

LLMs for Causal Reasoning in Medicine? A Call for Caution

Saurabh Mathur*[†]
TU Darmstadt
Darmstadt, Germany
saurabh.mathur@tu-darmstadt.de

Predrag Radivojac
Northeastern University
Boston, USA
predrag@northeastern.edu

Ranveer Singh*
University of Texas at Dallas
Richardson, USA
ranveer.singh@utdallas.edu

David M. Haas
Indiana University School of Medicine
Indianapolis, USA
dahaas@iu.edu

Michael Skinner
University of Texas at Dallas
Richardson, USA
mas140130@utdallas.edu

Lakshmi Raman
University of Texas Southwestern
Medical Center
Dallas, USA
lakshmi.raman@utsouthwestern.edu

Sriraam Natarajan
University of Texas at Dallas
Richardson, USA
sriraam.natarajan@utdallas.edu

Abstract

Large Language Models (LLMs) are increasingly being used for medical advice by patients and healthcare providers. These models capture knowledge from their training data, which consists of vast medical corpora. However, they lack the ability to use this knowledge to causally reason about the underlying physiological processes. Moreover, they are unable to deal with uncertainty, generating responses that are confidently presented yet factually incorrect. Acting on such factually incorrect medical advice can be dangerous. Mitigating these risks requires rethinking the role of LLMs in medicine. In this work, we present an evaluation scheme for LLMs in three roles: direct clinical decision support, exact medical knowledge base, and approximate medical knowledge base. We evaluate six LLMs on two clinical studies, in obstetrics and pediatric critical care, respectively. Our results indicate that LLMs are much better suited to the approximate knowledge base role. Based on these observations, we request caution when directly employing LLMs in safety-critical domains such as medicine.

CCS Concepts

• **Computing methodologies** → **Natural language processing**;
Causal reasoning and diagnostics; • **Applied computing** →
Health informatics.

Keywords

Large Language Models, Causal Bayesian Networks, Theory Refinement

*Equal contribution

[†]Work done while at The University of Texas at Dallas



This work is licensed under a Creative Commons Attribution 4.0 International License.
CODS 2025, Pune, India

© 2025 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-2355-1/25/12

<https://doi.org/10.1145/3799830.3799848>

ACM Reference Format:

Saurabh Mathur, Ranveer Singh, Michael Skinner, Predrag Radivojac, David M. Haas, Lakshmi Raman, and Sriraam Natarajan. 2025. LLMs for Causal Reasoning in Medicine? A Call for Caution. In *13th ACM IKDD International Conference on Data Science (CODS 2025)*, December 17–20, 2025, Pune, India. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3799830.3799848>

1 Introduction

Large language models (LLMs), such as Generative Pre-trained Transformer (GPT) and Gemini, have generated significant interest in their potential to assist or even replace aspects of medical practice, with speculation regarding their ability to generate differential diagnoses and treatment plans and especially to reduce administrative burdens [26]. Indeed, internal medicine residents already perceive a current and future role for LLMs in medicine and use these tools in professional settings, often without formal guidance [8]. Moreover, patients are increasingly using LLMs to obtain medical advice [17]. However, despite the compelling fluency of LLM-generated text, LLMs cannot reason [14, 49]. Rather, LLMs are designed to mimic human utterances by identifying linguistic patterns from large corpora; they lack an explicit logical or causal reasoner, a stark contrast to the way clinicians manage patients.

Medical practice is fundamentally rooted in sophisticated cognitive processes, and especially causal reasoning [19]. Clinicians develop causal models to understand physiological mechanisms, evaluate hypotheses, construct explanations, and devise physiological interventions. Such causal concepts are very challenging for LLMs; they struggle to perform reasoning tasks not represented in their training data, such as simple math problems involving infrequently used numbers [47]. In practice, this results in “hallucinations” or “confabulations”: coherent and confident yet factually incorrect statements; these are particularly dangerous in a medical context where accuracy is critical. In one report soliciting advice from an LLM to manage a serious infection, the LLM suggested dangerously incorrect management plans contradicting clinical guidelines [39].

The persuasive nature of LLM outputs can also exploit human automation bias, potentially leading clinicians to over-rely on machine suggestions and make errors. Moreover, inserting factually incorrect text directly into medical records could diminish the quality of information, impede clinical reasoning, and even hinder the development of future AI tools [26]. Given these substantial risks and limitations in areas critical for causal reasoning, it is dangerous and unethical to rely on current LLMs to diagnose and manage human disease. However, it is likely that less ambitious tasks may be assigned to LLMs to aid in medical practice and clinical decision support.

To this effect, we evaluate the efficacy of pre-trained LLMs as approximate sources of causal knowledge, focusing on two clinical studies in obstetrics and pediatric critical care. Specifically, we propose a three-stage evaluation scheme for such systems, consisting of pairwise question answering, full causal graph construction, and the refinability of the constructed causal graph. Our evaluation shows that while LLMs struggle to answer causal questions, their answers can be used as initial hypotheses to construct models more amenable to causal reasoning, such as Causal Bayesian Networks (CBNs).

The rest of the paper is organized as follows: first, we provide the necessary background for the paper on LLMs, causality, and theory refinement, along with the clinical tasks that we evaluate the LLMs on. We then present the approach taken for this evaluation before presenting our empirical evaluation and the key findings. We then conclude the paper by discussing areas for future research.

2 Background

2.1 Large Language Models

Large Language Models (LLMs [30]) are a class of generative models that represent the probability distribution over natural language text using neural networks, typically based on the transformer architecture [45]. They are distinguished from other language models, such as Hidden Markov Models (HMMs [37]), by their unprecedented scale, with some of the largest models boasting hundreds of billions of trainable parameters. This large size allows them to capture intricate statistical patterns from large corpora of natural language text.

LLMs can be categorized into encoder-only, decoder-only, and encoder-decoder models. Of these, decoder-only models have demonstrated impressive capabilities across a wide range of natural language processing tasks. These models process text in one direction, modeling its probability autoregressively, that is, the probability of each word is conditioned on all the words before it. This autoregressive structure allows decoder-only LLMs to be efficiently trained on vast amounts of unlabeled text through simple tasks such as next-word prediction, allowing them to generate highly realistic text. As a result, decoder-only LLMs, including models fine-tuned on medical data, demonstrate impressive performance on medical benchmarks such as question answering, clinical note summarization, patient report generation, and diagnostic reasoning [23, 39, 41].

While LLMs can generate clinically relevant and accurate text that mimics causal reasoning, they do not perform true causal inference. LLMs rely on statistical correlations rather than causal

understanding, making them prone to blending genuine causal relationships with spurious associations [49]. Moreover, the decoder-only architecture makes them inherently stochastic and prone to cascading errors [12, 27]. These limitations – an inability to reason, stochasticity, and cascading errors – result in these models generating confident-sounding yet factually incorrect text, especially about topics less represented in training data. This phenomenon is referred to as a hallucination or a confabulation. Since confabulations are the result of inherent limitations of decoder-only LLMs, commonly used mitigation strategies such as Retrieval-Augmented Generation (RAG [21]) are inadequate; they may even result in unsafe text generation [2].

2.2 AI-in-the-loop

LLM use in medicine can be analyzed by characterizing the nature of human-AI interaction. Since clinical practice requires the human clinician to be the primary decision-maker [5], it is an example of an AI-in-the-loop domain [32]. Here, the AI system’s role is to support the clinician by providing them with accurate and actionable information. The efficacy of such a system depends on its ability to improve the clinician’s decision-making, such as by alleviating their cognitive load by automating mechanical aspects of clinical reasoning.

However, LLMs’ use in medicine can potentially deviate from the AI-in-the-loop paradigm. Fig. 1 depicts this deviation. Clinicians could treat LLMs as an expert system, asking them questions that typically require causal inference. Since LLMs are unable to reason, stochastic, and prone to generating erroneous output, the human clinician would need to validate the LLM-generated answers, increasing their burden [15]. Moreover, since the erroneous output is often confidently phrased, there is a distinct possibility that clinicians under pressure might miss one or more errors, which might carry over in their ultimate decision. Mitigating these risks requires a fundamental rethinking of the role LLMs might play in clinical decision support.

2.3 Causal Bayesian Networks

LLMs can be contrasted with another class of generative models called Causal Bayesian Networks (CBNs [34]); unlike LLMs, CBNs are interpretable, amenable to causal reasoning, and naturally deal with uncertainty, satisfying the desiderata for AI-in-the-loop in medicine. CBNs are closely related to causal diagrams, which clinicians have used for causal reasoning [19].

CBNs are a subclass of Bayesian Networks (BNs). BNs represent the joint probability distribution over a set of variables by factorizing it over a directed acyclic graph (DAG). This DAG consists of nodes corresponding to each variable; each directed edge between two variables denotes direct influence. If the edges also denote direct causal relationships, then the BN is considered a CBN. In a CBN, each edge $X \rightarrow Y$ means that X is a cause of Y . These causal edges can be interpreted interventionally: intervening on the component corresponding to X should change the distribution over Y . Such targeted interventions might not be reasonable in some cases, such as when modeling the effect of the family history of a medical condition; in such cases, the edges can be interpreted historically or etiologically [9].

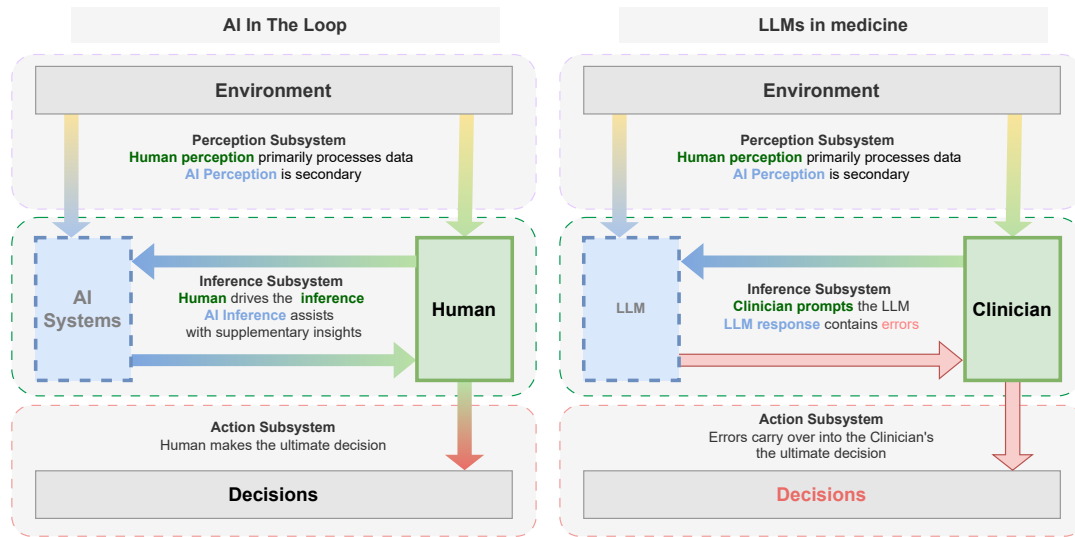


Figure 1: In AI-in-the-loop domains (left), humans must make the ultimate decisions, while AI systems assist with perception, inference, and action. The fluency of LLM-generated text blurs the lines between AI assistance and human decision-making. As a result, using LLMs as AI-in-the-loop systems in medicine (right) risks introducing new errors to clinical practice. In response to a difficult query, the LLM might generate confident-sounding text with errors that the human clinician might fail to catch, carrying them over in their final decision.

	nuMoM2b	PELICAN
Subfield	Obstetrics	Pediatric Critical Care
Time scale	8 to 9 months	< 1 month
Condition(s) considered	Adverse pregnancy outcomes	Neurological injury on life support
Rarity of condition(s)	Common (15% of US pregnancies)	Extremely rare (20% out of < 2500 cases a year in the US)
Existing research	2.7M results on google scholar	34k results on google scholar
Subject inclusion criteria	First-time mothers (Nulliparous) without pregestational diabetes	Pediatric patients supported by ECMO, but not having congenital heart disease
Average age of subjects	27.79 years	4.32 years
Risk factors	Demographics, existing conditions, family history, and lifestyle factors recorded at start of pregnancy	Abnormal events identified from high-frequency physiological measurements and laboratory test results up to 24 hours on ECMO
Data set size	3,856	71

Table 1: Comparison of the two medical domains and the corresponding subsets considered in this work

CBNs can be constructed by eliciting them from domain experts or clinical guidelines such as the Quick Medical Reference (QMR [40]). This approach falls short when modeling medical conditions that are less well-understood, such as rare diseases and conditions involving complex causal relationships. As a result, considerable research has been performed to devise ways to generate causal graphs from observational data [10].

Data-driven causal discovery methods typically use large amounts of data to exclude non-causal edges and rely on several assumptions to decide the causal direction of the remaining edges; examples of such methods include Peter-Clark (PC [42]), Greedy Equivalence

Search (GES [18]), and Fast Causal Inference (FCI [42]). The assumptions used by these methods include the causal Markov condition, which states that each variable is independent of its non-effects (non-descendants) given its direct causes; faithfulness, which states that any conditional independencies in the data arise from the structure of the causal graph itself; and causal sufficiency, which requires all common causes of observed variables to be included in the dataset.

Data-driven causal discovery in medicine is challenging. First, the complex dynamics underlying medical domains make it difficult

to make assumptions like causal sufficiency without expert knowledge. Second, human physiology is naturally cyclic, while causal graphs represent acyclic relationships [6]. This requires the causal modeling domain and its variables to be carefully designed. Finally, expensive data annotation, the temporal nature of the data, and the missingness of variables due to data collection issues further complicate data-driven causal discovery in medicine [16, 38, 48]. CBN construction in medical domains requires a hybrid approach, combining expert domain knowledge with empirical patterns.

2.4 Theory Refinement

Theory refinement is one such hybrid method that combines expert knowledge with data-driven learning. In it, an expert-specified model is refined to better fit empirical evidence [31]. This method has been used to improve the structure of BNs derived from incomplete or imperfect domain knowledge [4], such as LLM output [24, 25]. Refining a BN involves adding, removing, or reversing edges from it to maximize a score measuring the empirical validity of the relationships between variables. Commonly used scores include Bayesian-Dirichlet (BD [11]) and the Minimal Description Length (MDL [20]) scores. The MDL score has been used in causal discovery since it approximates the Kolmogorov complexity of the causal graph [13, 28]. It consist of two components – the log-likelihood of the data under the CBN ($\mathcal{L}(\mathcal{M}, \mathcal{D})$) and the description length cost of the CBN ($Cost(\mathcal{M})$), $Score(\mathcal{M}) = \mathcal{L}(\mathcal{M}, \mathcal{D}) - Cost(\mathcal{M})$.

3 Assessment of LLMs for Causal Reasoning

LLMs excel at capturing statistical patterns from large amounts of textual data, allowing them to synthesize highly coherent text and achieve impressive performance on medical benchmarks. This has generated significant interest in using them for clinical decision support, as AI-in-the-loop. However, these models are stochastic and lack a reasoner, making their responses a blend of accurate and inaccurate information. We aim to empirically evaluate the performance of these models on real-world medical domains.

3.1 Data sets

We consider two medical domains: nuMoM2b and PELICAN. nuMoM2b is an obstetrics domain, based on a study that aims to understand Adverse Pregnancy Outcomes (APOs) in nulliparous subjects (first-time mothers). APOs are common, affecting 15% of U.S. pregnancies, and the study covers a time scale of 8 to 9 months. In contrast, PELICAN is based on a pediatric critical care study that aims to understand neurological injury in pediatric subjects supported by Extracorporeal Membrane Oxygenation (ECMO). This condition is extremely rare, and the study deals with a much shorter period, encompassing the duration during which the patient was on life support (less than one month). The dataset size for PELICAN is considerably smaller, with only 71 subjects, and existing research on neurological injury in pediatric patients on ECMO is less abundant, with 34,000 results on Google Scholar, 80 times fewer results than APOs. Table 1 summarizes the differences between the two domains, and tables 2 and 3 summarize the specific variables considered for our evaluation. These variables were selected by our domain experts.

Variable	Value	% subjects
Age	≤ 21	15.9%
	21-35	77.0%
	> 35	7.0%
BMI	≤ 18	1.3%
	18-25	54.1%
	> 25	44.6%
Race	Non-Hispanic Asian	4.2%
	Non-Hispanic Black	11.2%
	Non-Hispanic White	67.8%
	Hispanic	12.3%
	Others	4.43%
DiabHist	TRUE	20.7%
HTNHist	TRUE	45.4%
HiBP	TRUE	2.7%
PCOS	TRUE	4.8%
METS	TRUE	66.1%
Smoking	TRUE	14.9%
PREc	TRUE	5.9%
NewHTN	TRUE	17.7%
PTB	TRUE	7.7%
GDM	TRUE	3.8%
Total		3,856

Table 2: Risk Factors and Outcomes for nuMoM2b. We consider four adverse outcomes: Preeclampsia (PREc), New Hypertension (NewHTN), Gestational Diabetes Mellitus (GDM), and Preterm Birth (PTB). For these adverse outcomes, we consider nine risk factors: Age, Body Mass Index (BMI), Race, Family History of Diabetes and Hypertension (DiabHist and HTNHist, respectively), Hypertension (HiBP), Polycystic Ovary Syndrome (PCOS), physical activity measured in Metabolic Equivalents of Time (METs), and Smoking in the three months before start of pregnancy.

Variable	Value	% subjects
HighVIS	TRUE	21.1%
Hypotension	TRUE	23.9%
Hypertension	TRUE	4.2%
LowPlatelet	TRUE	32.4%
HighLactate	TRUE	59.2%
LowpH	TRUE	9.86%
RelativepCO2	TRUE	29.6%
NeurologicalInjury	TRUE	23.9%
Total		71

Table 3: Risk Factors and Outcomes for PELICAN. We consider the adverse outcome of Neurological Injury. We consider seven of its risk factors: High Vasoactive-inotropic score (HighVIS), Hypotension, Hypertension, Low Platelets, High Lactate, Low pH, as well as the high relative change in pCO2 24 hours post-cannulation compared to pre-cannulation levels.

3.2 Evaluation scheme and metrics

To empirically assess each LLM’s performance on these real-world medical domains, we use a three-stage evaluation scheme. First, to assess the LLM’s ability to answer **direct causal questions**, we prompt it with queries about every pair of variables. We use their answers to construct a causal graph and evaluate this graph. For the second stage, we evaluate the **LLM as an exact knowledge source** and prompt it to construct a full causal graph from the given list of variables. Finally, the third stage evaluates the **LLM as an approximate knowledge source**, focusing on the refinability of the LLM-constructed causal graph. Here, we combine the LLM-generated graph with indirect expert knowledge like anticausal relations based on temporal order, and refine it using empirical data to further eliminate incorrect edges. We evaluate this refined graph.

Each of these graphs is compared against graphs constructed by our domain experts. Since both domains are being actively researched, these expert graphs do not fully capture all the causal relationships, but they do capture known causal relationships, to the best of our experts’ knowledge.

We quantify the differences between each candidate causal graph \mathcal{M} and the expert-constructed graph \mathcal{G} using three criteria. First, we use the **Structural Hamming Distance (SHD)** [1], which measures the number of edge additions, deletions, or reversals required to transform \mathcal{M} into \mathcal{G} . SHD is conceptually similar to the Levenshtein (edit) distance [29], providing an intuitive and computationally efficient way to compare causal graph structures. However, SHD has two key limitations: it treats all structural differences the same, and it fails to account for the differences in the causal conclusions drawn from the two graphs [46].

We address these limitations of SHD by considering two additional measures. First, we define spurious edges as the edges that are present in the candidate graph but not in the expert graph; we use this to compute the **number of spurious edges (SE)** as $|\text{Edges}(\mathcal{M}) \setminus \text{Edges}(\mathcal{G})|$, where Edges outputs the edge set of a given graph. Second, we quantify the difference in causal conclusions by using the **Structural Intervention Distance (SID)** [35]. SID captures the number of incorrect inferences about intervention effects made by a candidate graph, as compared to the expert graph. Unlike SHD and SE, which only capture structural differences, SID evaluates whether two graphs encode the same causal dependencies under interventions, thus providing a more semantically meaningful notion of causal accuracy. Finally, note that three measures would be 0 if the candidate graph is identical to the expert graph, that is, $\text{SID}(\mathcal{G}, \mathcal{G}) = \text{SHD}(\mathcal{G}, \mathcal{G}) = \text{SE}(\mathcal{G}, \mathcal{G}) = 0$. [36].

4 Results of Empirical Evaluation

We now present the results of our empirical investigation¹ and try to answer the following questions

- (1) **How do LLMs perform as causal question answering systems?**
- (2) **How do LLMs perform as exact knowledge bases?**
- (3) **How do LLMs perform as approximate knowledge bases?**

¹Additional details on the experimental setup—including data preprocessing, LLM prompts, and responses—can be found in the supplementary material: <https://github.com/s-ranveer/LLM-Causal-Medicine-Eval>

We evaluate our results using 10 bootstrap samples for each dataset. To construct each LLM’s representative causal graph, the model was prompted five times with the same full prompt for the full-graph generation. For the question-based graph construction, each question was similarly posed five times. These responses were aggregated to construct a DAG, adding edges in decreasing order of frequency across the five runs, excluding any edge that would introduce a cycle. Ties in edge frequency were resolved lexicographically based on the name of the source node.

Results are presented in Table 6 for six LLMs: Claude [3], DeepSeek [7], Gemini [43], GPT-4o [33], LLaMA [44], and the LLM fine-tuned on medical text, OpenBioLLM (OBLLM) [22]. For the data-only baselines, we consider the constraint-based Peter-Clark (PC) algorithm that starts with a fully connected undirected graph and uses statistical independence tests to remove or orient edges. Additionally, we consider the score-based Greedy Search and Score (GSS) algorithm, which is based on Greedy Equivalence Search (GES), where we evaluate graph structures by optimizing a score function such as the Bayesian Information Criterion (BIC). Finally, we consider Fast Causal Inference (FCI), which is designed to handle latent confounders and learn causal features that remain consistent across all graphs in an equivalence class.

4.1 LLMs for Causal Question Answering

To evaluate whether LLMs can be used for causal reasoning in obstetrics and pediatric critical care, we posed questions about the causal relationships between pairs of variables to the LLMs using the prompt illustrated in Figure 2 (left). The LLM responses were then compared to an expert-constructed causal graph. Table 6 presents the results of this experiment.

Both domains present unique challenges to the LLMs. For pediatric critical care, the limited availability of relevant literature on ECMO likely restricts the LLM’s ability to identify accurate causal relationships from its training data. Conversely, the obstetrics domain has a lot more literature, but it involves a significantly larger number of variables. As a result, asking pairwise causal questions without sufficient contextual information results in spurious associations, driven by hidden confounders. This issue is evident in the pairwise results for both domains, where the performance is close to purely data-driven baselines. The best-performing models in these scenarios tend to be those that responded more conservatively, affirming fewer causal relationships like Claude and Deepseek, as seen in table 4, resulting in lower SHD, SID, and SE values.

Additionally, the prompt responses often exhibit inconsistencies in determining the direction of causal relationships between variable pairs. When asked using the pairwise prompt format, LLMs frequently respond affirmatively to both directions—i.e., $A \rightarrow B$ and $B \rightarrow A$ —thereby introducing cycles. This can be seen in Table 4, which presents the number of edges deleted to enforce acyclicity for each case. Therefore, LLMs by themselves perform poorly as causal question answering systems.

4.2 LLMs as Exact Knowledge Bases

To evaluate the potential of LLMs as medical knowledge bases, we provided each model with a prompt for one-shot full causal graph construction, including the domain description, variable definitions,

<p>Extracorporeal Membrane Oxygenation, or ECMO for short, is an advanced therapy that is sometimes used to work the heart and lungs when a patient’s organs are too sick or weak to work on their own. It is effectively a modified heart-lung bypass machine—a machine that takes over heart and lung function (meaning it adds oxygen to and removes carbon dioxide from a patient’s blood supply).</p> <p>For such patients, is there a causal edge from HighLactate to NeurologicalInjury? Here, HighLactate is lactate acidosis in the first 24 hours of the ECMO run, and NeurologicalInjury is neurological injury confirmed after the end of the ECMO run. Provide the answer as a single word, Yes or No (with No also for cases when the answer is unknown).</p>	<p>## Setup You are a knowledge engineer working on a study on mitigating the risk of Adverse Pregnancy Outcomes (APOs). Think carefully and logically, explaining the reasons for your answer.</p> <p>## Available Information Your team has collected variables representing clinical and demographic information. The study data is collected on Nulliparous women and consists of variables representing clinical and demographic features that might influence the risk of Adverse Pregnancy outcomes. The variables being considered are as follows: {Variable_Descriptions}</p> <p>## Task Make a list of direct causal relations between these variables and explain the rationale behind each decision. Please provide the answers in the format: Variable 1 -> Variable 2, followed by an explanation on a different line.</p>
---	---

Figure 2: Prompts used for pairwise (left) and full causal graph elicitation (right) for PELICAN and nuMoM2b, respectively

Domain	Method	SHD	SID	SE
PELICAN	GSS	9.5 ± 1.6	19.2 ± 4.6	4.2 ± 1.5
	PC	8.5 ± 1.1	18.7 ± 4.0	1.2 ± 1.0
	FCI	8.0 ± 0.5	14.9 ± 0.3	0.1 ± 0.3
nuMoM2b	GSS	33 ± 1.5	90 ± 5.5	10.8 ± 1.1
	PC	33.7 ± 1.7	91.9 ± 7.5	7.5 ± 2.4
	FCI	31.8 ± 1.2	80.6 ± 4.1	2.3 ± 1.7

Table 5: Data-driven causal discovery PELICAN and nuMoM2b; the difference is quantified in terms of Structural Hamming Distance (SHD), Structural Interventional Distance (SID), and the number of spurious edges (SE).

Domain	LLM	Deleted/Total	
		Pairwise	Full
PELICAN	Claude	0/4	2/13
	Deepseek	1/13	3/17
	Gemini	13/35	2/17
	GPT 4o	1/10	1/11
	LLaMA	20/46	7/25
	OpenBioLLM	17/43	1/9
nuMoM2b	Claude	1/18	0/32
	Deepseek	0/27	0/31
	Gemini	0/40	0/34
	GPT 4o	1/37	0/25
	LLaMA	32/95	1/32
	OpenBioLLM	39/99	2/43

Table 4: The number of edges deleted to eliminate cycles during causal graph construction, for both prompt types and across PELICAN and nuMoM2b domains.

	LLM	LLM output			Subtractive refinement		
		SHD	SID	SE	SHD	SID	SE
Pair	Claude	6	6	1	6.3 ± 0.5	11.9 ± 1.2	0
	DeepSeek	9	7	7	8 ± 1.2	13.4 ± 2.2	1.8 ± 1.1
	Gemini	21	8	18	14.3 ± 1.6	11.6 ± 1.6	9 ± 1.5
	GPT 4o	9	14	6	6.6 ± 0.9	12.8 ± 1.9	0.4 ± 0.5
	LLaMA	24	17	23	14.3 ± 1.4	16.7 ± 3.9	8.7 ± 1.2
	OBLLM	23	11	22	14.8 ± 3	14.6 ± 2.5	8.5 ± 2.5
Full	Claude	4	5	4	4.5 ± 0.7	5 ± 1.2	2.1 ± 0.3
	DeepSeek	6	0	6	4.8 ± 0.9	4.2 ± 1.3	3.2 ± 0.4
	Gemini	9	10	9	5.9 ± 1.0	6.5 ± 2.4	3.5 ± 0.7
	GPT 4o	8	15	6	7 ± 1.0	12 ± 3.7	1.6 ± 0.9
	LLaMA	14	12	13	8.3 ± 1.0	7.1 ± 1.7	5.1 ± 0.9
	OBLLM	9	19	5	7.8 ± 1.2	14.3 ± 1.5	0.9 ± 0.9
	LLM	LLM output			Subtractive refinement		
		SHD	SID	SE	SHD	SID	SE
Pair	Claude	23	63	6	27.5 ± 0.8	66.9 ± 0.7	4.8 ± 0.6
	DeepSeek	23	56	10	22.1 ± 0.7	56.3 ± 1.2	7.5 ± 0.5
	Gemini	32	49	21	30 ± 0.5	49.7 ± 1.1	17.6 ± 4.8
	GPT 4o	25	45	15	23.8 ± 0.4	41.6 ± 1.8	13.7 ± 0.5
	LLaMA	45	52	41	24.6 ± 0.9	37.9 ± 3.3	15.8 ± 1.3
	OBLLM	43	35	38	31.8 ± 1.0	26.7 ± 0.6	23.5 ± 0.5
Full	Claude	17	44	9	18.5 ± 0.8	49.5 ± 0.8	8.2 ± 0.6
	DeepSeek	21	50	10	22.7 ± 1.3	53.1 ± 1.1	9.3 ± 1.0
	Gemini	16	38	9	16.1 ± 0.9	39.2 ± 1.7	8.7 ± 0.6
	GPT 4o	20	53	5	21.9 ± 0.8	58.5 ± 3.4	4.1 ± 0.3
	LLaMA	26	53	11	26.7 ± 1.5	55.6 ± 1.5	9.6 ± 0.5
	OBLLM	32	54	22	26.2 ± 1.0	45 ± 4.4	13.5 ± 0.5

Table 6: Evaluation results comparing graphs constructed by each LLM-based method to corresponding expert graphs on the PELICAN (top) and nuMoM2b (bottom) domains; the difference is quantified in terms of Structural Hamming Distance (SHD), Structural Interventional Distance (SID), and the number of spurious edges (SE).

and the overall task description. As in the pairwise evaluation, the generated graphs were compared against expert-established causal structures. Table 6 presents the results of full causal graph construction in both domains.

While the domain-specific challenges discussed earlier are not entirely resolved by providing the full set of available variables, they are significantly mitigated. As a result, LLM performance improves notably when using the full prompt compared to pairwise prompting. As in pairwise prompting, conservative models, such

as Claude and DeepSeek, produced fewer edges, performing better than the other models.

Full prompt-based causal graphs outperform data-driven causal discovery methods. However, the number of spurious edges, SHD, and SID is still too high to trust LLM-generated graphs fully. Therefore, while LLMs show promise, they cannot yet be relied upon as standalone exact medical knowledge bases.

4.3 LLMs as Approximate Knowledge Bases

To evaluate the use of LLMs as approximate knowledge bases, we refine their outputs using data and indirect domain knowledge about temporally impossible edges. The refinement procedure deletes edges to maximize the MDL score. Table 6 shows the results after refining the LLM-generated graphs.

For the pediatric critical care domain, we see a reduction in the SHD and the number of spurious edges (SE) across most LLMs, both in pairwise and full prompting, with a reduction in SID for some of the LLMs. However, in the obstetrics domain, performance improves only for the pairwise prompt. For the full prompt, only OpenBioLLM shows significant improvement, while the others have similar or slightly worse performance. This indicates the lower refinability of LLM-generated causal graphs in the obstetrics domain.

Overall, refinement appears to be more effective for the PELICAN domain than for nuMoM2b. This is likely due to the relatively limited literature available on pediatric critical care, which limits the LLM’s exposure during training. As a result, the model is more prone to generate non-causal edges that are non-associational, and hence easier to remove through the refinement process. In contrast, obstetrics is a well-studied domain with a broad body of research covering diverse populations. Moreover, obstetrics-related discussions are more prevalent in public discourse, often drawing from a mix of high- and low-quality sources. As a result, LLM outputs in this domain may include specific causal claims that are either not credible or not applicable to the nuMoM2b study population. This can lead to suboptimal refinement, including the unintended removal of valid edges from the graph.

Despite these limitations, the graphs constructed from LLM-responses are more accurate than those discovered from the limited data using algorithms like GSS, PC, and FCI, as seen in table 5. The difference becomes more pronounced after refinement, especially for LLMs that output a lot of causal edges like OpenBioLLM and LLaMA. Therefore, LLMs have utility as approximate knowledge sources.

4.4 Qualitative Evaluation

Figure 3 shows the expert-constructed graphs (left), along with the ones obtained by first prompting Gemini to create a full causal graph and then refining it using domain knowledge and data (right). Spurious edges that persist after refinement are labeled in **red**. These edges were generated by the LLM and persisted after refinement, indicating strong empirical correlation. However, correlation does not necessarily imply causation. So, these edges can be categorized into at least two types:

(1) Edges **oversimplifying a complex causal** relationship. The LLM-generated edge $X \rightarrow Y$ persists after refinement because of the correlation between X and Y due to a more complex causal

relationship, which is missing in the LLM-generated graph, such as one mediated by another factor U ($X \rightarrow U, U \rightarrow Y$).

(2) Edges **representing (potentially) correct causal** relationships. The LLM-generated edge $X \rightarrow Y$ persists after refinement because there is either a direct causal relationship or one mediated by an unobserved factor. This can be the case when the refined graph covers all the other known indirect relationships between X and Y , along with the direct causal edge.

In the nuMoM2b domain, the refined graph correctly identifies the two causes of Hypertension at the start of pregnancy (HiBP): age at the start of the pregnancy (Age) and level of physical activity (METs). But, it also identifies four additional causes: Race, family history of hypertension (HTNHist), body mass index at the start of pregnancy (BMI), and polycystic ovary syndrome (PCOS). Of these, the causal relations involving PCOS and BMI are complicated by the edge $PCOS \rightarrow BMI$, which is missing in the refined graph. Similarly, for preeclampsia (PreEc), the refined graph identifies PCOS as an additional cause beyond the other eight causes in the expert graph: Race, Age, BMI, smoking in three months before the start of pregnancy (Smoking), Gestational Diabetes (GDM), HTNHist, HiBP, and new hypertension (NewHTN). This is also complicated by the missing $PCOS \rightarrow BMI$ edge. Finally, for NewHTN, the refined graph identifies METs as an additional cause, beyond its three causes: HTNHist, HiBP, and BMI.

In the PELICAN domain, the refined graph correctly identifies the four causes of Neurological Injury: Low pH, High Lactate, Low Platelet, and Hypertension. It also identifies three additional causes: Relative pCO₂, Hypotension, and High VIS. Of these, Relative pCO₂ and Hypotension might be genuine direct causes, since the refined graph already captures their indirect influence through Low pH and High Lactate, respectively. On the other hand, High VIS is less likely to be a direct cause of Neurological Injury, since the refined graph fails to capture its indirect influence through Hypotension.

In both domains, while the refined graphs do not exactly match the expert-constructed graphs, they contain fewer spurious edges than the LLM-generated graphs. Moreover, some of the spurious edges that persist in the refined graphs might be genuinely causal. Hence, while LLMs perform poorly as causal expert systems, they might be powerful tools for causal discovery, especially when their approximate knowledge is combined with empirical data and indirect expert knowledge. As approximate knowledge sources, LLMs can reduce the time and effort required to build reliable causal expert systems such as Causal Bayesian Networks, enabling their wider adoption.

5 Conclusion

We considered the use of LLMs in medical practice. These models capture intricate statistical patterns from vast medical corpora to generate fluent text, achieving high performance on medical benchmarks. Indeed, there has been significant interest in their potential to assist or even replace aspects of medical practice. However, their lack of an explicit causal reasoner, along with their stochasticity, raises concerns about their suitability as clinical decision support systems. We proposed an evaluation scheme to evaluate LLMs in three different roles in the clinical decision support pipeline. We evaluated six LLMs on two medical domains. Our results indicate that while LLMs do capture medical domain knowledge from

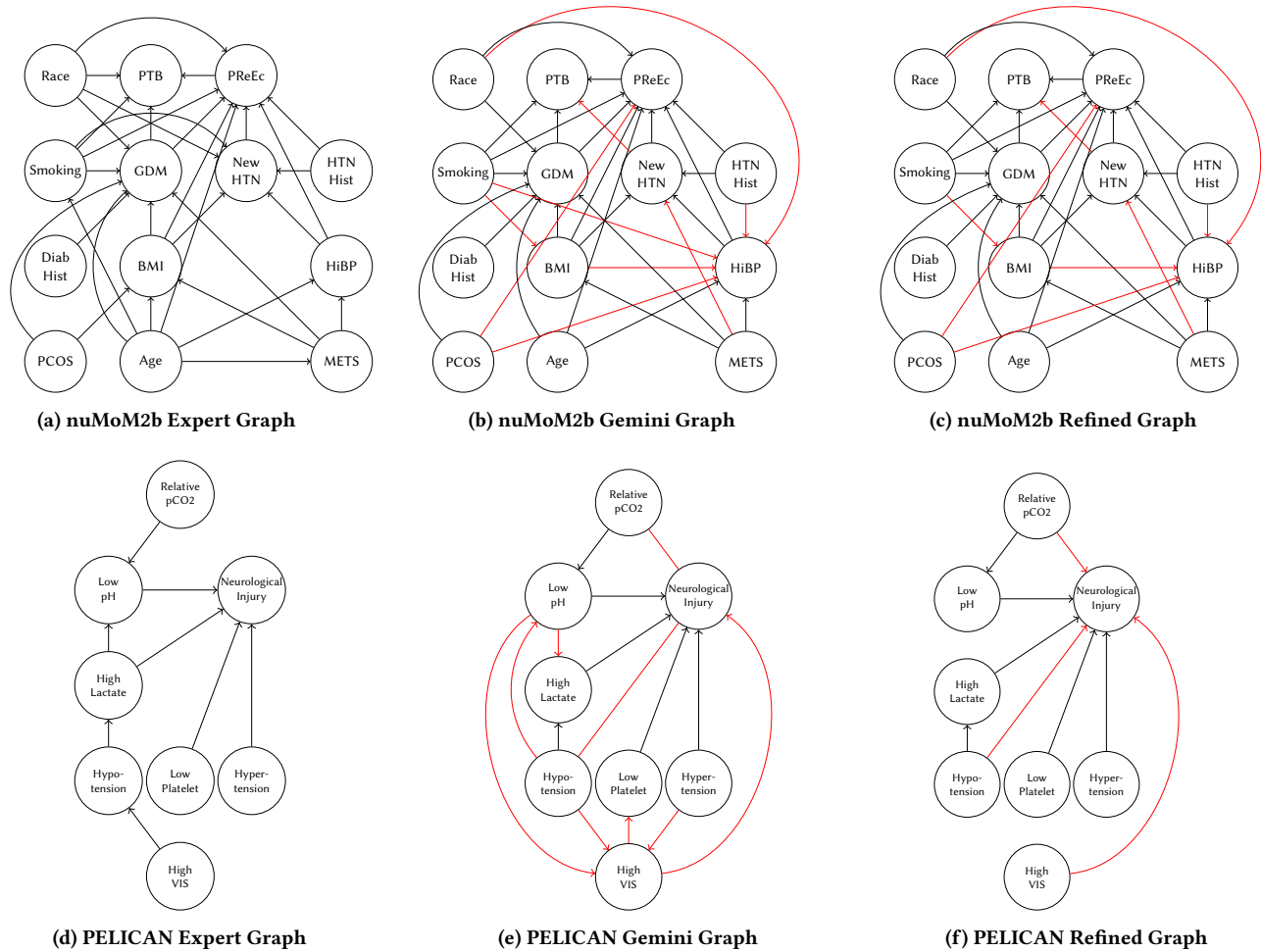


Figure 3: Causal Bayesian networks for nuMoM2b (a–c) and Pelican (d–f). (a) and (d) show the networks specified by domain experts. (b) and (e) show networks obtained by prompting an LLM (Gemini). (c) and (f) show networks obtained by refining the corresponding LLM-elicited networks. Spurious edges are shown in red; these are the edges that were present in the LLM-elicited or refined graphs but not in the corresponding expert graph.

their training data, they fail to accurately answer causal questions. LLM-use requires caution, especially in high-stakes domains like medicine, but these models might be used as approximate knowledge sources to construct models more amenable to causal reasoning, like Causal Bayesian Networks.

There are a number of directions for future work. First, this evaluation can be extended to additional medical domains that differ in the depth of existing research and data availability, encompassing both well-studied and niche areas. Second, the scope of the evaluation can be expanded to include ensembles of LLMs. Finally, imposing validity constraints on the LLM responses, such as acyclicity and partial order of nodes, could reduce variance and potentially improve the quality of the generated graphs.

Acknowledgments

The authors gratefully acknowledge the support of NIH grants R01NS133142 and R01HD101246.

References

- [1] S. Acid and L. M. De Campos. 2003. Searching for Bayesian network structures in the space of restricted acyclic partially directed graphs. *JAIR* 18 (May 2003), 445–490.
- [2] Bang An, Shiyue Zhang, and Mark Dredze. 2025. RAG LLMs are Not Safer: A Safety Analysis of Retrieval-Augmented Generation for Large Language Models. arXiv:2504.18041 [cs.CL] <https://arxiv.org/abs/2504.18041>
- [3] Anthropic. 2023. Claude: An AI Assistant. <https://www.anthropic.com/claude>. Accessed: 2025-05-26.
- [4] Wray Buntine. 1991. Theory refinement on Bayesian networks. In *Uncertainty in Artificial Intelligence*. Elsevier, 52–60.
- [5] Benjamin Chin-Yee and Ross Upshur. 2018. Clinical judgement in the era of big data and predictive analytics. *Journal of Evaluation in Clinical Practice* 24, 3 (2018), 638–645.

- [6] Jurgen AHR Claassen, Dick HJ Thijssen, Ronney B Panerai, and Frank M Faraci. 2021. Regulation of cerebral blood flow in humans: physiology and clinical implications of autoregulation. *Physiological reviews* 101, 4 (2021), 1487–1559.
- [7] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, et al. 2025. DeepSeek-V3 Technical Report. arXiv:2412.19437 [cs.CL] <https://arxiv.org/abs/2412.19437>
- [8] Aaron J Fried, Spencer D Dorn, William J Leland, Emily Mullen, Donna M Williams, Aimee K Zaas, Jack MacGuire, and Debra L Bynum. 2024. Large language models in internal medicine residency: current use and attitudes among internal medicine residents. *Discover Artificial Intelligence* 4, 1 (2024), 70.
- [9] Clark Glymour and Madelyn R Glymour. 2014. Commentary: race and sex are causes. *Epidemiology* 25, 4 (2014), 488–490.
- [10] Ruocheng Guo, Lu Cheng, Jundong Li, P. Richard Hahn, and Huan Liu. 2020. A Survey of Learning Causality with Data: Problems and Methods. *ACM Comput. Surv.* 53, 4, Article 75 (July 2020), 37 pages. doi:10.1145/3397269
- [11] David Heckerman, Dan Geiger, and David M Chickering. 1995. Learning Bayesian networks: The combination of knowledge and statistical data. *Machine learning* 20, 3 (1995), 197–243.
- [12] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The Curious Case of Neural Text Degeneration. arXiv:1904.09751 [cs.CL] <https://arxiv.org/abs/1904.09751>
- [13] Dominik Janzing and Bernhard Schölkopf. 2010. Causal inference using the algorithmic Markov condition. *IEEE Transactions on Information Theory* 56, 10 (2010), 5168–5194.
- [14] Subbarao Kambhampati. 2024. Can large language models reason and plan? *Annals of the New York Academy of Sciences* 1534, 1 (2024), 15–18.
- [15] Mert Karabacak and Konstantinos Margetis. 2023. Embracing large language models for medical applications: opportunities and challenges. *Cureus* 15, 5 (2023).
- [16] Neville Kenneth Kitson, Anthony C Constantinou, Zhigao Guo, Yang Liu, and Kiattikun Chobtham. 2023. A survey of Bayesian Network structure learning. *Artificial Intelligence Review* 56, 8 (2023), 8721–8814.
- [17] Isaac S Kohane. 2024. Compared with what? Measuring AI against the health care we have. *New England Journal of Medicine* 391, 17 (2024), 1564–1566.
- [18] Daphne Koller and Nir Friedman. 2009. *Probabilistic Graphical Models*. MIT press.
- [19] Benjamin Kuipers and Jerome P Kassirer. 1984. Causal reasoning in medicine: analysis of a protocol. *Cognitive Science* 8, 4 (1984), 363–385.
- [20] Wai Lam and Fahiem Bacchus. 1993. Using causal information and local measures to learn Bayesian networks. In *Uncertainty in Artificial Intelligence*. Elsevier, 243–250.
- [21] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems* 33 (2020), 9459–9474.
- [22] Fenglin Liu, Hongjian Zhou, Boyang Gu, Xinyu Zou, Jinfa Huang, Jing Wu, Yiru Li, Sam S Chen, Yining Hua, Peilin Zhou, et al. 2025. Application of large language models in medicine. *Nature Reviews Bioengineering* (2025), 1–20.
- [23] Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. BioGPT: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics* 23, 6 (2022).
- [24] Saurabh Mathur, Veerendra P Gadekar, Rashika Ramola, Peixin Wang, Ramachandran Thiruvengadam, David M Haas, Shinjini Bhatnagar, Nitya Wadhwa, Garbhini Study Group, Predrag Radivojac, Himanshu Sinha, Kristian Kersting, and Sriraam Natarajan. 2024. Modeling Multiple Adverse Pregnancy Outcomes: Learning from Diverse Data Sources. In *AIME*. Springer.
- [25] Saurabh Mathur, Ranveer Singh, Michael Skinner, Ethan Sanford, Neel Shah, Phillip Reeder, Lakshmi Raman, and Sriraam Natarajan. 2025. LLM-Guided Causal Bayesian Network Construction for Pediatric Patients on ECMO. In *International Conference on Artificial Intelligence in Medicine*. Springer, 255–260.
- [26] Liam G McCoy, Arjun K Manrai, and Adam Rodman. 2024. Large language models and the degradation of the medical record. *The New England journal of medicine* 391, 17 (2024), 1561–1564.
- [27] R Thomas McCoy, Shunyu Yao, Dan Friedman, Mathew D Hardy, and Thomas L Griffiths. 2024. Embers of autoregression show how large language models are shaped by the problem they are trained to solve. *Proceedings of the National Academy of Sciences* 121, 41 (2024), e2322420121.
- [28] Osman Mian, Michael Kamp, and Jilles Vreeken. 2023. Information-theoretic causal discovery and intervention detection over multiple environments. In *AAAI*, Vol. 37.
- [29] Frederic P. Miller, Agnes F. Vandome, and John McBrewhster. 2009. *Levenshtein Distance: Information theory, Computer science, String (computer science), String metric, Damerau-Levenshtein distance, Spell checker, Hamming distance*. Alpha Press.
- [30] 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* (2024).
- [31] Raymond J Mooney and Jude W Shavlik. 2021. A Recap of Early Work on Theory and Knowledge Refinement.
- [32] Sriraam Natarajan, Saurabh Mathur, Sahil Sidheekh, Wolfgang Stammer, and Kristian Kersting. 2025. Human-in-the-loop or AI-in-the-loop? Automate or Collaborate?. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 39. AAAI Press, 28594–28600.
- [33] OpenAI. 2025. ChatGPT: GPT-4o Language Model. <https://openai.com/chatgpt>.
- [34] Judea Pearl. 2009. *Causality*. Cambridge University Press.
- [35] Jonas Peters and Peter Bühlmann. 2015. Structural Intervention Distance for Evaluating Causal Graphs. *Neural Computation* 27, 3 (2015), 771–799.
- [36] Maxime Peyrard and Robert West. 2021. A Ladder of Causal Distances. In *Proceedings of International Joint Conference on Artificial Intelligence*.
- [37] Lawrence Rabiner and Biinghwang Juang. 1986. An introduction to hidden Markov models. *IEEE ASSP Magazine* 3, 1 (1986), 4–16.
- [38] Pedro Sanchez, Jeremy P Voisey, Tian Xia, Hannah I Watson, Alison Q O’Neil, and Sotirios A Tsiftaris. 2022. Causal machine learning for healthcare and precision medicine. *Royal Society Open Science* 9, 8 (2022), 220638.
- [39] Ilan S Schwartz, Katherine E Link, Roxana Daneshjoo, and Nicolás Cortés-Penfield. 2024. Black box warning: large language models and the future of infectious diseases consultation. *Clinical infectious diseases* 78, 4 (2024), 860–866.
- [40] Michael A Shwe, Blackford Middleton, David E Heckerman, Max Henrion, Eric J Horvitz, Harold P Lehmann, and Gregory F Cooper. 1991. Probabilistic diagnosis using a reformulation of the INTERNIST-1/QMR knowledge base. *Methods of information in Medicine* 30, 04 (1991), 241–255.
- [41] 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* 31, 3 (2025), 943–950.
- [42] Peter Spirtes, Clark N Glymour, and Richard Scheines. 2000. *Causation, prediction, and search*. MIT press.
- [43] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, et al. 2025. Gemini: A Family of Highly Capable Multimodal Models. arXiv:2312.11805 [cs.CL] <https://arxiv.org/abs/2312.11805>
- [44] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. <https://arxiv.org/abs/2302.13971>
- [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [46] Jonas Wahl and Jakob Runge. 2025. Separation-Based Distance Measures for Causal Graphs. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 3412–3420.
- [47] Razeghi Yasaman, Robert Logan IV, Gardner Matt, and S Sameer. 2022. Impact of pretraining term frequencies on few-shot numerical reasoning. *Findings of the Association for Computational Linguistics: EMNLP 2022* (2022), 840–854.
- [48] Alessio Zanga, Elif Ozkirimli, and Fabio Stella. 2022. A survey on causal discovery: theory and practice. *International Journal of Approximate Reasoning* 151 (2022), 101–129.
- [49] Matej Zečević, Moritz Willig, Devendra Singh Dhami, and Kristian Kersting. 2023. Causal Parrots: Large Language Models May Talk Causality But Are Not Causal. *Transactions on Machine Learning Research* (2023).