quad N_{\tau}(c) = \begin{cases} \{u : (u, c) \in E, \text{type}(u) = \tau\}, & \tau = \text{Sym}, \\ \{u : (c, u) \in E, \text{type}(u) = \tau\}, & \tau \in \{\text{Treat}, \text{FollowUp}\}. \end{cases} \quad (4)$$

576 For severity classification,

$$577 \quad S = \{u \in V : \text{type}(u) = \text{Severity}\}. \quad (5)$$

578 The final distractor set is

$$579 \quad D = \text{sample}(P_{\tau, \alpha}, k), \quad (6)$$

580 where $k = 3$.