Investigating the Large Language Models' Awareness of Changing Medical Knowledge

Anonymous ACL submission

Abstract

The growing capabilities of Large Language Models (LLMs) can enhance healthcare by assisting medical researchers, physicians, and improving access to health services for patients. LLMs encode extensive knowledge within their parameters, including medical knowledge derived from many sources. However, the knowledge in LLMs can become outdated over time, posing challenges in keeping up with evolving medical recommendations and research. This can lead to LLMs providing outdated health advice or failures in medical reasoning tasks. To address this gap, our study introduces two novel biomedical question-answering (QA) datasets derived from medical systematic literature reviews: MedRevQA, a general dataset of 16,501 biomedical QA pairs, and MedChangeQA, a subset of 512 QA pairs whose verdict changed though time. By evaluating the performance of eight popular LLMs, we find that all models exhibit memorization of outdated knowledge to some extent. We provide deeper insights and analysis, paving the way for future research on this challenging aspect of LLMs.

1 Introduction

003

004

012

014

016

017

018

037

041

The advent of pre-trained Large Language Models (LLMs) has revolutionized the field of Natural Language Processing (NLP) (Minaee et al., 2024). One of the most promising domains for LLM application is healthcare, where they hold the potential to democratize access to health services and improve patient treatment (Thirunavukarasu et al., 2023; Patel and Lam, 2023; Ayers et al., 2023). LLMs are trained to predict the next token on massive amounts of text data, which results with deeply encoding a lot of knowledge in their weights (Dhingra et al., 2022). Recent studies suggest that LLMs also encode clinical knowledge effectively (Singhal et al., 2023), by being trained on diverse medical texts like patient records and clinical trial reports.

World knowledge is dynamic and evolves with time, especially in quickly paced domains like entertainment or politics, where presidents rotate every few years. Still, this is also the case in the world of science and medicine (Solomon, 2015). Medical recommendations can get outdated as new, higherquality evidence from scientific research emerges. Therefore, knowledge encoded within LLMs can get outdated with time and they can struggle to keep up with the ever-changing world knowledge (Zhang et al., 2023). Presence of outdated medical knowledge in LLMs poses a major problem when digital assistants provide consumers outdated health advice (Li et al., 2023; Ong et al., 2024) or fail in clinical settings on complex reasoning tasks (Hager et al., 2024). Even when augmented with retrieved context, LLMs can reject it and resort to internal knowledge in so-called knowledge conflicts (Marjanovic et al., 2024).

042

043

044

047

054

056

060

061

062

063

064

065

067

068

069

070

071

072

073

074

076

077

079

While recent research has started exploring the problem of outdated knowledge in the general and encyclopedic domains (Vu et al., 2024; ChenghaoZhu et al., 2025), less research has been devoted to outdated medical knowledge in LLMs. To bridge this gap, we introduce a novel biomedical questionanswering dataset MedRevQA constructed from medical systematic reviews, with a subset of questions that changed their verdict over the years. We benchmark the performance of eight popular LLMs on the dataset, revealing that the outdated knowledge is present in all of them, especially GPT and Qwen, while Llama and Mistral had the most up-to-date medical knowledge.

2 Related Work

A popular task within NLP for healthcare is biomedical question answering (BQA) (Jin et al., 2022; Nentidis et al., 2024). BQA is seen as a good proxy of evaluating how well LLMs encode and recall medical knowledge (Subramanian et al.,

081

105 106

107

108 109

110

111 112

113

114

115 116

117 118

119

120

122

123 124

125

126

127 128

2024; Singhal et al., 2023, 2025) – therefore, we use it as our main task. Most similar QA dataset in construction is MedREQAL (Vladika et al., 2024), but we majorly expand the scope and the purpose.

Recent work has explored how to measure memorized training examples in LLMs (Jagielski et al., 2023; Maini et al., 2024). Similarly, temporal QA datasets have been constructed to investigate quickly changing knowledge, mostly focusing on the general, encyclopedic domain (Kasai et al., 2023; Vu et al., 2024; Li et al., 2024).

To the best of our knowledge, we introduce the first QA dataset focusing on knowledge change specifically for the medical domain and the first investigation of how much outdated medical knowledge do popular LLMs encode.

3 Dataset

Systematic Literature Reviews (SLRs). SLRs aim to bring evidence together to answer a predefined research question. This involves the identification of primary research relevant to the defined question, the critical appraisal of this research, and the synthesis of the findings (Kolaski et al., 2023). SLRs are considered highest quality evidence in the medical "hierarchy of evidence" (Wallace et al., 2022). We use SLRs to construct a OA dataset because their clear structure and strict criteria used for decisions make them a well suited tool for evaluating the state of encoded medical knowledge within LLMs. We use the SLRs from Cochrane Collaboration (Cumpston et al., 2022), which is the most well-known and respected international organization specializing in construction of SLRs.

Dataset Construction. PubMed, the largest database of medical research publications (White, 2020), contains all the Cochrane systematic review abstracts from 2000 to 2024 (until January 2024 when we scraped). We built a Python scraping project using Scrapy and scraped all the Cochrane SLR abstracts.¹ Every SLR consists of same sections: Background, Objectives, Search methods, Selection criteria, Data collection and analysis, Main results, and Authors' conclusions.

The final QA dataset consists of questions and labels. It was semi-automatically constructed from the abstract text, using gpt-4o-mini-2024-07-18. We utilize the Objectives section as the source of questions, by automatically rewriting them from a

¹https://pubmed.ncbi.nlm.nih.gov/?term= %22Cochrane+Database+syst+rev%22%5BJournal%5D

Question: Does long-term antibiotic use help prevent recurrent urinary tract infections in children?

Conclusion: Long-term antibiotics may reduce the risk of repeat symptomatic UTI in children who have had one or more previous UTIs but the benefit may be small and must be considered together with the increased risk of microbial resistance. (...) [Williams, 2019]

Verdict: Supported

Question: Does long-term antibiotic use help prevent recurrent urinary tract infections in children?

Conclusion: Large, properly randomised, double blinded studies are needed to determine the efficacy of long-term antibiotics for the prevention of UTI in susceptible children. (...) [Williams, 2011]

Verdict: Not Enough Information

Table 1: Example of two instances from our dataset, showing how the verdict changed through time as new, higher quality evidence was discovered.

declarative to an interrogative form using an LLM prompt. Similarly, we leverage an LLM to generate one of the three labels from Authors' conclusion, namely Supported, Refuted, and Not Enough Information as the label. These labels align with common labels in other medical QA and fact-checking datasets (Guo et al., 2022; Glockner et al., 2024b). In total, this dataset has 16,501 QA pairs, spanning virtually all medical disciplines and covering a wide array of important biomedical questions for benchmarking. We call this dataset MedRevQA.

129

130

131

132

133

134

135

136

137

138

145

157

Changed Knowledge. Out of total 16,501 SLR 140 records, 12,122 of them are unique SLRs that have 141 never had an update. The remaining 4379 SLRs 142 constitute 1535 groups (with a minimum 2 SLRs 143 in a group, maximum 9, and the mean 2.85) that 144 researched the same question. This means there are 1535 research objectives that have had multiple 146 SLR iterations written about them. Out of 1535, 147 512 have had a verdict change throughout their 148 time, meaning that the authors changed the conclu-149 sion of the investigated medical research question 150 in a follow-up SLR study, when they acquired updated evidence from research. This follows find-152 ings from medical research studies that have shown 153 how 20 to 30% of Cochrane reviews change their 154 conclusions throughout time (Hughes et al., 2012; 155 Babić et al., 2022). We consolidate these questions 156 with changed verdicts into the MedChangeQA dataset and collect all their verdicts through dif-158 ferent iterations. MedChangeQA has questions, latest label, and (the most recent) outdated label 160

161 for those studies where the label changed.

Annotation Quality. Two annotators, one who 162 is our in-house physician from the university clinic 163 and another an author with background in biomed-164 ical engineering, evaluated a random subset of 100 examples for correctness of generated questions and verdicts. They found 95% of questions and 92% of labels to be correct. We deem 168 is relatively high, since even the human annota-169 tion process is imperfect with errors due to in-170 correct problem understanding or loss of concen-171 tration. A common source of label errors was 172 conflating Refuted and NEI labels. On the other 173 hand, all 512 labels in MedChangeQA were manually checked and corrected by the two annotators. 175 Therefore, MedRevQA has silver labels, while 176 MedChangeQA has gold labels.

Dataset Description. In total, MedRevQA has 178 16,501 questions, of which 6499 are Supported, 3124 are Refuted, and for 6878 there is Not enough information. In MedChangeQA, for the 512 ques-181 tions with changed verdict, the newest labels have a 221/131/160 ratio for S/R/NEI, and the outdated labels are at 152/123/237 for S/R/NEI, showing 184 how the most common change is from not hav-185 ing enough information to becoming supported or refuted by relevant research. Datasets will be released publicly under a CC 4.0 license. 188

4 Experimental Setup

189

190

191

193

194

196

197

198

199

The models were instructed to predict a final label (one of the S/R/NEI) and to explain their output (so generate a long-form answer). For evaluation, we extract the predicted label and compare to groundtruth labels from the dataset, using precision, recall, and F1 score – all macro-averaged because we deem all three classes equally important.

We test multiple LLMs, starting with **GPT-40** (*2024-08-06*), as the most popular proprietary LLM. We also benchmark four open-weights LLMs: **Mistral 24B**, the latest LLM from Mistral AI; **Llama 3.3** (**70B**), the latest LLM from Meta AI; **Qwen 2.5** (**7B**); **DeepSeek-V3** (685B); and finally **OLMo 2** (**13B**) from AllenAI.

We additionally benchmark the performance of two biomedical models: **PMC-LLaMa 13B** (Wu et al., 2023), an extension of Llama 2 and **BioMistral 7B** (Labrak et al., 2024), an extension of Mistral-v0.2. Both were further pre-trained on research papers from PubMed Central. All prompts can be found in Table 6. GPT 40 was prompted through the OpenAI API. The four general-purpose models were prompted via the API of Together AI. Two biomedical LLMs were run locally (in 8-bit quantized version) on one Nvidia V100 GPU with 16 GB VRAM, for two computation hours. The token limit was set to 512 and the temperature parameter to 0 to maximize deterministic outputs. 210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

227

228

229

231

232

233

235

236

237

239

240

241

242

243

244

245

247

248

249

251

252

253

254

255

256

257

259

5 Results and Discussion

Experiment Rounds. We first (1) test the full dataset, MedRevQA. We also did two experiments on MedChangeQA, first with (2) outdated labels as ground truth, then with (3) latest labels as ground truth. We use the difference between the scores of (2) and (3) as a way to show the extent of outdated medical knowledge in LLMs. Final results are systematized in Table 2, measured by macro P, R, and F1. The last column shows the percentage of answers in the 3rd experiment (3) where the outdated label was predicted (and not the correct latest label or an incorrect label altogether).

Performance. On the full dataset, Mistral exhibited the best R and F1, showing it has the best overview of the overall medical knowledge landscape. Precision was the highest in Deepseek-V3. Nevertheless, none of the models has a very high performance, pointing to the challenging nature of MedRevQA as a general biomedical QA testbed.

When it comes to outdated knowledge, Llama 3.3 had the highest degree of latest knowledge as compared to the outdated labels (+7.4), while OLMo also had a positive difference (+2.9). Mistral showed an almost identical performance, while GPT, Qwen, DeepSeek, and PMC-Llama all struggled. Qwen was also the smallest and least capable model, which could explain low scores in general and low awareness of recent knowledge.

An example of outdated and incorrect knowledge is shown in Table 7 in Appendix. Additionally, Figure 1 shows how the average F1 across LLMs on questions from different years on MedRevQA declines in more recent years, as all post-2016 average scores are lower than any before 2016. A similar drop in performance on more recent biomedical questions was found by Park et al. (2025).

Pretraining Data. Most popular LLMs (proprietary and open) do not fully disclose their pretraining data, making it difficult to assert if concrete medical studies were memorized. Still, recent stud-

		Full Dataset (16,501)			Changed Knowledge Dataset (512)							
	Release				(a) O	utdated	Labels	(b) L	atest L	abels	F1	Outdated
	Date	Р	R	F1	P	R	F1	P	R	F1	diff.	Answ. (%)
GPT 40	2024-05-13	52.6	45.1	42.9	45.5	38.9	34.1	35.2	34.5	31.1	-3.0	39.4
Mistral 24B	2025-01-30	50.6	46.3	45.7	38.2	37.6	33.9	36.9	35.5	33.7	-0.2	38.7
Llama 3.3 70B	2024-12-06	52.7	45.9	39.3	38.9	36.6	26.7	42.8	39.3	34.1	+7.4	32.2
Qwen 2.5 7B	2024-09-19	46.4	42.3	38.7	42.6	37.1	30.8	27.1	30.8	26.0	-4.8	35.4
Deepseek V3	2024-12-26	56.2	46.2	43.8	43.2	38.6	33.9	40.2	35.1	32.2	-1.7	40.6
OLMo 2 13B	2024-11-24	43.5	42.5	37.9	36.2	35.3	29.3	35.5	35.7	33.2	+2.9	32.0
PMC-Lm 13B	2023-08-28	39.5	37.6	36.5	41.9	39.8	35.9	34.5	34.3	33.1	-2.8	37.3
BioMistral 7B	2024-02-19	41.2	41.5	40.9	36.8	37.2	36.3	35.4	35.5	35.3	-1.5	37.1

Table 2: Final results of eight LLMs, measured by Precision (P), Recall (R), and F1 (macro-averaged). Experiments include the full dataset, and the changed knowledge dataset, using the (a) outdated labels and (b) latest labels as ground truth, respectively. Final column is the percentage of answers in (a) where an outdated label was predicted.

ies demonstrated empirically the presence of memorized medical datasets (Gallifant et al., 2024; Yang et al., 2024). We also saw a tendency of models to explicitly mention specific studies in their answers, 263 including Cochrane reviews, many of which were decade-old (see Table 8), thus displaying outdated memorized knowledge (see Table 4). We outline pre-training corpora of used LLMs in Appendix B, and for the fully open OLMo, we show the presence of all used SLRs in its pre-training corpus, with earlier ones being more prevalent (Figure 2). Potential Explanations. We hypothesize some reasons for the presence of strongly encoded outdated knowledge. Firstly, older scientific findings have been around for a longer time and have already permeated the Internet, discussions, follow-up studies, and other resources present in pre-training corpora. Additionally, scientific findings are often misrepresented online (Glockner et al., 2024a; Wührl et al., 2024), so faulty medical knowledge could get encoded. Secondly, LLM memorization rate has been correlated in past work with various training parameters, such as learning rate (Tirumala et al., 2022), model size (Biderman et al., 2023), or frequency of appearance in training data (Carlini et al., 2023). Therefore, it is possible Llama had highest data quality and more weight during training put on 286 more recent text, leading to the best performance. Finally, the cutoff of all models is 2023 and the vast majority of "latest labels" originated from before 2023 (see Figure 3). Cutoff could explain the drop in 2023/2024 (Fig. 1) but not earlier years.

260 261

262

270

271

272

273

276

277

278

282

287

290

Future Directions. One potential improvement is using retrieval-augmented generation (RAG). We 294 demonstrate how a simple retrieval strategy (with the first PubMed results added to prompt) can al-295 ready bring decent improvement in Appendix A. Still, LLMs can hallucinate extra information even within RAG settings (Adlakha et al., 2024) or not 298



Figure 1: Average F1-Macro performance for questions originating from each year in the dataset across five LLMs, showing decline in more recent years.

follow the provided references (Liu et al., 2023). Therefore, advanced RAG techniques are needed (Yu et al., 2024). Other promising future directions for outdated medical knowledge, where our datasets could be used, include: resolving knowledge conflicts (Wang et al., 2024b), machine unlearning (Yao et al., 2024; Gao et al., 2025), knowledge editing (Wang et al., 2024a; Jiang et al., 2024). 299

300

301

302

303

304

305

307

308

309

310

311

312

313

314

315

316

317

318

319

320

6 Conclusion

We introduce two new QA datasets constructed from rigorous biomedical SLRs for benchmarking general biomedical knowledge of LLMs. With a subset of 512 questions where the answer (verdict) changed over the years, we showcase how eight popular LLMs fare better on older medical knowledge and encode a considerable amount of outdated knowledge labels, which can hinder their usability in healthcare settings, including helping physicians, researchers, and patients. We outline future work directions and hope our datasets will serve as a challenging testbed for addressing the problem of outdated LLM knowledge.

Limitations

321

323

324

325

326 327

331

332

336

338

341

342

344

345

352

370

Most of the MedRevQA dataset introduced in this study was constructed semi-automatically, by scraping the content and using an LLM to generate the question and label. It is possible that some of the generated questions and labels are imperfect. Our manual analysis of 100 randomly selected instances showed that the performance is 92–95% correct, for labels and questions respectively. We considered this to be a good enough performance, considering that even human annotation is not always perfect. MedRevQA should be interpreted as having silver labels and used as such.

We use the difference between the predicted labels when using "outdated labels" and "latest labels" as ground truth, as a proxy for evaluating the degree of encoded outdated medical knowledge. This is not a perfect measure since it is possible that the LLM predicted an incorrect label due to some logical error or misinterpreting the question, not necessarily because of outdated knowledge. Still, our manual inspection of a big number of generated labels and explanations showed that outdated references were indeed the most common explanation for the label misprediction and models often referred to old SLRs and meta-analysis, ranging many years in the past.

We do not benchmark all relevant biomedical LLMs; some, like Med-Gemini are either not available or computationally too expensive for us to run. Due to lack of financial resources, our study also lacks deeper human evaluation of generated model labels with medical experts, which could have given a more rigorous evaluation.

Ethics Statement

The work presented in this study focuses on the delicate field of healthcare and medical natural language processing (NLP). We predict answers to question in a zero-shot setting to uncover their internal encoded medical knowledge for research purposes, but this approach is not suitable for end users or patients. Some responses can include inaccuracies and misleading medical advice, which should be critically evaluated and verified with reliable sources or medical professionals.

The original text of Cochrane's systematic literature review abstracts belongs to the Cochrane Collaboration. We will release only the generated questions and label pairs, and respect the copyright of the original text of Cochrane authors.

References

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. 2024. Evaluating correctness and faithfulness of instructionfollowing models for question answering. Transactions of the Association for Computational Linguistics, 12:681–699.
- John W. Ayers, Adam Poliak, Mark Dredze, Eric C. Leas, Zechariah Zhu, Jessica B. Kelley, Dennis J. Faix, Aaron M. Goodman, Christopher A. Longhurst, Michael Hogarth, and Davey M. Smith. 2023. Comparing physician and artificial intelligence chatbot responses to patient questions posted to a public social media forum. JAMA Internal Medicine, 183(6):589-596.
- Andrija Babić, Tina Poklepović Peričić, Dawid Pieper, and Livia Puljak. 2022. When is the evidence conclusive? analysis of systematic reviews for which cochrane declared that conclusions will not change with further studies. Research Synthesis Methods, 13(4):478-488.
- Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Gregory Anthony, Shivanshu Purohit, and Edward Raff. 2023. Emergent and predictable memorization in large language models. In Thirty-seventh Conference on Neural Information Processing Systems.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2023. Quantifying memorization across neural language models. In The Eleventh International Conference on Learning Representations.
- ChenghaoZhu ChenghaoZhu, Nuo Chen, Yufei Gao, Yunyi Zhang, Prayag Tiwari, and Benyou Wang. 2025. Is your LLM outdated? a deep look at temporal generalization. In Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 7433-7457, Albuquerque, New Mexico. Association for Computational Linguistics.
- Miranda S Cumpston, Joanne E McKenzie, Vivian A Welch, and Sue E Brennan. 2022. Strengthening systematic reviews in public health: guidance in the cochrane handbook for systematic reviews of interventions. Journal of Public Health, 44(4):e588e592.
- Bhuwan Dhingra, Jeremy R Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W Cohen. 2022. Time-aware language models as temporal knowledge bases. Transactions of the Association for Computational Linguistics, 10:257-273.
- Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. 2021. Documenting

413

414

415

416

417

418

419

420

421

422

423

424

425

426

371 372

373

535

536

537

538

539

540

1: 1: 1: 1: 2. g-22. pys ni, 483

484

large webtext corpora: A case study on the colossal clean crawled corpus. In *Proceedings of the* 2021 Conference on Empirical Methods in Natural Language Processing, pages 1286–1305, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

- Jack Gallifant, Shan Chen, Pedro José Ferreira Moreira, Nikolaj Munch, Mingye Gao, Jackson Pond, Leo Anthony Celi, Hugo Aerts, Thomas Hartvigsen, and Danielle Bitterman. 2024. Language models are surprisingly fragile to drug names in biomedical benchmarks. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 12448– 12465, Miami, Florida, USA. Association for Computational Linguistics.
- Chongyang Gao, Lixu Wang, Kaize Ding, Chenkai Weng, Xiao Wang, and Qi Zhu. 2025. On large language model continual unlearning. In *The Thirteenth International Conference on Learning Representations*.
- Max Glockner, Yufang Hou, Preslav Nakov, and Iryna Gurevych. 2024a. Missci: Reconstructing fallacies in misrepresented science. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4372– 4405.
- Max Glockner, Ieva Staliūnaitė, James Thorne, Gisela Vallejo, Andreas Vlachos, and Iryna Gurevych. 2024b. AmbiFC: Fact-checking ambiguous claims with evidence. *Transactions of the Association for Computational Linguistics*, 12:1–18.
- Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. *Transactions of the Association for Computational Linguistics*, 10:178–206.
- Paul Hager, Friederike Jungmann, Robbie Holland, Kunal Bhagat, Inga Hubrecht, Manuel Knauer, Jakob Vielhauer, Marcus Makowski, Rickmer Braren, Georgios Kaissis, et al. 2024. Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nature medicine*, 30(9):2613– 2622.
- EG Hughes, M van Wely, and CM Farquhar. 2012. Cochrane reviews in perspective: the importance of appropriate conclusions and timing of publication.
- Matthew Jagielski, Om Thakkar, Florian Tramer, Daphne Ippolito, Katherine Lee, Nicholas Carlini, Eric Wallace, Shuang Song, Abhradeep Guha Thakurta, Nicolas Papernot, and Chiyuan Zhang. 2023. Measuring forgetting of memorized training examples. In *The Eleventh International Conference* on Learning Representations.
- Yuxin Jiang, Yufei Wang, Chuhan Wu, Wanjun Zhong, Xingshan Zeng, Jiahui Gao, Liangyou Li, Xin Jiang, Lifeng Shang, Ruiming Tang, et al. 2024. Learning to edit: Aligning Ilms with knowledge editing. In

Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4689–4705.

- Qiao Jin, Zheng Yuan, Guangzhi Xiong, Qianlan Yu, Huaiyuan Ying, Chuanqi Tan, Mosha Chen, Songfang Huang, Xiaozhong Liu, and Sheng Yu. 2022. Biomedical question answering: a survey of approaches and challenges. *ACM Computing Surveys* (*CSUR*), 55(2):1–36.
- Jungo Kasai, Keisuke Sakaguchi, yoichi takahashi, Ronan Le Bras, Akari Asai, Xinyan Velocity Yu, Dragomir Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui. 2023. Realtime QA: What's the answer right now? In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track.*
- Kat Kolaski, Lynne Romeiser Logan, and John PA Ioannidis. 2023. Guidance to best tools and practices for systematic reviews. *Journal of Pediatric Rehabilitation Medicine*, 16(2):241–273.
- Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. BioMistral: A collection of opensource pretrained large language models for medical domains. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 5848–5864, Bangkok, Thailand. Association for Computational Linguistics.
- Yucheng Li, Frank Guerin, and Chenghua Lin. 2024. Latesteval: addressing data contamination in language model evaluation through dynamic and timesensitive test construction. In Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'24/IAAI'24/EAAI'24. AAAI Press.
- Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge. *Cureus*, 15(6).
- Jiacheng Liu, Sewon Min, Luke Zettlemoyer, Yejin Choi, and Hannaneh Hajishirzi. 2024. Infini-gram: Scaling unbounded n-gram language models to a trillion tokens. In *First Conference on Language Modeling*.
- Nelson F Liu, Tianyi Zhang, and Percy Liang. 2023. Evaluating verifiability in generative search engines. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7001–7025.
- Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. S2ORC: The semantic scholar open research corpus. In *Proceedings of the* 58th Annual Meeting of the Association for Computational Linguistics, pages 4969–4983, Online. Association for Computational Linguistics.

651

652

653

654

655

- 541 542
- 545
- 546

- 554 555

564

- 565 566
- 567
- 569
- 570 571
- 572
- 573 574
- 578 579
- 580

582 583

585

- 586
- 587

590 591

593

594

597

- Pratyush Maini, Hengrui Jia, Nicolas Papernot, and Adam Dziedzic. 2024. LLM dataset inference: Did you train on my dataset? In The Thirty-eighth Annual Conference on Neural Information Processing Systems.
- Sara Vera Marjanovic, Haeun Yu, Pepa Atanasova, Maria Maistro, Christina Lioma, and Isabelle Augenstein. 2024. DYNAMICQA: Tracing internal knowledge conflicts in language models. In Findings of the Association for Computational Linguistics: EMNLP 2024, pages 14346–14360, Miami, Florida, USA. Association for Computational Linguistics.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. arXiv preprint arXiv:2402.06196.
- Anastasios Nentidis, Georgios Katsimpras, Anastasia Krithara, Salvador Lima-López, Eulàlia Farré-Maduell, Martin Krallinger, Natalia Loukachevitch, Vera Davydova, Elena Tutubalina, and Georgios Paliouras. 2024. Overview of bioasq 2024: the twelfth bioasq challenge on large-scale biomedical semantic indexing and question answering. In International Conference of the Cross-Language Evaluation Forum for European Languages, pages 3–27. Springer.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Michal Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2025. 2 olmo 2 furious. Preprint, arXiv:2501.00656.
 - Jasmine Chiat Ling Ong, Shelley Yin-Hsi Chang, Wasswa William, Atul J. Butte, Nigam H. Shah, Lita Sui Tjien Chew, Nan Liu, Finale Doshi-Velez, Wei Lu, Julian Savulescu, and Daniel Shu Wei Ting. 2024. Medical ethics of large language models in medicine. NEJM AI, 1(7):AIra2400038.
- Yein Park, Chanwoong Yoon, Jungwoo Park, Donghyeon Lee, Minbyul Jeong, and Jaewoo Kang. 2025. Chroknowledge: Unveiling chronological knowledge of language models in multiple domains. In The Thirteenth International Conference on Learning Representations.
- Sajan B Patel and Kyle Lam. 2023. Chatgpt: the future of discharge summaries? The Lancet Digital Health, 5(3):e107-e108.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales,

Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023. Large language models encode clinical knowledge. Nature, 620(7972):172-180.

- Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin Clark, Stephen R Pfohl, Heather Cole-Lewis, et al. 2025. Toward expert-level medical question answering with large language models. Nature Medicine, pages 1-8.
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Jha, Sachin Kumar, Li Lucy, Xinxi Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Evan Walsh, Luke Zettlemoyer, Noah Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. 2024. Dolma: an open corpus of three trillion tokens for language model pretraining research. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15725–15788, Bangkok, Thailand. Association for Computational Linguistics.
- Luca Soldaini and Kyle Lo. 2023. peS2o (Pretraining Efficiently on S2ORC) Dataset. Technical report, Allen Institute for AI. ODC-By, https://github. com/allenai/pes2o.
- Miriam Solomon. 2015. Making medical knowledge. Oxford University Press.
- Anand Subramanian, Viktor Schlegel, Abhinav Ramesh Kashyap, Thanh-Tung Nguyen, Vijay Prakash Dwivedi, and Stefan Winkler. 2024. M-QALM: A benchmark to assess clinical reading comprehension and knowledge recall in large language models via question answering. In Findings of the Association for Computational Linguistics: ACL 2024, pages 4002-4042, Bangkok, Thailand. Association for Computational Linguistics.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. Nature medicine, 29(8):1930-1940.
- Kushal Tirumala, Aram H. Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. Memorization without overfitting: Analyzing the training dynamics of large language models. In Advances in Neural Information Processing Systems.
- Juraj Vladika, Phillip Schneider, and Florian Matthes. 2024. MedREQAL: Examining medical knowledge recall of large language models via question answering. In Findings of the Association for Computational Linguistics: ACL 2024, pages 14459-14469, Bangkok, Thailand. Association for Computational Linguistics.

Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. 2024. Fresh-LLMs: Refreshing large language models with search engine augmentation. In Findings of the Association for Computational Linguistics: ACL 2024, pages 13697–13720, Bangkok, Thailand. Association for Computational Linguistics.

656

657

672

674

675

682

683

688

697

701

703

704

706

707

710

- Sowdhamini S Wallace, Gal Barak, Grace Truong, and Michelle W Parker. 2022. Hierarchy of evidence within the medical literature. Hospital Pediatrics, 12(8):745-750.
- Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, and Jundong Li. 2024a. Knowledge editing for large language models: A survey. ACM Computing Surveys, 57(3):1–37.
- Yike Wang, Shangbin Feng, Heng Wang, Weijia Shi, Vidhisha Balachandran, Tianxing He, and Yulia Tsvetkov. 2024b. Resolving knowledge conflicts in large language models. In First Conference on Language Modeling.
- Jacob White. 2020. Pubmed 2.0. Medical reference *services quarterly*, 39(4):382–387.
- Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Pmc-llama: Towards building open-source language models for medicine. Preprint, arXiv:2304.14454.
- Amelie Wührl, Dustin Wright, Roman Klinger, and Isabelle Augenstein. 2024. Understanding fine-grained distortions in reports of scientific findings. In Findings of the Association for Computational Linguistics ACL 2024, pages 6175-6191.
- Xinyu Yang, Zichen Wen, Wenjie Qu, Zhaorun Chen, Zhiying Xiang, Beidi Chen, and Huaxiu Yao. 2024. Memorization and privacy risks in domain-specific large language models. In ICLR 2024 Workshop on Reliable and Responsible Foundation Models.
- Yuanshun Yao, Xiaojun Xu, and YangLiu. 2024. Large language model unlearning. In Advances in Neural Information Processing Systems, volume 37, pages 105425-105475. Curran Associates, Inc.
- Yue Yu, Wei Ping, Zihan Liu, Boxin Wang, Jiaxuan You, Chao Zhang, Mohammad Shoeybi, and Bryan Catanzaro. 2024. RankRAG: Unifying context ranking with retrieval-augmented generation in LLMs. In The Thirty-eighth Annual Conference on Neural Information Processing Systems.
- Zihan Zhang, Meng Fang, Ling Chen, Mohammad-Reza Namazi-Rad, and Jun Wang. 2023. How do large language models capture the ever-changing world knowledge? a review of recent advances. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 8289-8311, Singapore. Association for Computational Linguistics.

Simple RAG Improvement Α

Common ways to address the outdated knowledge 712 with knowledge editing include continual learning 713 methods and external search augmentation. We do 714 a simple experiment using a retrieval-augmented 715 method. For each of the 512 questions in Med-716 ChangeQA, we query the PubMed API and take 717 the abstract of the Top 1 result, and append it to the 718 main prompt as an additional context. Results are 719 shown in Table 3. This improves the F1 scores by 720 margin of 3–16 and partially closes the gap, but still 721 leaves a lot of room for improvement. This serves 722 as a simple demonstration of one way to address 723 the outdated knowledge – future work could focus 724 on retrieving more documents, using structured and 725 focused search queries (like searching for SLRs), semantic search, graph RAG, learning to re-rank and avoid conflicts, etc. Additionally, methods of continual learning and fine-tuning could be used, with MedChangeQA serving as a testbed to measure the rate of success of the proposed techniques. 731

	P	R	F1	Improv. F1
GPT 40	43.4	40.2	39.8	+8.7
Mistral	47.5	41.5	39.6	+5.9
Llama 3.3	44.1	42.3	38.8	+4.7
Qwen	43.3	43.7	42.2	+16.2
Deepseek	40.7	39.3	35.4	+3.2

Table 3: Performance improvements with the abstract of the top PubMed paper included in prompt.

Pre-Training Data of LLMs B

B.1 General LLMs

This section provides an overview of what is publicly known about pre-training data for each of the used LLMs, as reported in their respective technical reports or official documentation:

- Llama 3.3: Pretrained on approximately 15 trillion tokens of data sourced from publicly available online sources. The exact composition and breakdown of the dataset are not detailed, but Meta emphasizes that the data is "a new mix" of public internet data. The data cutoff for pretraining is December 2023.
- Mistral 24B The official technical report and available documentation do not provide explicit details about the pre-training corpus for 747 Mistral 24B. However, Mistral's models are generally known to be trained on large-scale, diverse datasets, often including filtered web 750

711

data, code, and other standard sources, but no 751 specifics are publicly disclosed for the 24B 752 version in the sources provided.

754

757

761

767

772

773

775

776

781

782

790

792

796

799

- GPT 40 It was trained on data up to October 2023, sourced from a "wide variety of mate-755 rials," including: (a) publicly available data 756 (web pages, ML datasets, common crawls), (b) proprietary data (obtained via data partner-758 ships, e.g., paywalled content, archives, metadata), (c) key dataset components (web data, code and math data, multimodal data). The dataset underwent safety filtering to remove harmful content, personal information, and explicit material. OpenAI does not provide a detailed breakdown of dataset proportions or specific sources. 766
 - Qwen 2.5 It was was trained on up to 18 trillion tokens of data. The dataset is described as "large-scale" and "high-quality," but the technical report does not specify the exact sources. The data is designed to provide a strong foundation for common sense, expert knowledge, and reasoning. Qwen 2.5 also supports multilingual capabilities across more than 29 languages.
 - DeepSeek V3 It was trained on 14.8 trillion tokens of "diverse and high-quality" data. The dataset construction focused on: increased ratio of mathematical and programming samples, multilingual coverage, data processing pipeline optimized for diversity and minimal redundancy. Used document packing and fillin-middle (FIM) strategies for code and text infilling tasks. The technical report does not provide a granular breakdown of data sources but highlights the focus on math, code, and multilingual content.
 - PMC-Llama and BioMistral use the base models of Llama and Mistral as described before, but were then further pre-trained on abstracts of biomedical publications from PubMed and full publications from PubMed Central. As described in our paper, PubMed contains all the abstracts of Cochrane systematic reviews, which means this data was surely seen.

B.2 Inspection of OLMo

OLMo 2 (OLMo et al., 2025) is trained on the Dolma corpus (Soldaini et al., 2024), a fully open dataset containing around 3 trillion tokens. This is a high-level breakdown of the composition of the pre-training corpus:

- Common Crawl: 2,479 billion tokens 803 • GitHub: 411 billion tokens 804 • Reddit: 89 billion tokens 805 • Semantic Scholar: 70 billion tokens
- Project Gutenberg: 6.0 billion tokens 807
- Wikipedia, Wikibooks: 4.3 billion tokens

In particular, the Semantic Scholar part consists of peS2o (Soldaini and Lo, 2023) and S2ORC (Lo et al., 2020) corpora that constitute the academic knowledge base Semantic Scholar. This database also indexes all of Cochrane's systematic literature reviews.² Therefore, we can with high certainty say that the OLMo models have seen Cochrane's SLRs during its pre-training process. Other than in the two academic corpora, there is a wide presence of these SLRs in other parts of the dataset, especially various websites of the Internet found in Common Crawl (Dodge et al., 2021).



Figure 2: N-gram counts per study year in Dolma.

We used Infini-gram (Liu et al., 2024), an ngram language model that can be used to query Dolma and other pre-training corpora,³ to inspect the presence of Cochrane's SLRs. Searching for "Cochrane Database of Systematic Reviews" (exact journal name, case-sensitive search) returns 144,493 hits for Dolma v1.7 used for OLMo 2 pre-training. Additionally, we queried the title of

828

800

801

802

808

809

810

811

812

813

814

815

816

817

818

819

820

²https://www.semanticscholar.org/venue?name= Cochrane%20Database%20of%20Systematic%20Reviews ³https://huggingface.co/spaces/liujch1998/

infini-gram

each of the SLR studies found in our MedRevQA (n=16,501) and report on the mean and median 830 amount of n-gram counts per year in Figure 2. The figure shows how the mean and median almost steadily decrease throughout the years, meaning that the most mentioned and discussed studies are the earliest ones since they have had more time to spread throughout the web. The higher frequency of mentions can lead to to stronger encoding of this outdated knowledge in LLM weights.

829

833

835

839

842

844

847

848

852

854

861

863

С **Mentions of Common Terms**

In Table 4, we show the count of the number of mentions of some common terms referring to specific medical studies across all LLM answers on MedRevQA questions. It shows how most models tend to cite specific studies when providing some of their answers and this can become problematic when the referred studies are outdated and deprecated. Other than specific studies, the term "studies" in the last row is the only generic term included here. It is notable how GPT overwhelmingly resorted to using general phrasing such as "studies have shown a positive effect..." without specifying what studies exactly it is referring to. This likely comes from its final alignment and preference-tuning phase where a particular style of answers is learned.

D **Prompts and Examples**

This Appendix section provides additional material for the study, including the model prompts in full length (Tables 5 and 6) and example questions and model answers (Tables 7 and 8). Figure 3 shows distribution of "latest label" through years. Additionally, Figure 4 shows a larger version of the plot of average F1 score for tested LLMs on questions through years.



Figure 3: Distribution of the year of "Latest Label", either most recent label from a particular SLR study group or label of a standalone SLR (n=12,972)

	Llama 3.3	Mistral	GPT 40	Qwen 2.5	Deepseek	OLMo 2	BioMistral	PMC-L
"Cochrane"	51	783	2	629	901	283	2067	2344
"systematic	221	1664	623	3194	3046	531	3990	4956
review"								
"meta-	844	3511	714	4180	2776	981	4253	4618
analysis"								
"journal"	53	689	7	4620	448	196	574	624
"studies"	7024	12419	13493	12516	13421	6615	7598	9720

Table 4: Number of answers (out of 16,501) by each tested LLM where the respective terms were mentioned, showing the tendency to refer to and cite relevant medical studies that were memorized during pre-training. Two biomedical models, that were further pre-trained on biomedical publications, also refer to specific studies the most.

Use Case	Prompt Content
Question & Label	SYSTEM: You're a helpful assistant. Your task is to help with generating questions
generation	and labels in the medical and clinical domain.
	USER You will be given an excerpt of an abstract of a clinical systematic review.
	Based on the given background, objectives, and author's conclusions, generate
	only ONE SINGLE question, answerable with yes/no/uncertain, that sums up the
	main medical objective that was investigated. Please keep the question short and
	general and use the "Objectives" section to construct the question. The question
	should be about a general medical hypothesis, not about this specific review.
	Afterwards, please also give a label for the author's conclusions. Label tries
	to answer the objective by looking at the conclusion. The label may be ONLY
	from one of the following three: (1) SUPPORTED: (2) REFUTED: (3) NOT
	ENOLIGH INFORMATION Do not try to make up a new label. Please only
	select the third label if not enough evidence was found to reach a verdict not
	if certainty of the conclusion is low! Place aim to predict "SUPPORTED" or
	"DEFLITED" over if containing of these conclusions by outputs is low or moderate
	REFUTED even in certainty of these conclusions by authors is low of moderate.
	Please structure the output in two lines, as:
	QUESTION: (question)
	LABEL: (label)
	The documents begins now.

Table 5: Overview of applied prompts for data generation and annotation.

Model	Prompt
PMC-LLaMa	Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request. ### Instruction: Based on your best current knowledge, please answer the following medical question. If you think there is not enough evidence to answer, then say so. Please answer the question with "SUPPORTED" or "REFUTED" or "NOT ENOUGH INFORMATION". Briefly explain your answer. ### Input: {question} ### Response:
BioMistral	<s>[INST] Based on your knowledge, please answer this clinical question only with SUPPORTED (if the question is supported by the clinical research) or REFUTED (if the hypothesis is refuted by the current clinical research) or NOT ENOUGH INFORMATION (if there is insufficient evidence for the question in current research). Please give your output in form of LABEL: (label) . Briefly explain your answer. QUESTION: {question} [/INST]</s>
Mistral 24B, Llama 3.3, GPT- 40, Qwen 2.5, DeepSeek-V3	 SYSTEM You are an AI assistant helping answer clinical and medical questions based on your best knowledge. AGENT Please answer this clinical question only with SUPPORTED (if the question is supported by the clinical research) or REFUTED (if the hypothesis is refuted by the current clinical research) or NOT ENOUGH INFORMATION (if there is insufficient evidence for the question in current research). Please give your output in form of LABEL: (label) . Briefly explain your answer. QUESTION: {question} Table 6: Prompts used for each of the LLMs in experiments.

Question	Is routine malaria chemoprophylaxis beneficial for people with sickle cell disease in malaria-endemic areas?				
Conclusion	() It is beneficial to give routine malaria chemoprophylaxis in sickle cell disease in areas where malaria				
(Latest)	is endemic. (Oniyangi, 2019) [label: Supported]				
Conclusion	() There is very little direct evidence to support or refute giving routine chemoprophylaxis in sickle ce				
(Outdated)	disease in areas where malaria is endemic. (Oniyangi, 2006) [label: Not Enough Information]				
GPT-40	Routine malaria chemoprophylaxis is beneficial for people with sickle cell disease in malaria-endemic areas. Clinical research indicates that individuals with sickle cell disease are at increased risk for severe malaria and its complications. Chemoprophylaxis can reduce the incidence of malaria and its associated morbidity in this vulnerable population, thereby improving health outcomes. [predicted: Supported]				
Llama 3.3	Current clinical research suggests that routine malaria chemoprophylaxis is not recommended for people with sickle cell disease in malaria-endemic areas, as the risks associated with long-term prophylaxis may outweigh the benefits, and other preventive measures such as insecticide-treated bed nets and prompt treatment of malaria episodes are preferred. [predicted: Refuted]				
Mistral 24B	Routine malaria chemoprophylaxis is generally recommended for people with sickle cell disease living in malaria-endemic areas. This is supported by clinical research indicating that malaria infections can trigger sickle cell crises and other complications. Chemoprophylaxis helps reduce the risk of malaria, thereby decreasing the frequency of these crises and improving overall health outcomes. [predicted: Supported]				

Table 7: Example of a question and two labels, one outdated and one recent from our dataset. Answers show how GPT and Mistral capture the latest and up-to-date knowledge for this question, while Llama predicted incorrect and outdated recommendations.

Question	Do probiotics reduce the risk of necrotising enterocolitis in very preterm or low birth weight infants?
Conclusion (Outdated)	Enteral supplementation of probiotics prevents severe NEC and all cause mortality in preterm infants. Our updated review of available evidence strongly supports a change in practice. () (AlFaleh et al., 2014) [label: Supported]
Conclusion (Latest)	Given the low to moderate certainty of evidence for the effects of probiotic supplements on the risk of NEC and associated morbidity and mortality for very preterm or VLBW infants, and particularly for extremely preterm or ELBW infants, there is a need for further large, high-quality trials to provide evidence of sufficient validity and applicability to inform policy and practice. (Sharif et al., 2023) [label: Not Enough Information]
Mistral 24B	Several clinical studies and meta-analyses have shown that probiotic supplementation in very preterm or low birth weight infants can reduce the risk of necrotizing enterocolitis (NEC). A Cochrane review from 2014 has supported this finding, indicating that probiotics can be beneficial in preventing NEC in this vulnerable population. [predicted : Supported]

Table 8: Example of a study objective that had the *Supported* label in its 2014 version and an updated *Not enough information* label in its 2023 version. Mistral cited the more than decade-old study to supports its answer.



Figure 4: Larger version of the previous figure: Average F1-Macro performance for questions originating from each year in the dataset across five LLMs, showing decline in more recent years.