42

44

46

48

52

59

61

62

63

64

65

66

67

68

69

70

71

72

73

MedG–KRP: Medical Graph Knowledge Representation Probing

Anonymous Author

1

Abstract

Recently, large language models (LLMs) have 2 emerged as powerful tools, finding many ap-3 plications in medicine. LLMs' ability to coa-4 lesce vast amounts of information from many 5 sources in order to come to a response-a pro-6 cess similar to that of a human expert—has 7 led many to see potential in deploying LLMs 8 for clinical use. However, medicine is a setting 9 where accurate reasoning is paramount. Many 10 researchers are questioning the effectiveness of 11 multiple choice question answering (MCQA) 12 benchmarks, frequently used to test LLMs. Re-13 searchers and clinicians alike must have com-14 plete confidence in LLMs' abilities for them to 15 be deployed in a medical setting. In order to 16 address this need for understanding, we intro-17 duce a knowledge graph (KG)-based method 18 to evaluate the biomedical reasoning abilities 19 of LLMs. Essentially, we map how LLMs link 20 medical concepts in order to better understand 21 how they reason. We test GPT-4, Llama3-22 70b, and PalmyraMed-70b, a medical model. 23 We enlist a panel of medical students to review 24 a total of 60 LLM-generated graphs and com-25 pare these graphs to BIOS, a large biomedical 26 KG. We observe GPT-4 to perform best in our 27 human review but worst in our ground truth 28 comparison; vice-versa with PalmyraMed, the 29 medical model. Our work provides a means of 30 visualizing the medical reasoning pathways of 31 LLMs, so they can be implemented in clinical 32 settings safely and effectively. 33

Keywords: Knowledge Graph, Large 34 Language Models, Healthcare, Biomedical 35 Database, Causal Graph 36

Data and Code Availability Prompts, generated 37 graphs, code and human evaluations are available at 38 https://tinyurl.com/35c9aeap. 39

Institutional Review Board (IRB) Our re-40 search does not require IRB approval. 41

1. Introduction

The increasing use of large language models 43 (LLMs) has diversified their applications beyond standard natural language processing (NLP) tasks 45 such as text generation, translation, and summarization (Wu et al., 2021; OpenAI, 2023; Dubey and 47 et.al., 2024; Xu et al., 2023). This advancement has led to a growing interest among researchers and 49 healthcare professionals in leveraging LLMs for med-50 ical applications. The capacity of LLMs to handle 51 extensive volumes of clinical data, medical records, and scientific literature (Huang et al., 2019; Alsentzer 53 et al., 2019; Bolton et al., 2022) introduces the po-54 tential for advancements in clinical decision support, 55 diagnostics, and patient management (Yang et al., 56 2022; Jiang et al., 2023a; Singhal et al., 2023b; Mc-57 Duff et al., 2023; Tu et al., 2024). For safety-critical 58 applications such as healthcare, the performance of LLM must be vigorously validated. 60

Benchmarking LLMs' medical abilities is a challenging task, however. Medical knowledge, even when limited to common diseases, is vast, making it difficult to design benchmarks that capture the breadth of information clinicians rely on daily (Jain, 2024). Additionally, due to the large volume of training data that LLMs memorize, there is concern that their performance on traditional benchmarks may be artificially inflated by memorization (Carlini et al., 2023). As a result, developing more rigorous and comprehensive benchmarks is essential to accurately evaluate LLMs' true medical understanding and ensure their safe and effective deployment in clinical settings.

The medical capabilities of LLMs are often eval-74 uated through multiple-choice question answering 75 (MCQA) benchmarks (Pal et al., 2022). Datasets 76 such as, MedQA, which is based on questions from 77 the USMLE, draw directly from standardized med-78 ical examinations, while other benchmarks, such as 79 MultiMedQA, aggregate data from a variety of medi-80 cal knowledge sources (Jin et al., 2020; Singhal et al., 81 2023a). However, recent findings by Griot et al. 82 (2024) raise concerns that these MCQA benchmarks 83

may not adequately evaluate the depth of LLMs' 84 medical understanding or reasoning ability, suggest-85 ing that performance may be influenced by surface-86 level pattern recognition rather than genuine clini-87 cal reasoning. Moreover, prior studies have demon-88 strated that some state-of-the-art LLMs exhibit bi-89 ases in medical reasoning and perform poorly in es-٩n sential tasks such as medical coding, highlighting fur-91

- $_{92}$ ther limitations in their practical utility and accuracy
- 93 (Omiye et al., 2023; Soroush et al., 2024).

Confidence in LLMs' medical capabilities and the
methods used for their evaluation must be ensured
before their deployment in clinical settings. It is
therefore essential to develop alternative methods for
a comprehensive assessment of LLMs' performance.

Our research is guided by the objective of increas-99 ing the transparency of LLMs by structuring their 100 medical reasoning processes. This approach aims to 101 offer a deeper understanding of LLMs' medical per-102 formance that extends beyond the capabilities of tra-103 ditional MCQA benchmarks. To address the limi-104 tations inherent in existing evaluation methods, we 105 propose a novel technique for visualizing the con-106 nections between medical concepts and understand-107 ing pathways in medicine by generating knowledge 108 graphs (KGs). This method reduces the risk of LLMs 109 relying on verbatim memorization from pretraining 110 data and circumvents issues related to the overlap of 111 benchmark data with the training corpus. 112

Our work is motivated by the dual nature of LLMs: 113 their potential to automate complex medical tasks 114 while also presenting challenges due to their black-115 box nature and susceptibility to errors. Our ap-116 proach *MedG-KRP* leverages LLMs to systematically 117 structure and visualize their parametric knowledge. 118 We begin with a single medical concept and use the 119 LLM to identify and generate a knowledge graph of 120 "causes" and "effects" associated with this concept. 121 To the best of our knowledge, this is the first work to 122 leverage an LLM to systematically generate a knowl-123 edge graph from a single, specified medical concept. 124 The knowledge graphs generated by the MedG-125 *KRP* process offer several potential applications. By 126 interpreting these graphs as proxies for LLMs' inter-127 nal knowledge structures, we can enhance the inter-128 pretability of the models by examining their grasp 129 of medical pathways. Additionally, LLM-generated 130 graphs could be employed to augment or correct ex-131 isting biomedical knowledge graphs. Future research 132 could also explore using MedG-KRP for chain-of-133 thought (COT) prompting, as proposed by Wei et al. 134

(2023). In this approach, models could first generate ¹³⁵ knowledge graphs to inform their reasoning process ¹³⁶ when addressing medical questions. ¹³⁷

We generate a total of sixty graphs for twenty medi-138 cal concepts using three LLMs: GPT-4, Llama3-70b, 139 and PalmyraMed. We enlist a panel of medical ex-140 perts to score each graph in terms of accuracy and 141 comprehensiveness to the current medical literature. 142 Additionally, we benchmark our graphs against the 143 BIOS KG (Yu et al., 2022) as ground truth. The re-144 sults from expert evaluation indicate that the accu-145 racy of the generated graphs is generally higher than 146 their comprehensiveness. Additionally, both general-147 ist and specialized medical models show a tendency 148 to incorporate public knowledge, which may influence 149 the graphs' content and affect the representation of 150 clinical information. 151

Our contributions can be summarized as follows:

152

164

- We proposed *MedG-KRP* to map the medical ¹⁵³ knowledge embedded in LLMs, aiming to enhance explainability. ¹⁵⁵
- Analyze LLMs' understanding of causal pathways in medicine by utilizing both human reviewers and comparison to a current state-ofthe-art biomedical KG.
- We observe an interesting where medical finetuned models performed unexpectedly worse in human evaluations, despite being specialized for the domain. 160

2. Related Works

Biomedical knowledge graphs are designed to 165 integrate and categorize extensive medical concepts 166 and their interrelationships. Bodenreider (2004) pro-167 posed the Unified Medical Language System (UMLS), 168 which categorizes hundreds of thousands of medical 169 concepts and millions of relationships between these 170 concepts. KGs vary in scope; while the UMLS is quite 171 general, databases such as Orphanet focus specifically 172 on rare diseases (Weinreich et al., 2008). Biomed-173 ical KGs can be generated in various ways. Some 174 have used probabilistic models to extract data from 175 patient notes (Rotmensch et al., 2017), others use 176 named entity recognition and other NLP techniques 177 (Yu et al., 2022), and many are built by reconcil-178 ing a set of various sources (Chandak et al., 2023). 179 We use Yu et al. (2022)'s biomedical informatics on-180

tology system (BIOS) as ground truth to compare
LLM-generated graphs to.

LLMs and Graphs Causal graphs have found 183 use in general medicine (Greenland and Brum-184 back, 2002), epidemiology (Greenland et al., 1999), 185 and bioinformatics (Kleinberg and Hripcsak, 2011). 186 LLMs have been shown to find pairwise relationships 187 (Kiciman et al., 2023), accurately determine edge di-188 rection, (Naik et al., 2023), hypothesize missing vari-189 ables (Sheth et al., 2024), and be capable of generat-190 ing small causal graphs with reasonable accuracy and 191 efficiency (Long et al., 2024; Jiralerspong et al., 2024) 192 and large KGs from texts (Hao et al., 2022; Melnyk 193 et al., 2022; Zhang et al., 2023). LLMs have been 194 also combined with statistical methods for generation 195 (Ban et al., 2023; Abdulaal et al., 2024; Vashishtha 196 et al., 2023). One reason why LLMs are so appealing 197 for graph generation is that they are able to lever-198 age metadata similarly to how a human expert would 199 go about generating a causal graph (Kiciman et al., 200 2023; Abdulaal et al., 2024; Choi et al., 2022). Aug-201 menting LLM with KGs have been shown to improve 202 task performance (Jin et al., 2023; Soman et al., 2023; 203 Jiang et al., 2023b). To the best of our knowledge, our 204 work is the first work to build a complete graph from 205 one given concept (going beyond pairwise compari-206 son and partial graphs), which is used for evaluating 207 LLMs for medical use. 208

²⁰⁹ 3. Methodology

210 3.1. Preliminaries

A knowledge graph can be mathematically denoted as 211 G = (V, E), where V defines a finite set of vertices or 212 nodes and E is a set of ordered pairs of vertices. The 213 vertices, V is represented as $\{v_1, v_2, \ldots, v_n\}$, with 214 each v_i signifying a distinct entity or concept within 215 the graph. The cardinality of V, denoted |V|, indi-216 cates the total number of entities or concepts rep-217 resented in the graph. The edges are denoted as 218 $\{(v_i, v_j) \mid v_i, v_j \in V \text{ and } (v_i \neq v_j)\}, \text{ where each pair}$ 219 (v_i, v_j) represents a directed edge from node v_i to 220 node v_i . The presence of an edge (v_i, v_i) signifies a 221 relationship or interaction between the entities rep-222 resented by v_i and v_j . 223

224 3.2. MedG–KRP

We introduce an algorithm, based on the process of sequentially expanding from a given medical concept, for the generation of biomedical KGs using227LLMs. After generations are complete, a panel of228medical students scores each graph based on accuracy229and comprehensiveness. We also compare our LLM-230generated KGs to the biomedical KG BIOS, comput-231ing precision and recall.232

3.3. Generation Algorithm

233 234

240

268

We divide our graph generation process into two primary stages: **node expansion** and **edge refinement**. In the first stage, nodes are recursively hypothesized by querying the LLM for relevant medical concepts, while the second stage involves validating and refining the edges between these nodes. 239

3.3.1. Node Expansion

Our node expansion algorithm (Algorithm 1) aims 241 to explore the causal relationships between medical 242 concepts. The process begins with a root node r, 243 representing an initial medical concept, and recur-244 sively prompts an LLM for concepts that are either 245 caused by or cause the root concept. The objective 246 of this stage is to identify which medical concepts the 247 LLM associates with r, thereby capturing the model's 248 understanding of the causal pathways surrounding a 249 given medical condition or concept. 250

Formally, let r denote the root node, and i represent the current recursion depth. We expand the graph G by exploring both forward (causal) and backward (caused-by) relationships. The algorithm proceeds recursively, with each newly identified node vbeing further expanded to find related concepts. 256

To prevent unbounded expansion and ensure the 257 graphs remain interpretable, we impose a maximum 258 recursion depth of $d_{\text{max}} = 2$. Additionally, to main-259 tain legibility and minimize the risk of hallucination, 260 we limit the LLM to returning at most $n_{\text{max}} = 3$ con-261 cepts in response to each query. Importantly, there is 262 no lower bound on the number of concepts an LLM 263 may return; the LLM can indicate that there are no 264 concepts either causing or caused by a given node, 265 which helps maintain the algorithm's reliability and 266 reduces over-expansion. 267

3.3.2. Edge refinement

In the second stage (Algorithm 2), we perform an exhaustive check for additional causal connections that the LLM may infer should exist between the concepts already present in the graph. This step is crucial for 270 271 272 273

Algorithm 1: Recursive Node Exploration

```
EXPAND-FROM-NODE(r, i, d)
if i > maximum depth then
| return
end
if d = \rightarrow then
    Ask LLM for a maximum of n concepts, a_1 \ldots a_n,
   caused by r
   G := G + a_1 \dots a_n
   for c in a_1 \ldots a_n do
        Add edge r \to c to G
        Expand-From-Node(c, i + 1, \rightarrow)
       Expand-From-Node(c, i + 1, \leftarrow)
   \mathbf{end}
end
if d = \leftarrow then
    Ask LLM for a maximum of n concepts, b_1 \dots b_n,
   causing r
    G := G + b_1 \dots b_n
   for c in b_1 \ldots b_n do
        Add edge c \to r to G
        Expand-From-Node(c, i+1, \rightarrow)
        Expand-From-Node(c, i + 1, \leftarrow)
   end
end
```

ensuring the completeness of the knowledge graph by 273 identifying all potential relationships between nodes. 274 Given a pair of nodes a and b, we query the LLM to 275 determine if a directed edge $a \rightarrow b$ exists. Similarly, 276 the reverse direction $b \rightarrow a$ is also queried, treating 277 these two directions as distinct. This approach allows 278 for the possibility of bidirectional edges, representing 279 mutual causality or interdependence in medical con-280 texts. 281

Let G denote the graph of concepts obtained af-282 ter the expansion stage. For each pair of distinct 283 nodes $i, j \in G$, where $i \neq j$, we query the LLM 284 for the existence of a directed edge $i \rightarrow j$. If the 285 LLM confirms that such an edge should exist, it is 286 added to the graph. This process is repeated for ev-287 ery pair of nodes and for both directions, ensuring 288 that the graph captures all potential causal relation-289 ships based on the LLM's understanding. 290

We opt to query each direction separately, rather than including all possible edge directions in a single query, in order to reduce the cognitive load on the LLM. By isolating each query to a single direction, we hypothesize that the LLM can provide more accurate predictions regarding the presence of specific edges. Algorithm 2: Edge RefinementEDGE-REFINEMENT (G)foreach $i \in G.nodes()$ doforeach $j \in G.nodes()$ doif $i \neq j$ thenQuery LLM if edge $i \rightarrow j$ existsif LLM confirms $i \rightarrow j$ thenAdd edge $i \rightarrow j$ to Gendendend

4. Experimental Setup

4.1. Concept Selection

We selected twenty conditions from various subdisciplines of medicine to act as the root nodes for our graphs. We chose a list of conditions that would vary vastly in prevalence and level of study. We include both conditions with clear causal pathways and unclear ones. A full list of root concepts, verified by a board-certified physician, can be found in Table 1.

4.2. Models

We tested our benchmark on diverse models-the pro-307 priety GPT-4 model (OpenAI, 2023), open source 308 Llama3-70b (Dubey and et.al., 2024), and finally 309 the current state-of-art medical model PalmyraMed-310 70b (Writer Engineering Team, 2024). PalmyraMed-311 70b is a Llama base model fine-tuned for medical us-312 age which displays very good performance on medical 313 LLM benchmarks. We aimed to compare the perfor-314 mance of a medical finetune model in comparison to 315 its base model counterpart. 316

4.3. Hyperparameters

We run Algorithm 1 in both directions for the graph 318 G - exploring concepts that either cause or are caused 319 by the root medical concept r. Starting with iteration 320 i = 1, we limit the maximum depth of recursion to 321 2, meaning the EXPAND-FROM-NODE() function calls 322 itself only once. After completing Algorithm 1, the 323 graph G will contain all relevant nodes. To identify 324 additional directed edges between the concepts in G, 325 we then execute Algorithm 2. 326

All models are evaluated with a temperature setting of 0.05 and a top_p value of 1.0. The low tem-

```
298
```

297

```
306
```

317

329 perature ensures that results primarily reflect the 330 models' reasoning abilities, enhancing reproducibil-

 $_{\rm 331}\,$ ity. However, we do not set the temperature to $0.0\,$

332 to allow for slight variations in responses during re-

³³³ prompting, should formatting issues arise.

334 4.4. Prompting

In this paper, we aim to evaluate the ability of LLMs
to hypothesize knowledge graphs using their zeroshot prompting abilities. Three main prompts were
used, one system prompt and one general prompt for
each algorithm.

System Prompt The system prompt (see Ap-340 pendix \mathbf{B} , section \mathbf{B} .1) is designed to enhance the 341 LLM's reasoning ability by focusing on distinguishing 342 direct and indirect causality—an area where LLMs 343 often struggle. To improve response quality, we in-344 struct the model to employ counterfactual reasoning, 345 asking it to evaluate causal relationships by consid-346 ering hypothetical scenarios. 347

Expansion Prompts Two expansion prompts (see 348 Appendix B, Section B.2, Section B.3) are used in 349 Algorithm 1 to discover concepts related to the root 350 node, one for "causes" and one for "caused by." Sim-351 ilar to the system prompt, we emphasize counterfac-352 tual reasoning to improve accuracy. We employ a 353 zero-shot chain-of-thought (CoT) approach, following 354 Kojima et al. (2023), which enhances performance in 355 medical QA tasks and pairwise edge-checking. The 356 current graph state is passed into each prompt to 357 maintain context during the expansion process. 358

Edge Check Prompt The edge check prompt (see 359 Appendix B, Section B.4), as used in Algorithm 2, 360 queries the LLM to determine if a directed causal re-361 lationship exists between two medical concepts. Like 362 the system and expansion prompts, we emphasize dis-363 tinguishing direct from indirect causality using coun-364 terfactual reasoning. To isolate the causal connec-365 tion, the prompt assumes no external risk factors are 366 influencing the relationship, ensuring that the LLM 367 focuses on the specific medical concepts being tested. 368

369 4.5. Metrics

Human Evaluation: Graph Accuracy and
Comprehensiveness We enlisted a panel of medical students to manually comment on and score all
generated graphs in terms of accuracy and comprehensiveness. We defined accuracy as medical correct-

ness of all concepts, relationships, and implied causal 375 pathways in a given graph. Comprehensiveness re-376 ferred to how comprehensive a graph was to the cur-377 rent medical understanding of the causal pathways 378 surrounding a disease. Thus, graphs with many miss-379 ing nodes that would be present in an ideal graph (of 380 the current medical understanding) would have low 381 comprehensiveness scores. Accuracy was scored from 382 a scale of 1-4; [Completely accurate (4), Mostly Accu-383 rate (3), Inaccurate (2), Completely inaccurate (1)]. 384 Comprehensiveness was scored similarly, on the scale 385 [Completely Comprehensive (4), Mostly Comprehen-386 sive (3), Poorly Comprehensive (2), Not At All Com-387 prehensive (1)]. Three reviewers scored each graph. 388 We report all reviewer scores and an average thereof 389 for each graph. 390

Ground Truth Comparison Generated graphs 391 were also compared to the Biomedical Informatics 392 Ontology System (BIOS). BIOS is a large knowledge 303 graph composed from numerous sources and contain-394 ing hundreds of thousands of nodes and edges. We 395 chose it because BIOS appeared more complete and 396 suitable for our use case than other biomedical KGs. 397 We calculated the precision and recall of generated 398 edges using algorithm 3. For each generated graph, 399 we iterate through all edges and check if there is a 400 short path (≤ 5) between the two corresponding con-401 cepts in the ground truth. If a path in the ground 402 truth satisfies this condition, it is marked as a hit. 403 Otherwise, it is marked as a miss. The intent of 404 checking for paths instead of a direct edge is to avoid 405 the case where an edge in a generated graph would 406 be deemed as medically correct but have intermedi-407 ary concepts in the ground truth. 408

Mapping node names We directly match BIOS 409 graph and LLM generated graphs by building a vec-410 tor database with embeddings of all BIOS parent con-411 cepts using the sentence transformers model e5-base-412 v2 (Wang et al., 2022). We retrieved the five near-413 est neighbors in the vector database, then prompted 414 GPT as to which (if any) names of the nearest neigh-415 bor nodes matched a given node in an LLM generated 416 graph in meaning. The prompt used can be found in 417 Appendix B, Section B.5. We iterated through all 418 generated nodes and created a json file with every 419 LLM generated node and its BIOS counterpart. If a 420 node had no counterpart, we simply set its value to 421 "none". 422

450

476

477

478

480

482

484

486

Algorithm 3: Precision and recall **Input:** $G = (V_G, E_G)$, the generated graph. $G' = (V_{G'}, E_{G'})$, the BIOS graph. d=5, the path length threshold **Output:** precision: accuracy of predicted edges recall: completeness of predicted edges Let $G'' = (V_{G''}, E_{G''}) \subseteq G'$, where $V_{G''} = V_G \bigcap V_{G'}$ and $E_{G''} = \{ (i'', j'') \in E_{G''} : (i', j') \in E_{G'} \}.$ foreach $(i, j) \in E_G$ do if $\{i, j\} \subseteq V_{G''}$ and $\exists P(i,j) \subseteq E_{G''}$ s.t. $|P(i,j)| \leq d$ then $\mid n_{\rm hit} \rightarrow n_{\rm hit} + 1$ end end Precision = $n_{\rm hit}/|E_G|$ $\text{Recall} = n_{\text{hit}} / |E_{G''}|$

Edge types While all edges in MedG-KRP gener-423 ated graphs specify "cause", edges in BIOS contain 424 many specific relationship labels. We do not recog-425 nize any edges in BIOS labeled as "is a" or "reverse 426 is a" due to the fact that both are used only for sub-427 classes or superclasses of a given concept and because 428 of performance constraints. A consequence of this is, 429 due to BIOS's incompleteness, some nodes are not 430 reachable. The remaining edges are a mix of bidirec-431 tional and directional edges, so we interpret all edges 432 as bidirectional for the sake of consistency, and due 433 to the fact that the directions of all edges in BIOS are 434 implied by their labels, rather than explicitly stated. 435

5. Results 436

5.1. Overview 437

We generated sixty graphs across three models 438 for twenty different conditions from various fields of 439 medicine. We observe that all LLMs perform gen-440 erally well in terms of average reviewer scores 441 (see Table 1). GPT-4 displays the strongest perfor-442 mance in the human review, while PalmyraMed dis-443 plays the weakest. Human reviewers generally found 444 that PalmyraMed's graphs are more specific than 445 those generated by Llama3–70b and GPT–4. Even 446 for the same model, generated graphs have a wide 447 variety of density values, reciprocity values, and sim-448 ple cycle counts. 449

5.2. Human Evaluation

Accuracy, as rated by human reviewers, is gener-451 ally strong, with all averages of all reviewer scores 452 for each model being between 3 and 4, "mostly accu-453 rate" and "completely accurate" (see Table 1). Com-454 prehensiveness scores range from just under 3 to 4. 455 We attribute comprehensiveness scores consistently 456 being lower than accuracy scores to us limiting re-457 sponses and recursion depth in our recursive node 458 exploration algorithm (see Algorithm 1). 459

GPT-4 performed best in accuracy, with an aver-460 age accuracy score across all graphs of 3.37 (see Ta-461 ble 1. Llama3 was close behind with an average ac-462 curacy of 3.28 and PalmyraMed displayed the worst 463 performance with a score of 3.13. Both Llama3 and 464 PalmyraMed performed similarly in comprehensive-465 ness, with average comprehensiveness scores across 466 all graphs of 3.00 and 2.97 (see Table 1. GPT-4 467 displayed the best comprehensiveness, with a score 468 of 3.23–a significantly stronger performance than all 469 other models. 470

In their comments, reviewers mentioned that 471 PalmyraMed's graphs were generally more specific 472 than those of GPT-4 and Llama3-70b. We specu-473 late this to be a result of PalmyraMed being aligned 474 for medical usage. 475

Llama3–70b having weaker overall performance than GPT–4 follows its generally weaker performance on traditional QA benchmarks. PalmyraMed, however, has been shown to have better average perfor-479 mance on QA benchmarks than GPT-4, yet it performed worse overall on our benchmark. Reviewers 481 noticed that PalmyraMed appeared much more prone to hallucination than other models, with it naming 483 multiple graph nodes "myra-med" or "PalmyraMed", and having trouble with instruction following 485

5.3. Ground Truth Comparison

We observe notable results in the ground truth com-487 parison metric, where models demonstrated behavior 488 nearly opposite to that observed in human evalua-489 tions (see tables 2, 1). PalmyraMed performed ex-490 ceptionally well, with the highest precision and recall 491 scores across the board. In particular, PalmyraMed 492 displayed more than three times the average recall 493 score of GPT-4, which displayed the worst perfor-494 mance. Interestingly enough, Llama3–70b, which is 495 usually surpassed by GPT-4 on almost all major QA 496 benchmarks, outperformed GPT-4 in both precision 497 and recall in objective evaluation. 498

	Llam	a3–70b	Palm	yraMed	GF	PT-4	Average	
	Acc.	Comp.	Acc.	Comp.	Acc.	Comp.	Acc.	Comp.
Acute flaccid myelitis	3.67	3.00	2.67	3.33	3.33	3.67	3.22	3.33
Arthritis	$3.00 \\ \pm 0.00$	3.33 ± 0.33	2.67	$3.00 \\ \pm 1.00$	3.33	3.33 + 33	3.00	3.22
Asthma	3.33	$3.00 \\ \pm 1.00$	2.33	$3.00 \\ \pm 1.00$	3.33	4.00	3.00	3.33
Creutzfeldt–Jakob disease	2.67	$2.67 \\ \pm 0.33$	2.67 ± 0.33	$3.33 \\ \pm 0.33$	3.33 +0.33	$3.00 \\ \pm 0.00$	2.89	3.00
Dementia	$3.33 \\ \pm 0.33$	$3.33 \\ \pm 1.33$	$3.33 \\ \pm 0.33$	$2.33 \\ \pm 0.33$	4.00	3.67 +0.33	3.56	3.11
Diabetes Mellitus	$4.00 \\ \pm 0.00$	3.67 ± 0.33	$3.67 \\ \pm 0.33$	$3.00 \\ \pm 1.00$	$3.00 \\ \pm 0.00$	$2.83 \\ \pm 0.58$	3.56	3.17
Esophageal achalasia	$3.00 \\ \pm 0.00$	$3.00 \\ \pm 0.00$	2.67	$3.00 \\ \pm 1.00$	3.67	$3.00 \\ \pm 1.00$	3.11	3.00
Glioblastoma	2.67	2.00 ± 1.00	3.00 + 0.00	$3.00 \\ \pm 1.00$	2.67	3.33 +1.33	2.78	2.78
HIV	$3.33 \\ \pm 0.33$	$2.33 \\ \pm 0.33$	$3.33 \\ \pm 0.33$	3.00 ±1.00	$3.33 \\ \pm 0.33$	$2.33 \\ +2.33$	3.33	2.56
Hyperparathyroidism	$3.33 \\ \pm 0.33$	2.67	3.00 + 0.00	3.00 + 1.00	4.00	3.00 +1.00	3.44	2.89
Ischemic Stroke	3.67	4.00	4.00	$3.00 \\ \pm 0.00$	3.67	2.67 +2.33	3.78	3.22
Lung Cancer	3.67	$2.33 \\ \pm 0.33$	3.67	2.67 ± 0.33	$3.00 \\ \pm 0.00$	4.00	3.44	3.00
Malignant neoplasms of liver	3.67	3.33	3.67	3.33 +1.00	4.00	2.67 ± 0.33	3.78	2.89
Myocardial infarction	3.33	2.67	4.00	3.00 ±0.00	3.33	3.00 ±1.00	3.56	2.89
Myocarditis	3.67	3.00 ±1.00	$3.33 \\ \pm 0.33$	2.67 ± 0.33	3.00 + 1.00	3.00 ±1.00	3.33	2.89
Parkinson's disease	$3.00 \\ \pm 0.00$	3.00 ±0.00	3.33	$2.33 \\ \pm 0.33$	3.00 + 1.00	3.00 ±1.00	3.11	2.78
Renal artery stenosis	3.33	3.67	3.67	3.67	4.00	$3.33 \\ \pm 0.33$	3.67	3.56
SARS-CoV-2	3.00	$3.33 \\ \pm 0.33$	3.00 + 1.00	3.67 ± 0.33	3.33	3.67	3.11	3.56
Spontaneous coronary artery dissection	3.00	$3.00 \\ \pm 3.00$	2.00	2.67	3.00	3.67	2.67	3.11
Ulcerative colitis	$3.00 \\ \pm 1.00$	3.00 ± 1.00	$2.67 \\ \pm 0.33$	$2.67 \\ \pm 0.33$	$3.00 \\ \pm 1.00$	$3.33 \\ \pm 1.33$	2.89	3.00
Average Score	3.28	3.00	3.13	2.97	3.37	3.23		
Average Variance	0.42	0.72	0.43	0.57	0.32	0.76		

Table 1: Mean Reviewer Scores (from 1–4) per Graph per Model

Table 2: Average Precision and Recall per Model

	Llama3-70b		Palmy	raMed	GPT-4		
	Pre.	Rec.	Pre.	Rec.	Pre.	Rec.	
Mean	.201	.012	.243	.033	.163	.011	
Min.	.000	.000	.026	.004	.018	.003	
Max.	.486	.034	.527	.393	.359	.031	
SD	.123	.008	.150	.085	.106	.007	

499 5.4. Graph Attributes

⁵⁰⁰ Overall graph node and edge counts (see Ta-⁵⁰¹ ble 11) varied between models and between graphs. ⁵⁰² PalmyraMed was generally the most conservative ⁵⁰³ when creating nodes and edges while GPT-4 was the least, possibly contributing to PalmyraMed's low comprehensiveness.

504

505

We observe an inverse relationship between per-506 formance based on reviewer scores and average reci-507 procity, density, and simple cycle count across all 508 graphs for a given model. GPT-4 displayed the high-509 est performance and the lowest reciprocity, density, 510 and simple cycle counts, while PalmyraMed displayed 511 the highest. Llama3–70b's values for these three met-512 rics were in-between those of the other models. 513

Simple cycle counts for graphs for a given model varied widely. Each model consistently displayed one or two graphs which were outliers in terms of simply cycle count. PalmyraMed had the most extreme outlier, with its graph for "Malignant neoplasms of liver" containing more than one million cycles, while

all other graphs contained under 2500 and all graphs 520 in the bottom 75th percentile contained less than 521 one hundred. Llama3's graph for "creutzfeldt-jakob 522 disease" has a simple cycle count of greater than 523 3500 and all other graphs display cycle counts less 524 than 225. GPT–4 displays the most reasonable cycle 525 counts out of all models in its generated graphs, with 526 one outlier of 340 and a bottom 50th percentile of less 527 than or equal to eight cycles. 528

529 5.5. Direct and Indirect Causality

Reviewers found that PalmyraMed often had dif-530 ficulty distinguishing direct and indirect causality. 531 Some reviewers mentioned that PalmyraMed often 532 listed nodes as "causes" that would be much more 533 appropriately labeled as "risk factors" GPT-4, on the 534 other hand, was observed by reviewers to display the 535 strongest ability to distinguish between directy and 536 indirect causality, an ability crucial in medicine. 537

538 6. Discussion

539 6.1. Conclusions

Our algorithm, MedG-KRP, is able to generate 540 KGs representing the medical reasoning abilities of 541 LLMs. Coupling MedG-KRP with human reviewers 542 allowed insights into model behavior that were not 543 covered by traditional QA benchmarks. We found 544 that PalmyraMed was generally more specific in its 545 reasoning, but also had a weaker understanding of 546 the differences between direct and indirect causality. 547 while GPT-4 covered more broad concepts and was 548 often able to correctly determine between direct and 549 550 indirect causes of concepts.

Although PalmyraMed displayed worse perfor-551 mance in our human review compared to other mod-552 els, its KG—while flawed—was more specific than 553 that of other models. This is supported by Palmyra-554 Med's exceptionally high recall on our KG compar-555 ison task. We hypothesize that PalmyraMed, as a 556 medical model, was trained on similar sources to 557 which BIOS was constructed from than other LLMs 558 we tested. This would lead to more frequent matches 559 between nodes generated by PalmyraMed and BIOS 560 nodes. Since nodes without mappings would have all 561 adjacent edges generated counted as misses, it fol-562 lows that a model that produced nodes more similar 563 to those in BIOS would have much higher recall. 564

⁵⁶⁵ Clinicians may see PalmyraMed's specificity as a ⁵⁶⁶ desirable trait. GPT-4 and Llama3-70b using more vague terms may signal that they are more influ-567 enced by public knowledge than by clinical knowledge 568 since they are generalist models. It is worth noting 569 that models were asked to be particularly specific and 570 to stay to only medical—as opposed to colloquial-571 terminology. A human doctor whose reasoning was 572 based on public discourse over medical understanding 573 would not be trusted. Likewise, although expected, 574 generalist LLMs having less specific KGs may suggest 575 value in aligning models for clinical use. We wish to 576 once again stress that the ability to find these obser-577 vations is possible with our method, but not neces-578 sarily covered by traditional QA benchmarks. 579

6.2. Future Work

Given that reviewers observe generalist models have a 581 better causal reasoning ability compared to the med-582 ical model we tested but are lacking in domain speci-583 ficity, the question of how we can build models that 584 display both of these abilities naturally arises. Fu-585 ture works may seek to supplement the training cor-586 pora of traditional medical models with information 587 on causal inference and causal reasoning to improve 588 models' medical understanding and viability for real-580 world application. 590

580

We also believe that attempting to explore LLMs' ⁵⁹¹ internal KGs that are unrelated to medicine may ⁵⁹² yield interesting results. The topics of KGs could ⁵⁹³ be from any field, and seeing how LLMs' reasoning ⁵⁹⁴ changes when encountering vastly different subjects ⁵⁹⁵ could give deeper insight into LLMs' behaviors. ⁵⁹⁶

Using MedG-KRP or a similar algorithm as a prompting technique may also be possible. An LLM could generate a reasoning graph then be prompted to make inferences or answer questions given the graph it produced like CoT prompting.

Other pathways that may be worthwhile to explore 602 include, in no specific order: exploring the effect of an 603 LLM's training data on its reasoning KGs, using KG 604 generation to determine the effect (if any) of pretrain-605 ing data order on LLM behavior, revising the MedKG 606 algorithm or developing new algorithms to efficiently 607 use directly prompted LLMs for biomedical KG gen-608 eration or repair, and building very large reasoning 609 KGs with LLMs to probe behavior at a larger scale 610 and how and when LLMs connect interdisciplinary or 611 seemingly unrelated concepts. 612

613 References

adamos hadjivasiliou, Ahmed Abdulaal, Nina 614 Montana-Brown, Tiantian He, Avodeji Ijishakin, 615 Ivana Drobnjak, Daniel C. Castro, and Daniel C. 616 Alexander. Causal modelling agents: Causal 617 graph discovery through synergising metadata- and 618 data-driven reasoning. In The Twelfth Interna-619 tional Conference on Learning Representations, 620 2024. URL https://openreview.net/forum?id= 621 pAogR1TBtY. 622

Emily Alsentzer, John R Murphy, Willie Boag,
Wei-Hung Weng, Di Jin, Tristan Naumann, and

625 Matthew B A McDermott. Publicly available clin-

ical BERT embeddings. arXiv [cs. CL], April 2019.

Taiyu Ban, Lyvzhou Chen, Xiangyu Wang, and
Huanhuan Chen. From query tools to causal architects: Harnessing large language models for advanced causal discovery from data, 2023. URL
https://arxiv.org/abs/2306.16902.

⁶³² Olivier Bodenreider. The unified medical language
 ⁶³³ system (umls): integrating biomedical terminology.
 ⁶³⁴ Nucleic acids research, 32(Database Issue):D267–
 ⁶³⁵ 70, 2004.

Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang,
Michael Carbin, and Christopher D Manning.
BioMedLM: A 2.7B parameter language model
trained on biomedical text. arXiv [cs. CL], December 2022.

⁶⁴³ Nicholas Carlini, Daphne Ippolito, Matthew Jagiel⁶⁴⁴ ski, Katherine Lee, Florian Tramer, and Chiyuan
⁶⁴⁵ Zhang. Quantifying memorization across neural
⁶⁴⁶ language models, 2023. URL https://arxiv.org/
⁶⁴⁷ abs/2202.07646.

Payal Chandak, Kexin Huang, and Marinka Zit-648 nik. Building a knowledge graph to enable 649 precision medicine. Scientific Data, 10(1):67, 650 Feb 2023. ISSN 2052-4463. doi: 10.1038/651 s41597-023-01960-3. URL https://doi.org/10. 652 1038/s41597-023-01960-3. 653

Kristy Choi, Chris Cundy, Sanjari Srivastava, and
Stefano Ermon. Lmpriors: Pre-trained language
models as task-specific priors, 2022. URL https:
//arxiv.org/abs/2210.12530.

- Abhimanyu Dubey and et.al. The llama 3 herd of models. arXiv [cs.AI], July 2024.
- Sander Greenland and Babette Brumback. An overview of relations among causal modelling methods. International Journal of Epidemiology, 31(5):1030–1037, 10 2002. ISSN 0300-5771. doi: 10.1093/ije/31.5.1030. URL https://doi.org/ 10.1093/ije/31.5.1030. 665
- Sander Greenland, Judea Pearl, and James M. Robins. Causal diagrams for epidemiologic research. *Epidemiology*, 10(1):37–48, 1999. ISSN 10443983. URL http://www.jstor.org/stable/ 3702180. 670
- Maxime Griot, Jean Vanderdonckt, Demet Yuksel, and Coralie Hemptinne. Multiple choice questions and large languages models: A case study with fictional medical data, 2024.
- Shibo Hao, Bowen Tan, Kaiwen Tang, Bin Ni, Xiyan Shao, Hengzhe Zhang, Eric P Xing, and Zhiting Hu. BertNet: Harvesting knowledge graphs with arbitrary relations from pretrained language models. *arXiv [cs. CL]*, June 2022.
- Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. ClinicalBERT: Modeling clinical notes and predicting hospital readmission. *arXiv [cs. CL]*, April 2019.
- Sachin Jain. Clinical humility and the 683 limits of medical knowledge. https: 684 //magazine.hms.harvard.edu/articles/ 685 clinical-humility-and-limits-medical-knowledme June 2024. Accessed: 2024-9-11. 687
- Lavender Yao Jiang, Xujin Chris Liu, Nima Pour 688 Nejatian, Mustafa Nasir-Moin, Duo Wang, Anas 689 Abidin, Kevin Eaton, Howard Antony Riina, 690 Ilya Laufer, Paawan Punjabi, Madeline Miceli, 691 Nora C Kim, Cordelia Orillac, Zane Schnurman, 692 Christopher Livia, Hannah Weiss, David Kurland, 693 Sean Neifert, Yosef Dastagirzada, Douglas Kondzi-694 olka, Alexander T M Cheung, Grace Yang, Ming 695 Cao, Mona Flores, Anthony B Costa, Yindalon 696 Aphinyanaphongs, Kyunghyun Cho, and Eric Karl 697 Oermann. Health system-scale language models are 698 all-purpose prediction engines. Nature, June 2023a. 699
- Pengcheng Jiang, Cao Xiao, Adam Cross, and Jimeng Sun. GraphCare: Enhancing healthcare predictions with personalized knowledge graphs. *arXiv* 702 [cs.AI], May 2023b. 703

Bowen Jin, Gang Liu, Chi Han, Meng Jiang, Heng Ji, 704 and Jiawei Han. Large language models on graphs: 705 A comprehensive survey. arXiv [cs.CL], December 706

2023.707

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung 708 Weng, Hanyi Fang, and Peter Szolovits. What dis-709 ease does this patient have? a large-scale open do-710 main question answering dataset from medical ex-711 ams, 2020. URL https://arxiv.org/abs/2009. 712 13081. 713

Thomas Jiralerspong, Xiaoyin Chen, Yash More, 714 Vedant Shah, and Yoshua Bengio. Efficient causal 715 graph discovery using large language models, 2024. 716

Samantha Kleinberg and George Hripcsak. А 717 review of causal inference for biomedical in-718 Journal of Biomedical Informatics, formatics. 719 44(6):1102-1112, 2011.ISSN 1532-0464. doi: 720 https://doi.org/10.1016/j.jbi.2011.07.001. URL 721 https://www.sciencedirect.com/science/ 722

article/pii/S1532046411001195. 723

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, 724 Yutaka Matsuo, and Yusuke Iwasawa. Large lan-725 guage models are zero-shot reasoners, 2023. 726

Emre Kıcıman, Robert Ness, Amit Sharma, and 727 Chenhao Tan. Causal reasoning and large language 728 models: Opening a new frontier for causality, 2023. 729

Stephanie Long, Tibor Schuster, and Alexandre 730 Piché. Can large language models build causal 731 graphs?, 2024. 732

- Daniel McDuff, Mike Schaekermann, Tao Tu, Anil 733 Palepu, Amy Wang, Jake Garrison, Karan Sing-734 hal, Yash Sharma, Shekoofeh Azizi, Kavita Kulka-735 rni, Le Hou, Yong Cheng, Yun Liu, S Sara Mah-736 davi, Sushant Prakash, Anupam Pathak, Christo-737 pher Semturs, Shwetak Patel, Dale R Webster, 738 Ewa Dominowska, Juraj Gottweis, Joelle Barral, 739 Katherine Chou, Greg S Corrado, Yossi Matias, 740 Jake Sunshine, Alan Karthikesalingam, and Vivek 741 Natarajan. Towards accurate differential diagnosis 742 with large language models. arXiv [cs. CY], Novem-743 ber 2023. 744
- Igor Melnyk, Pierre Dognin, and Payel Das. Knowl-745 edge graph generation from text. arXiv [cs.CL], 746 November 2022.

747

- Narmada Naik, Ayush Khandelwal, Mohit Joshi, 748 Madhusudan Atre, Hollis Wright, Kavya Kannan, 749 Scott Hill, Giridhar Mamidipudi, Ganapati Srini-750 vasa, Carlo Bifulco, Brian Piening, and Kevin Mat-751 lock. Applying large language models for causal 752 structure learning in non small cell lung cancer, 753 2023. URL https://arxiv.org/abs/2311.07191. 754
- Jesutofunmi A Omiye, Jenna C Lester, Si-755 mon Spichak, Veronica Rotemberg, and Roxana 756 Daneshjou. Large language models propagate race-757 based medicine. NPJ Digital Medicine, 6(1):195, 758 2023.759
- OpenAI. GPT-4 technical report. arXiv [cs.CL], 760 March 2023. 761
- Ankit Pal, Logesh Kumar Umapathi, and Malaikan-762 nan Sankarasubbu. MedMCQA : A large-scale 763 multi-subject multi-choice dataset for medical do-764 main question answering. arXiv [cs.CL], March 765 2022.766
- Maya Rotmensch, Yoni Halpern, Abdulhakim Tli-767 mat, Steven Horng, and David Sontag. Learn-768 ing a health knowledge graph from electronic 769 medical records. Scientific Reports, 7(1):5994, 770 ISSN 2045-2322. Jul 2017. doi: 10.1038/771 s41598-017-05778-z. URL https://doi.org/10. 772 1038/s41598-017-05778-z. 773
- Ivaxi Sheth, Sahar Abdelnabi, and Mario Fritz. Hy-774 pothesizing missing causal variables with llms. 775 arXiv preprint arXiv:2409.02604, 2024. 776
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara 777 Mahdavi, Jason Wei, Hyung Won Chung, 778 Nathan Scales, Ajay Tanwani, Heather Cole-779 Lewis, Stephen Pfohl, Perry Payne, Martin 780 Seneviratne, Paul Gamble, Chris Kelly, Abubakr 781 Babiker, Nathanael Schärli, Aakanksha Chowd-782 hery, Philip Mansfield, Dina Demner-Fushman, 783 Blaise Agüera y Arcas, Dale Webster, Greg S. Cor-784 rado, Yossi Matias, Katherine Chou, Juraj Got-785 tweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, 786 Joelle Barral, Christopher Semturs, Alan Karthike-787 salingam, and Vivek Natarajan. Large language 788 models encode clinical knowledge. Nature, 620 789 (7972):172-180, Aug 2023a. ISSN 1476-4687. doi: 790 10.1038/s41586-023-06291-2. URL https://doi. 791 org/10.1038/s41586-023-06291-2. 792
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara 793 Jason Wei, Hyung Won Chung, Mahdavi, 794

- Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, 795 Stephen Pfohl, Perry Payne, Martin Seneviratne, 796 Paul Gamble, Chris Kelly, Abubakr Babiker, 797 Nathanael Schärli, Aakanksha Chowdhery, Philip 798 Mansfield, Dina Demner-Fushman, Blaise Agüera 799 Y Arcas, Dale Webster, Greg S Corrado, Yossi 800 Matias, Katherine Chou, Juraj Gottweis, Nenad 801 Tomasev, Yun Liu, Alvin Rajkomar, Joelle Bar-802 ral, Christopher Semturs, Alan Karthikesalingam, 803 and Vivek Natarajan. Large language models en-804 code clinical knowledge. Nature, 620(7972):172-805 180, August 2023b. 806
- Karthik Soman, Peter W Rose, John H Mor-807 ris, Rabia E Akbas, Brett Smith, Braian Pee-808 toom, Catalina Villouta-Reyes, Gabriel Cerono, 809 Yongmei Shi, Angela Rizk-Jackson, Sharat Is-810 rani, Charlotte A Nelson, Sui Huang, and Ser-811 gio E Baranzini. Biomedical knowledge graph-812 optimized prompt generation for large language 813 models. arXiv [cs.CL], November 2023. 814
- Ali Soroush, Benjamin S. Glicksberg, Eyal Zimlich-815 man, Yiftach Barash, Robert Freeman, Alexan-816 der W. Charney, Girish N Nadkarni, and Eyal 817 Klang. Large language models are poor medical 818 coders — benchmarking of medical code query-819 NEJM AI, 1(5), 2024.doi: 10.1056/ ing. 820 AIdbp2300040. URL https://ai.nejm.org/doi/ 821 abs/10.1056/AIdbp2300040. 822
- Tao Tu, Anil Palepu, Mike Schaekermann, Khaled 823 Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, 824 Brenna Li, Mohamed Amin, Nenad Tomasev, 825 Shekoofeh Azizi, Karan Singhal, Yong Cheng, 826 Le Hou, Albert Webson, Kavita Kulkarni, S Sara 827 Mahdavi, Christopher Semturs, Juraj Gottweis, 828 Joelle Barral, Katherine Chou, Greg S Corrado, 829 Yossi Matias, Alan Karthikesalingam, and Vivek 830 Natarajan. Towards conversational diagnostic AI. 831 arXiv [cs.AI], January 2024. 832
- Aniket Vashishtha, Abbavaram Gowtham Reddy,
 Abhinav Kumar, Saketh Bachu, Vineeth N Balasubramanian, and Amit Sharma. Causal inference using llm-guided discovery, 2023. URL https:
 //arxiv.org/abs/2310.15117.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by
 weakly-supervised contrastive pre-training. arXiv
 preprint arXiv:2212.03533, 2022.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL https://arxiv.org/abs/2201.11903.
- S S Weinreich, R Mangon, J J Sikkens, M E en Teeuw, and M C Cornel. [orphanet: a european database for rare diseases]. *Ned Tijdschr Geneeskd*, 152(9): 518–519, March 2008.
- Writer Engineering Team. Palmyra-Med-70b: A powerful LLM designed for healthcare. https://dev. writer.com, 2024.
- Jeff Wu, Long Ouyang, Daniel M Ziegler, Nisan Stiennon, Ryan Lowe, Jan Leike, and Paul Christiano. Recursively summarizing books with human feedback. arXiv [cs. CL], September 2021.
- Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv* [cs.CL], September 2023.
- Xi Yang, Aokun Chen, Nima PourNejatian, 864 Hoo Chang Shin, Kaleb E Smith, Christopher 865 Parisien, Colin Compas, Cheryl Martin, Mona G 866 Flores, Ying Zhang, Tanja Magoc, Christopher A 867 Harle, Gloria Lipori, Duane A Mitchell, William R 868 Hogan, Elizabeth A Shenkman, Jiang Bian, and 869 Yonghui Wu. GatorTron: A large clinical language 870 model to unlock patient information from unstruc-871 tured electronic health records. arXiv [cs.CL], 872 February 2022. 873
- Sheng Yu, Zheng Yuan, Jun Xia, Shengxuan Luo, Huaiyuan Ying, Sihang Zeng, Jingyi Ren, Hongyi Yuan, Zhengyun Zhao, Yucong Lin, Keming Lu, Jing Wang, Yutao Xie, and Heung-Yeung Shum. Bios: An algorithmically generated biomedical knowledge graph, 2022. URL https://arxiv. org/abs/2203.09975.
- Yuan Zhang, Xin Sui, Feng Pan, Kaixian Yu, Keqiao 881 Li, Shubo Tian, Arslan Erdengasileng, Qing Han, 882 Wanjing Wang, Jianan Wang, Jian Wang, Donghu 883 Sun, Henry Chung, Jun Zhou, Eric Zhou, Ben Lee, 884 Peili Zhang, Xing Qiu, Tingting Zhao, and Jin-885 feng Zhang. BioKG: a comprehensive, large-scale 886 biomedical knowledge graph for AI-powered, data-887 driven biomedical research. bioRxiv, October 2023. 888

Appendix A. Limitations

While we test our method on a diverse set of mod-890 els, others may show very different behavior. Al-891 though we aimed to make our list of diseases used for 892 graph generation broad, it is by no means compre-893 hensive. Due to its time complexity, our approach is 894 also only suitable for the generation of small graphs-895 sufficient for benchmarking purposes but not for full 896 generation of KGs. Our human review is subjec-897 tive, and only three reviewers go over a given graph. 898 In addition, we found there was often high variance 800 in reviewer opinions. The knowledge graph we use 900 as ground truth, BIOS, may also be quite incom-901 plete. We also only test using one knowledge graph as 902 ground truth. We believe that, in the case that a hu-903 man expert scored every graph edge, precision values 904 would be greater than the ones which we report. 905

⁹⁰⁶ Appendix B. Prompts

907 B.1. System Prompt

You are a helpful assistant for causal 908 inference and causal reasoning about medical 909 questions. You are always specific in your 910 answers. You always format your answers 911 consistently and name all medical terms in 912 the correct and accepted medical lexicon. 913 You understand the differences between 914 direct and indirect causality and 915 acknowledge these differences when 916 formulating an answer. You utilize a 917 counterfactual model of causal inference 918 when formulating a response. 919

920 B.2. Left Expansion Prompt

A directed knowledge graph that you 921 generated is surrounded in XML tags and 922 provided below. This directed knowledge 923 graph is formatted as a list of edges like 924 so: ['a causes b', 'b causes c', etc]. The 925 knowledge graph you generated is as follows: 926 927 <Begin Knowledge Graph> 928 {edges:} 929 </End Knowledge Graph> 930 931 Given the directed knowledge graph above 932 that you generated, up to three factors that 933 directly cause {concept:}. These factors do 934

not need to be in the knowledge graph above, 935 but can be. If a factor you answer with is 936 in the knowledge graph above, in your 937 response, name it exactly as it is named in 938 the graph above. Do not answer with any 939 factors that only indirectly cause 940 {concept:}. In your final answer, surround 941 the medical name of each cause in square 942 brackets characters. Do not include acronyms 943 or abbreviations in your answer. Utilize a 944 counterfactual model of causal inference 945 when formulating a response. Be as specific 946 as possible. Let's think step by step like a 947 medical expert. 948

B.3. Right Expansion Prompt

A directed knowledge graph that you 950 generated is surrounded in XML tags and 951 provided below. This directed knowledge 952 graph is formatted as a list of edges like 953 so: ['a causes b', 'b causes c', etc]. The 954 knowledge graph you generated is as follows: 955

949

956

957

958

959

960

979

<Begin Knowledge Graph> {edges:} </End Knowledge Graph>

Given the directed knowledge graph above 961 that you generated, List up to three medical 962 concepts directly caused by {concept:}. 963 These factors do not need to be in the 964 knowledge graph above, but can be. If a 965 factor you answer with is in the knowledge 966 graph above, in your response, name it 967 exactly as it is named in the graph above. 968 Do not answer with any factors that only are 969 indirectly caused by {concept:}. In your 970 final answer, surround the medical name of 971 each medical concept that {concept:} causes 972 in square brackets characters. Do not 973 include acronyms or abbreviations in your 974 answer. Utilize a counterfactual model of 975 causal inference when formulating a 976 response. Be as specific as possible. Let's 977 think step by step like a medical expert. 978

B.4. Edge Refinement Prompt

Does {node0:} directly cause {node1:}? Your 980 answer must be one of the following: [yes] / 981 [no]. Surround your final [yes] / [no] 982

```
answer in square brackets characters. If
983
    there is only an indirect causal
984
    relationship as opposed to a direct one,
985
    answer with [no]. Utilize a counterfactual
986
    model of causal inference. Assume no other
987
    risk factors are present. Let's think step
988
    by step. Be concise in your response.
989
    B.5. Nearest Neighbor Selection Prompt
990
    Is the concept ['{original}'] identical in
991
    meaning to any of the concepts in the
992
    following list?
993
994
    Concepts: {retrieved}
995
996
    If so, reply with the name of one concept
997
    in the list identical in meaning to
998
    {original} as it is written in the list. If
999
    there is more than one item of the same
1000
    meaning in the list, answer with the
1001
    concept which best fits and which is in
1002
    proper medical lexicon. Provide one and
1003
    only one answer. If no items in the list
1004
    are identical in meaning to {original},
1005
    provide an empty set of square brackets.
1006
    Surround your final answer in square
1007
    brackets characters. It is very important
1008
    that you do this or else your answer will
1009
    not be processed. It is also very important
1010
    that you provide only one answer and your
1011
    answer as it is written in the list.
1012
    .....
1013
```

Appendix C. Additional Tables

Please see the next page for double-column tables. 1015

1014

	Revi	iewer 1	Revi	ewer 2	Revi	ewer 3
	Acc.	Comp.	Acc.	Comp.	Acc.	Comp.
Acute flaccid myelitis	4	4	4	2	3	3
Arthritis	3	4	3	3	3	3
Asthma	3	3	4	2	3	4
Creutzfeldt–Jakob disease	2	3	4	2	2	3
Dementia	3	2	4	4	3	4
Diabetes Mellitus	4	3	4	4	4	4
Esophageal achalasia	3	3	3	3	3	3
Glioblastoma	4	2	2	1	2	3
HIV	3	2	3	2	4	3
Hyperparathyroidism	3	3	4	3	3	2
Ischemic Stroke	3	4	4	4	4	4
Lung Cancer	4	2	3	2	4	3
Malignant neoplasms of liver	3	3	4	3	4	3
Myocardial infarction	3	3	4	1	3	4
Myocarditis	3	3	4	2	4	4
Parkinson's disease	3	3	3	3	3	3
Renal artery stenosis	3	4	4	4	3	3
SARS-CoV-2	3	4	2	3	4	3
Spontaneous coronary artery dissection	3	4	3	1	3	4
Ulcerative colitis	3	2	2	3	4	4

Table 3: All Reviewer Scores for Llama3–70b Generations

Table 4: All Reviewer Scores for PalmyraMed–70b Generations

	Revi	iewer 1	Revi	iewer 2	Reviewer 3	
	Acc.	Comp.	Acc.	Comp.	Acc.	Comp
Acute flaccid myelitis	2	3	3	4	3	3
Arthritis	2	2	3	3	3	4
Asthma	2	2	2	3	3	4
Creutzfeldt–Jakob disease	3	3	3	4	2	3
Dementia	3	3	4	2	3	2
Diabetes Mellitus	4	3	4	4	3	2
Esophageal achalasia	4	4	1	2	3	3
Glioblastoma	3	3	3	2	3	4
HIV	3	4	4	3	3	2
Hyperparathyroidism	3	4	3	3	3	2
Ischemic Stroke	4	3	4	3	4	3
Lung Cancer	4	3	4	2	3	3
Malignant neoplasms of liver	4	4	4	2	3	3
Myocardial infarction	4	3	4	3	4	3
Myocarditis	4	3	3	2	3	3
Parkinson's disease	4	3	3	2	3	2
Renal artery stenosis	4	3	4	4	3	4
SARS-CoV-2	4	3	3	4	2	4
Spontaneous coronary artery dissection	2	3	1	2	3	3
Ulcerative colitis	3	3	2	2	3	3

	Revi	iewer 1	Revi	ewer 2	Revi	ewer 3
	Acc.	Comp.	Acc.	Comp.	Acc.	Comp.
Acute flaccid myelitis	3	3	3	4	4	4
Arthritis	3	3	4	3	3	4
Asthma	3	4	4	4	3	4
Creutzfeldt–Jakob disease	4	3	3	3	3	3
Dementia	4	3	4	4	4	4
Diabetes Mellitus	3	3.5	3	2	3	3
Esophageal achalasia	4	4	3	2	4	3
Glioblastoma	3	4	2	2	3	4
HIV	4	4	3	1	3	2
Hyperparathyroidism	4	3	4	4	4	2
Ischemic Stroke	4	3	3	1	4	4
Lung Cancer	3	4	3	4	3	4
Malignant neoplasms of liver	4	3	4	2	4	3
Myocardial infarction	3	4	3	2	4	3
Myocarditis	4	4	2	2	3	3
Parkinson's disease	4	4	2	3	3	2
Renal artery stenosis	4	3	4	4	4	3
SARS-CoV-2	3	4	4	3	3	4
Spontaneous coronary artery dissection	3	3	3	4	3	4
Ulcerative colitis	4	4	2	2	3	4

Table 5: All Reviewer Scores for GPT-4 Generations

Note: a reviewer answered with "3-4" for the comprehensiveness of GPT–4's graph for Diabetes Mellitus. With their approval, we reported the value as 3.5.

Table 6:	Conditions	Sorted by	Average	Accuracy	and	Comprehensiveness	Across a	ll Graphs
	0 0 0 0 0					• • • • • • • • • • • • • • • • • • •		

Condition (Sorted by Acc.)	Acc. ▼	Comp.	Condition (Sorted by Comp.)	Acc.	Comp. ▼
Ischemic Stroke	3.78	3.22	SARS-CoV-2	3.11	3.56
Malignant neoplasms of liver	3.78	2.89	Renal artery stenosis	3.67	3.56
Renal artery stenosis	3.67	3.56	Acute flaccid myelitis	3.22	3.33
Dementia	3.56	3.11	Asthma	3.00	3.33
Diabetes Mellitus	3.56	3.17	Arthritis	3.00	3.22
Myocardial infarction	3.56	2.89	Ischemic Stroke	3.78	3.22
Lung Cancer	3.44	3.00	Diabetes Mellitus	3.56	3.17
Hyperparathyroidism	3.44	2.89	Dementia	3.56	3.11
Myocarditis	3.33	2.89	Spontaneous coronary artery dissection	2.67	3.11
HIV	3.33	2.56	Esophageal achalasia	3.11	3.00
Acute flaccid myelitis	3.22	3.33	Lung Cancer	3.44	3.00
Esophageal achalasia	3.11	3.00	Creutzfeldt–Jakob disease	2.89	3.00
Parkinson's disease	3.11	2.78	Ulcerative colitis	2.89	3.00
SARS-CoV-2	3.11	3.56	Hyperparathyroidism	3.44	2.89
Arthritis	3.00	3.22	Malignant neoplasms of liver	3.78	2.89
Asthma	3.00	3.33	Myocardial infarction	3.56	2.89
Creutzfeldt–Jakob disease	2.89	3.00	Myocarditis	3.33	2.89
Ulcerative colitis	2.89	3.00	Glioblastoma	2.78	2.78
Glioblastoma	2.78	2.78	Parkinson's disease	3.11	2.78
Spontaneous coronary artery dissection	2.67	3.11	HIV	3.33	2.56

Condition (Sorted by Precision)	Precision \blacktriangledown	Condition (Sorted by Recall)	Recall \blacktriangledown
Ischemic Stroke	0.392919	Myocardial infarction	0.149441
Myocarditis	0.352433	Diabetes Mellitus	0.027118
Acute flaccid myelitis	0.329147	Ulcerative colitis	0.017456
Hyperparathyroidism	0.307151	Glioblastoma	0.016545
Asthma	0.299079	Esophageal achalasia	0.015814
Lung Cancer	0.275498	Arthritis	0.015560
Malignant neoplasms of liver	0.247299	Ischemic Stroke	0.013806
Ulcerative colitis	0.231435	Hyperparathyroidism	0.013094
Glioblastoma	0.225804	Creutzfeldt–Jakob disease	0.012317
HIV	0.216769	Dementia	0.011893
Renal artery stenosis	0.211895	Myocarditis	0.011819
Dementia	0.153292	Asthma	0.011401
Diabetes Mellitus	0.145785	Renal artery stenosis	0.010728
Arthritis	0.143932	Acute flaccid myelitis	0.010599
Creutzfeldt–Jakob disease	0.131771	Lung Cancer	0.010384
Myocardial infarction	0.120545	Malignant neoplasms of liver	0.010237
Esophageal achalasia	0.087591	HIV	0.009080
Spontaneous coronary artery dissection	0.087235	Spontaneous coronary artery dissection	0.007348
SARS-CoV-2	0.082722	SARS-CoV-2	0.004768
Parkinson's disease	0.021010	Parkinson's disease	0.003900

Table 7: Conditions Sorted by Average Precision and Recall Across all Graphs

Table 8: Graph Attributes for GPT–4 Generations, Sorted by Precision

Condition	Precision \blacksquare	Recall	Density	Reciprocity	Nodes	Edges	Cycles
Acute flaccid myelitis	0.359	0.012	0.075	0.085	36	94	5
HIV	0.353	0.011	0.059	0.255	31	55	22
Ischemic Stroke	0.312	0.015	0.071	0.043	37	94	151
Myocarditis	0.309	0.012	0.072	0.088	36	91	137
Ulcerative colitis	0.238	0.031	0.060	0.047	38	85	9
Glioblastoma	0.224	0.016	0.101	0.056	38	142	43
Renal artery stenosis	0.205	0.010	0.065	0.051	25	39	7
Dementia	0.155	0.014	0.102	0.040	32	101	4
Arthritis	0.149	0.018	0.092	0.125	30	80	10
Lung Cancer	0.147	0.006	0.083	0.034	38	117	3
Creutzfeldt–Jakob disease	0.138	0.017	0.130	0.014	34	146	32
Asthma	0.132	0.005	0.078	0.020	36	98	1
Spontaneous coronary artery dissection	0.130	0.009	0.059	0.109	31	55	7
Hyperparathyroidism	0.127	0.007	0.063	0.078	29	51	2
SARS-CoV-2	0.082	0.007	0.083	0.000	26	54	0
Myocardial infarction	0.075	0.022	0.078	0.130	32	77	340
Malignant neoplasms of liver	0.052	0.009	0.053	0.102	34	59	5
Parkinson's disease	0.036	0.003	0.083	0.054	37	111	4
Esophageal achalasia	0.034	0.004	0.057	0.029	35	68	10
Diabetes Mellitus	0.018	0.006	0.091	0.125	33	96	77
Mean	0.164	0.012	0.078	0.074	33.40	85.650	43.450
SD	0.106	0.007	0.019	0.058	3.872	29.462	82.453

Condition	Precision \blacktriangledown	Recall	Density	Reciprocity	Nodes	Edges	Cycles
Hyperparathyroidism	0.486	0.017	0.089	0.191	33	94	205
Malignant neoplasms of liver	0.375	0.013	0.121	0.183	32	120	26
Myocarditis	0.371	0.011	0.085	0.214	32	84	20
Ischemic Stroke	0.359	0.011	0.080	0.143	30	70	6
Lung Cancer	0.242	0.008	0.061	0.083	35	72	5
Glioblastoma	0.239	0.019	0.083	0.194	34	93	13
Asthma	0.238	0.010	0.125	0.091	27	88	27
Ulcerative colitis	0.211	0.009	0.052	0.178	30	45	7
HIV	0.202	0.004	0.101	0.338	27	71	98
Diabetes Mellitus	0.191	0.029	0.060	0.036	31	56	1
Renal artery stenosis	0.181	0.011	0.062	0.051	36	78	35
Acute flaccid myelitis	0.179	0.007	0.106	0.174	26	69	60
Myocardial infarction	0.155	0.034	0.091	0.125	27	64	5
Arthritis	0.145	0.013	0.075	0.180	35	89	40
Esophageal achalasia	0.142	0.023	0.049	0.129	36	62	14
Dementia	0.118	0.011	0.048	0.000	35	57	1
Creutzfeldt–Jakob disease	0.112	0.013	0.111	0.427	31	103	3553
SARS-CoV-2	0.089	0.004	0.117	0.146	27	82	10
Parkinson's disease	0.000	0.000	0.095	0.064	32	94	109
Spontaneous coronary artery dissection	0.000	0.000	0.066	0.135	34	74	5

Table 9: Graph Attributes for Llama3–70b Generations, Sorted by Precision

Table 10: Graph Attributes for PalmyraMed–70b Generations, Sorted by Precision

Condition	Precision \blacktriangledown	Recall	Density	Reciprocity	Nodes	Edges	Cycles
Asthma	0.527	0.019	0.075	0.143	31	70	11
Ischemic Stroke	0.508	0.016	0.070	0.264	28	53	35
Acute flaccid myelitis	0.449	0.013	0.134	0.253	26	87	2199
Lung Cancer	0.438	0.017	0.069	0.094	31	64	3
Myocarditis	0.377	0.012	0.068	0.182	29	55	5
Malignant neoplasms of liver	0.315	0.008	0.153	0.349	34	172	1106539
Hyperparathyroidism	0.308	0.015	0.113	0.343	31	105	1448
Renal artery stenosis	0.250	0.011	0.082	0.204	25	49	12
Ulcerative colitis	0.246	0.013	0.074	0.308	27	52	18
Diabetes Mellitus	0.228	0.046	0.058	0.074	31	54	14
Glioblastoma	0.214	0.015	0.081	0.240	31	75	25
Dementia	0.187	0.010	0.108	0.123	25	65	137
Creutzfeldt–Jakob disease	0.145	0.007	0.140	0.262	25	84	75
Arthritis	0.138	0.016	0.073	0.000	28	55	4
Spontaneous coronary artery dissection	0.132	0.013	0.128	0.338	23	65	83
Myocardial infarction	0.131	0.393	0.108	0.123	25	65	60
HIV	0.095	0.013	0.136	0.187	24	75	125
Esophageal achalasia	0.087	0.021	0.069	0.281	31	64	13
SARS-CoV-2	0.077	0.004	0.195	0.189	17	53	8
Parkinson's disease	0.027	0.009	0.171	0.329	22	79	62

	Llama3–70b		Palmy	raMed	$\mathbf{GPT}-4$		
	nodes	edges	nodes	edges	nodes	edges	
Mean	31.5	78.25	27	72.05	33.4	85.65	
Min.	26	45	17	49	25	39.	
Max.	36	120	34	172	38	146	
SD	3.32	17.93	4.09	27.49	3.87	29.46	

Table 11: Graph Node and Edge Counts per Model