

FILLING THE GAPS: LLMs FOR CAUSAL HYPOTHESIS GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Scientific discovery catalyzes human intellectual advances, driven by the cycle of hypothesis generation, experimental design, data evaluation, and iterative assumption refinement. This process, while crucial, is expensive and heavily dependent on the domain knowledge of scientists to generate hypotheses and navigate the scientific cycle. Central to this is causality, the ability to establish the relationship between the cause and the effect. Motivated by the scientific discovery process, in this work, we formulate a novel task where the input is a partial causal graph with missing variables, and the output is a hypothesis about the missing variables to complete the partial graph. We design a benchmark with varying difficulty levels and knowledge assumptions about the causal graph. With the growing interest in using Large Language Models (LLMs) to assist in scientific discovery, we benchmark open-source and closed models on our testbed. We show the strong ability of LLMs to hypothesize the mediation variables between a cause and its effect. In contrast, they underperform in hypothesizing the cause and effect variables themselves. We also observe surprising results where some of the open-source models outperform the closed GPT-4 model.

1 INTRODUCTION

Scientific discovery has been key to humankind’s advances. It is a dynamic process revolving around inquiry and constant refinements driven by new observations. Scientists adhere to a structured process that involves formulating a hypothesis and then collecting pertinent data (Wang et al., 2023a). They then draw inferences from experiments and the collected data, modify the hypothesis, formulate sub-questions, and repeat the process until the research question is answered (Kiciman et al., 2023).

Causality empowers scientists to assess the hypotheses and interpret the collected data beyond mere correlations and associations. Tools such as Randomised Control Trials (RCTs) (Kendall, 2003) allow for establishing causal relationships between variables. Naturally, the process of causal discovery heavily relies on human experts to guide the hypothesis formation and experimental design (Kiciman et al., 2023). Expert domain knowledge is crucial to narrow the search space of hypotheses, especially when it is expensive to collect data or when systematic exploration is infeasible. However, a possible impediment is that domain knowledge can be difficult to formulate and collect (Kiciman et al., 2023).

With the recent advancement of Large Language Models (LLMs) (Brown et al., 2020; OpenAI, 2023), there has been a growing interest in using them for scientific discovery (AI4Science and Quantum, 2023; Lu et al., 2024; Cory-Wright et al., 2024). Their potential is now studied in domains such as natural sciences (AI4Science and Quantum, 2023). Given the importance of causality in the scientific discovery process, we focus on how LLMs can

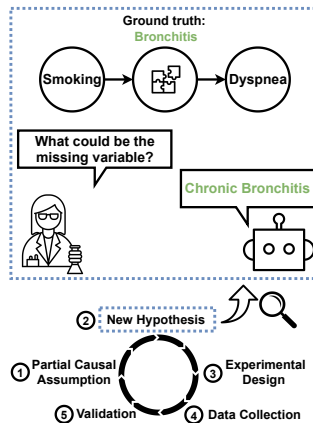


Figure 1: Scientific discovery iteratively generates hypotheses from assumptions using human expertise. We use LLMs as proxy experts to propose new hypotheses in causal DAGs.

054 assist with causal reasoning. LLMs have achieved state-of-the-art results for causal tasks such as
055 determining pairwise causal relationships by considering variable names (Kıcıman et al., 2023),
056 combined with causal discovery algorithms (Abdulaal et al., 2024; Ban et al., 2023a; Vashishtha et al.,
057 2023) for refinement.

058 Causal discovery, however, comes *after* hypothesizing the variables of interest (which require domain
059 knowledge), forming experiments, and potentially costly data collection. Our work, therefore,
060 extends LLM applications to assist in steps essential *before* causal discovery, specifically identifying
061 and hypothesizing *missing variables* in partially known causal graphs. This simulates the realistic
062 scientific discovery process of incremental hypothesis formation and testing. **By leveraging curated
063 causal graphs, we evaluate the feasibility and reliability of LLMs in generating hypotheses under
064 controlled yet realistic settings, ensuring reproducibility and providing a foundation for LLM-driven
065 scientific discovery.**

066 We break down causal hypothesis generation into smaller tasks, starting with baseline experiments,
067 and progressing to realistic scenarios where only treatments and outcomes are known. We leverage
068 LLMs’ large-scale training to propose memorized or inferred variables based on their general and
069 domain knowledge. This enables users to identify missing variables to guide data collection, fol-
070 lowed by subsequent downstream causal tasks. Importantly, we avoid requiring LLMs to determine
071 pairwise causal relations or perform numerical calculations, sidestepping their limitations in these
072 areas (Zečević et al., 2023; Jin et al., 2023a). Existing work explores the inductive hypothesis
073 generation capabilities of LLMs by using them as creative solution proposers with task-specific
074 means of verifying said solutions (Romera-Paredes et al., 2023; Wang et al., 2023b; Qiu et al., 2024).
075 In contrast, our work uniquely focuses on hypothesis generation within a causal paradigm.

076 **Contributions.** Our main contributions are: 1) We propose and formalize the novel task of LLM-
077 assisted causal variable identification and hypothesizing. 2) We propose a benchmark for hypothesiz-
078 ing missing variables across diverse domains of existing causal graphs. 3) We design experimental
079 tests with different difficulty levels and knowledge assumptions, such as open-world and closed-world
080 settings, the number of missing variables, etc. 4) **Our benchmark allows for allow for both grounded
081 evaluations and a reproducible framework to benchmark LLMs’ capabilities in hypothesis generation.**

082 2 RELATED WORK

083 **LLMs and Causality.** Our work is based on the framework of causality as proposed by Pearl
084 (2009). The intersection of language and causality is explored by Girju et al. (2002); Hassanzadeh
085 et al. (2020); Tan et al. (2023); Dhawan et al. (2024) to extract causal relationships from a large
086 corpus of text. With the advancements in LLMs and their ability to process large contexts, there
087 has been an interest in using them for causal reasoning (Kıcıman et al., 2023). Some works have
088 focused on commonsense causality (Frohberg and Binder, 2021; Singh et al., 2021) and temporal
089 causal reasoning (Zhang et al., 2020; 2022). More recently Kıcıman et al. (2023); Long et al. (2023);
090 Darvari et al. (2024) introduced methods to discover causal structures by prompting LLMs with
091 variable names. Ban et al. (2023b); Vashishtha et al. (2023); Ban et al. (2023a) extended this work by
092 introducing ancestral constraints to refine the causal structures derived from LLMs. Abdulaal et al.
093 (2024) combined data-based deep structural causal models, such as (Yu et al., 2019), with LLMs
094 generated causal structure. Jin et al. (2023b) focused on causal inference using LLMs. While a tool
095 on GitHub PyWhy-LLM used LLMs to propose confounders, our work formalizes such a task along
096 with detailed insights. Recent work attempted to train transformers for improved causal inference
097 and discovery (Vashishtha et al., 2024; Zhang et al., 2024). In contrast to previous work, we focus
098 on the novel task of identifying and hypothesizing missing variables, a task that comes before data
099 collection and evaluation, with LLMs as assistants. We test the hypothesizing abilities of generalist
100 pre-trained LLMs as our task is primarily linked with pre-training knowledge.

101 **LLMs and hypothesis generation.** Existing works tested inductive hypothesis generation with
102 LLMs in reasoning tasks or free-form scientific hypotheses from background knowledge provided in
103 the context (Gendron et al., 2023; Qi et al., 2023; Xu et al., 2023a;b; Qiu et al., 2024; Lu et al., 2024).
104 In contrast, we consider the structured task of causal hypothesis generation, where the ground-truth
105 variables are known and can be used for evaluation. We also assume a pertinent human-in-the-loop
106 assistive scientific discovery setup to counter LLMs’ limitations and hallucinations.

3 PRELIMINARIES: CAUSAL GRAPH

A causal relationship can be modeled via a Directed Acyclic Graph (DAG). A causal DAG represents relationships between a set of N variables defined by $\mathbf{V} = \{v_1, \dots, v_N\}$. The variables are encoded in a graph $\mathcal{G} = (\mathbf{V}, \mathbf{E})$ where \mathbf{E} is a set of directed edges between the nodes $\in \mathbf{V}$ such that no cycle is formed. Mathematically it can be expressed as:

$$\mathcal{G} = (\mathbf{V}, \mathbf{E}), \mathbf{E} = \{e_{i,j} \mid v_i, v_j \in \mathbf{V}, i \neq j\} \text{ and } v_i \rightarrow v_j$$

Each edge $e_{i,j} \in \mathbf{E}$ denotes causal relationship between v_i and v_j , $v_i \xrightarrow{e_{i,j}} v_j$, emphasizing the influence from v_i to v_j . Beyond visualization, causal DAGs allow for the mathematical characterization of different node types for a causal model to understand the influences and dependencies.

We define $d(v)$ as the degree of a node v , representing the total number of edges connected to v . $d_{\text{in}}(v)$ is the in-degree, representing the number of incoming edges to v . $d_{\text{out}}(v)$ is the out-degree, representing the number of outgoing edges from v .

Sources are variables v_s with no incoming edges. Mathematically sources are $d_{\text{in}}(v_s) = 0$ where d_{in} is the in-degree of the graph.

Sinks are variables v_k with no outgoing edges. Sinks are $d_{\text{out}}(v_k) = 0$ where d_{out} is the out-degree of the graph.

Treatment are variables v_t , characterized as nodes $d_{\text{in}}(v_t) = 0$ that are being intervened upon.

Outcome are variables v_y , characterized as the nodes $d_{\text{out}}(v_y) = 0$ that are observed for interventions from the treatments.

Mediator are variables v_m that have both incoming and outgoing edges ($d_{\text{in}}(v_m) > 0$ and $d_{\text{out}}(v_m) > 0$), acting as intermediaries in the causal pathways between treatment and outcome.

Confounder are variables v_k that influence both treatment and outcome, exhibiting edges directed towards the treatment and outcome nodes ($d_{\text{out}}(v_k) \geq 2$). Hence v_k is a confounder if it is a parent of both v_i and v_j .

Collider are variables v_l that have two edges meeting, and have an in-degree greater than one ($d_{\text{in}}(v_l) > 1$). Hence v_k is a collider if it is a child of both v_i and v_j .

Average Treatment Effect. Average Treatment Effect (ATE) quantifies the expected change in the outcome v_y caused by the unit change of the treatment v_t . ATE is a part of the causal do-calculus introduced by Pearl (2009). We consider binary causal DAGs, i.e., each variable can either take 0 or 1 as values.

$$\text{ATE} = \mathbb{E}[v_y | \text{do}(v_t = 1)] - \mathbb{E}[v_y | \text{do}(v_t = 0)]$$

where the $\text{do}(\cdot)$ operator, represents an intervention. The $E[v_y | \text{do}(v_t = 1)]$ represents the expected value of the outcome variable v_y when we intervene to set the treatment variable v_t to 1 (i.e., apply the treatment), and $E[v_y | \text{do}(v_t = 0)]$ represents the expected value of v_y when we set v_t to 0 (i.e., do not apply the treatment).

Mediation Analysis. Mediation analysis is implemented to quantify the effect of a treatment on the outcome via a third variable, the mediator. The total mediation effect can be decomposed into the Natural Direct Effect (NDE) and the Natural Indirect Effect (NIE). The Natural Direct Effect (NDE) is the effect of the treatment on the outcome variable when not mediated by the mediator variable. The Natural Indirect Effect (NIE) is the effect of the treatment variable on the outcome variable when mediated by the mediator variable.

$$\text{NDE} = \mathbb{E}[v_{t=1}, v_{m=0} - v_{t=0}, v_{m=0}]$$

Here, NDE is calculated by comparing the expected outcome when the treatment variable is set to 1 and the mediator is fixed at the level it would take under the control treatment $v_t = 0$, with the expected outcome when both the treatment and the mediator are set to the control level.

$$\text{NIE} = \mathbb{E}[v_{t=0}, v_{m=1} - v_{t=0}, v_{m=0}]$$

Here, NIE is calculated by comparing the expected outcome when the treatment variable is set to 1 and the mediator is allowed to change as it would under the treatment, with the expected outcome when the treatment variable is set to 1 but the mediator is fixed at the control level.

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

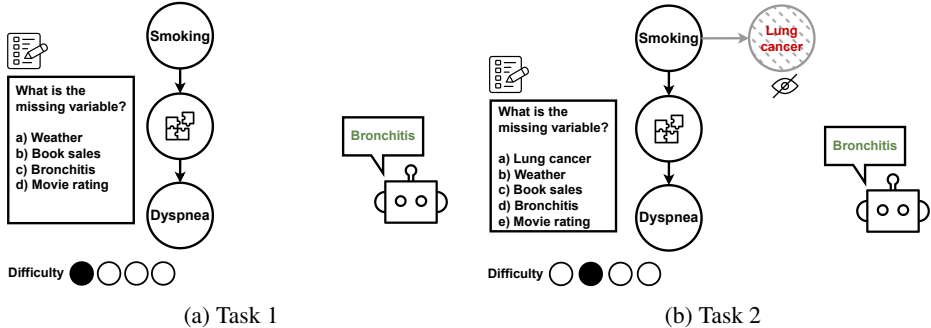


Figure 2: Leveraging LLM to identify the missing variable for a causal DAG in the presence of out-of-context distractors (a), an in-context distractor along with out-of-context distractors (b).

4 LLMs FOR IDENTIFYING AND HYPOTHESIZING CAUSAL VARIABLES

In this work, we aim to leverage language models to identify and hypothesize variables in a causal DAG. Motivated by the process of hypothesizing a causal graph from a partially known structure Glymour et al. (2019), this paper proceeds under the assumption that some elements of the graph are already known. The aim is to find additional variables that can be incorporated into the existing causal structure to enhance the underlying causal mechanism.

We assume a partially known causal DAG, defined as $\mathcal{G}^* = (V^*, E)$, where $V^* \subseteq V$. The objective is to identify the set of missing variables $V^* = V \setminus V_{\text{missing}}$ thereby expanding \mathcal{G}^* to \mathcal{G} . This implies that all causal relationships (edges) among variables in V^* are known and correctly represented in \mathcal{G}^* ; i.e., E is fully specified. **Here, “missing” variables are not latent or hidden by measurement error but known unknowns within the causal graph reflective of LLMs perspective.**

Our methodology is structured around progressively challenging scenarios, explores the ability of LLMs to identify and hypothesize causal variables. This starts from a restrictive and controlled exploration to an open-ended one. Initially, we restrict the exploration by providing the language models with a partially known causal DAG and a set of multiple choices for the missing variables. The complexity of the task is gradually increased by removing more than one node from the graph. Finally, we move to an open-ended scenario where the ground truth is not available to LLM. In this setting, LLM is required to hypothesize the missing variables of the causal DAG without any explicit hints. We evaluate the causal reasoning capability of LLMs through prompting. Given LLMs’ limitation to textual input, we represent the graph \mathcal{G}^* using a prompt template $P_{\text{LLM}}(\cdot)$ which enables LLMs to parse the causal relationships embedded within the DAG.

4.1 TASK 1: OUT-OF-CONTEXT CONTROLLED VARIABLE IDENTIFICATION

This task (depicted in Figure 2a) evaluates LLMs’ ability to identify missing variables in a causal graph from a list of multiple choices, thereby reconstructing the original graph. The partial DAG \mathcal{G}^* is created by removing one variable from the original DAG \mathcal{G} . Let us denote the removed node as v_x . Along with the partial graphs, we operate in the multiple-choice question answering (MCQA) paradigm. The role of the LLM is to select a variable from the multiple choices, MCQ_{v_x} , that can be used to complete the graph. The multiple choices include the missing variable v_x and out-of-context distractors. **The out-of-context distractors are unrelated to the causal domain of the given DAG, chosen to minimize any contextual and overlap with the true missing variable.** Let v_x^* represent the variable selected by the LLM to complete \mathcal{G}^* .

$$v_x^* = P_{\text{LLM}}(\mathcal{G}^*, \text{MCQ}_{v_x}) \quad \forall v_x \in V$$

4.2 TASK 2: IN-CONTEXT CONTROLLED VARIABLE IDENTIFICATION

In practical applications, such as healthcare (Robins, 1986) and finance (Hughes et al., 2019), dealing with missing data and unobserved latent variables is a major challenge (Tian and Pearl, 2012; Bentler, 1980). Therefore, it is important to identify the missing variables and their underlying causal

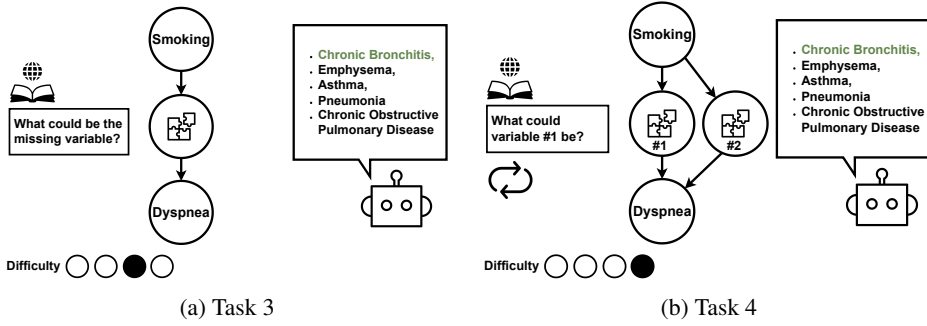


Figure 3: Leveraging LLM to hypothesize the missing variable in a causal DAG in an open-world setting for one variable (a), in an iterative fashion for multiple missing mediators (b).

mechanism. To simulate this, a more challenging task is introduced (see Figure 2b). Here, instead of removing one node from the ground truth DAG \mathcal{G} , two nodes, v_{x_1} and v_{x_2} , are now removed to create the partial graph, \mathcal{G}^* .

$$\mathcal{G}^* = \mathcal{G} \setminus \{v_{x_1}, v_{x_2}\} \quad \text{for } v_{x_1}, v_{x_2} \in \mathbf{V}$$

We use the MCQA paradigm to provide multiple choices that include the missing variables v_{x_1} and v_{x_2} . The task for the LLM here is to select the correct variable v_{x_1} only, given an in-context choice v_{x_2} and out-of-context choices. **The in-context variables are plausible variables within the same causal graph.** We introduce the non-parental constraint for v_{x_1} and v_{x_2} . This prevents the removal of both a parent node and its immediate child node in \mathcal{G}^* .

$$v_{x_1}^* = P_{\text{LLM}}(\mathcal{G}^*, \text{MCQ}_{v_{x_1}, v_{x_2}}) \quad \forall v_{x_1}, v_{x_2} \in \mathbf{V} \text{ and } v_{x_1} \not\rightarrow v_{x_2}, v_{x_2} \not\rightarrow v_{x_1}$$

4.3 TASK 3: HYPOTHESIZING IN OPEN WORLD

So far, we have described the testbeds for variable identification in a partial DAG given the controlled world knowledge in the form of distractors. This assumption allows for the evaluation of the language model’s ability to select the correct answer from a set of options. However, in the open-world setting, we increase the complexity to provide no choices, as shown in Figure 3a. Hence the task is to predict the missing node v_x given the partial graph \mathcal{G}^* to complete the ground truth graph \mathcal{G} . Here, the model returns a set of potential hypotheses, $\{v_{x,1}^*, \dots, v_{x,k}^*\}$ where k is the number of hypotheses.

$$\{v_{x,1}^*, v_{x,2}^*, \dots, v_{x,k}^*\} = P_{\text{LLM}}(\mathcal{G}^*) \quad \forall v_x \in \mathbf{V}$$

4.4 TASK 4: ITERATIVELY HYPOTHESIZING IN OPEN WORLD

In addition to the search space relaxation, we further relax the number of missing variables. The partial DAG here, is obtained for one or more missing node variables. $\mathcal{G}^* = \mathcal{G} \setminus \{v_{x_1} \dots v_{x_M}\}$. The fine-grained results from the open-world setting reveal that language models exhibit a particularly strong performance in identifying mediator variables. Thus, the LLM is used here to iteratively hypothesize mediator variables in a causal DAG given a treatment and an effect. The task (shown in Figure 3b) is set up as follows: given a partial graph \mathcal{G}^* , which includes observed treatment and outcome variables, we aim to hypothesize a set of mediators, denoted as $M = \{v_{m_1}, v_{m_2}, \dots, v_{m_H}\}$, that mediates the treatment v_t to the outcome v_y . Here, H represents the number of direct, and indirect mediators. A pair of treatments and outcomes are considered iteratively across the causal DAG. In the first iteration, the LLM generates a hypothesis for the mediator v_{m_1} . The hypothesized mediator, v_{m_1} is then added to the graph, updating $\mathcal{G}^* \rightarrow \mathcal{G}^* \cup \{v_{m_1}\}$. The partial graph that now also includes $v_{m_1}^*$ can be used to identify the second mediator $v_{m_2}^*$ and so on. Therefore, in each subsequent iteration i , the LLM is tasked to generate a hypothesis for the next missing mediator v_{m_i} given the updated graph $\mathcal{G}^* \cup \{v_{m_1}^*, \dots, v_{m_{i-1}}^*\}$.

$$v_{m_i}^* = P_{\text{LLM}}(\mathcal{G}^* \cup \{v_{m_1}^*, \dots, v_{m_{i-1}}^*\}) \text{ for } i = 1, \dots, H$$

The sequence of mediators $M = \{v_{m_1}, v_{m_2}, \dots, v_{m_H}\}$ is chosen at random. To formally investigate how the order of hypothesized mediators influences LLM performance, we borrow concepts from the

mediation analysis literature, specifically the Natural Direct Effect (NDE) and the Natural Indirect Effect (NIE). The NDE measures the effect of the treatment on the outcome that is not mediated by a particular mediator, while the NIE measures the effect of the treatment that is mediated by the mediator. We introduce a metric called Mediation Influence Score (MIS) that quantifies the influence of each mediator between a treatment and an effect. MIS defined as the ratio of NIE to NDE, provides a scale-free measure of a mediator’s relative influence, enabling prioritization. MIS is always positive, reflecting the absolute contribution of mediators.

$$\text{MIS}(v_{m_i}) = \left| \frac{\text{NIE}(v_{m_i})}{\text{NDE}(v_{m_i})} \right| \quad \text{for } i = 1, \dots, H.$$

This metric quantifies the relative importance of the indirect effect (through the mediator) compared to the direct impact. Mediators are then ranked and prioritized based on their MIS scores, with higher scores indicating a stronger mediation effect.

5 EVALUATION AND RESULTS

5.1 EXPERIMENTAL SETUP

We evaluate a variety of causal datasets spanning diverse domains. We use the semi-synthetic datasets from BNLearn repository - Cancer: $\mathcal{G}(5, 4)$ (Korb and Nicholson, 2010), Survey: $\mathcal{G}(6, 6)$ (Scutari and Denis, 2021), Asia: $\mathcal{G}(8, 8)$ (Lauritzen and Spiegelhalter, 1988), Child: $\mathcal{G}(20, 25)$ (Spiegelhalter, 1992), Insurance: $\mathcal{G}(27, 52)$ (Binder et al., 1997), and Alarm: $\mathcal{G}(37, 46)$ (Beinlich et al., 1989). We also evaluate our approach on a realistic Alzheimer’s Disease dataset: $\mathcal{G}(9, 16)$ (Abdulaal et al., 2024), developed by five domain experts. We also test on a legal causal graph, Law: $\mathcal{G}(8, 20)$ (VanderWeele and Staudt, 2011). See Appendix A.1 for further details.

We evaluate our setups across different open-source and closed models. The models we use are GPT-3.5 (Brown et al., 2020), GPT-4 (OpenAI, 2023), LLama2-chat-7b (Touvron et al., 2023), Mistral-7B-Instruct-v0.2 (Jiang et al., 2023), Mixtral-7B-Instruct-v0.1 (Jiang et al., 2024), Zephyr-7b-Beta (Tunstall et al., 2023) and Neural-chat-7b-v3-1 (Intel, 2023).

Implementation details are mentioned in Appendix A. Prompt templates are illustrated in Appendix F.

5.2 TASK 1: OUT-OF-CONTEXT CONTROLLED VARIABLE IDENTIFICATION

Our first experiment is designed to assess the fundamental ability of language models to identify missing variables in a partial causal graph, serving as a baseline for understanding their performance in variable identification tasks. Here, the input to the LLM is the ground truth variable name in addition to out-of-context multiple choices for the missing variable v_x and the partial DAG \mathcal{G}^* . We then calculate the models’ accuracy in correctly predicting v_x .

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(v_x^* = v_x^i)$$

Results. In Figure 4a, we report the accuracy of different LLMs in identifying the missing variable. GPT-4, followed by Mixtral, consistently performs well, achieving perfect accuracy on most of the datasets. GPT-3.5 also shows overall strong performance, apart from the Insurance and Alarm datasets. The other models, including Mistral-7b, Llama-7b, and Zephyr-7b, demonstrate varying degrees of success. Insurance is the most challenging dataset, which could potentially be due to the high number of edges present in the DAG. It is noteworthy that all models significantly outperform the random baseline, indicating that given out-of-context multiple choices along with the ground truth variable, the language model can pick out the missing causal variable to complete the partial graph \mathcal{G}^* . However, we may conjecture that the high performance could be attributed to the simplicity of the task. The models might be primarily inferring from the context of the dataset domain, rather than performing actual causal reasoning among multiple plausible choices. To further investigate this, we introduce an in-domain choice in the multiple choices in the next experiment. This can assess LLMs’ ability to choose a causal variable for a partial DAG beyond the highly evident correlations.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

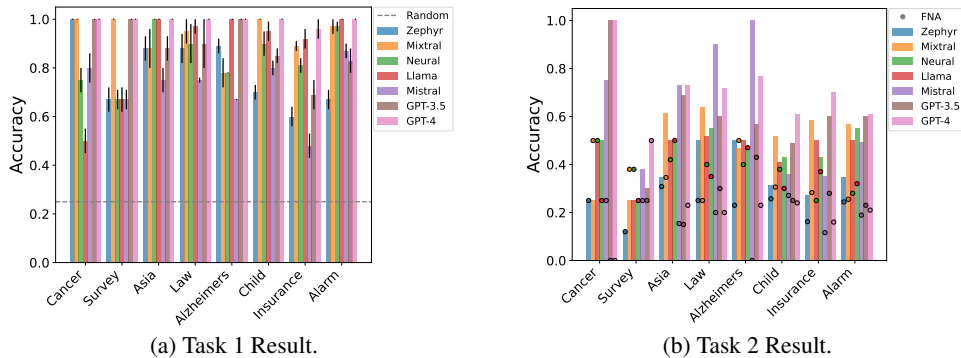


Figure 4: Accuracy of LLMs in identifying the missing causal variable from multiple choices with out-of-context distractors (a), and from both out-of-context and in-context distractors (b).

	Cancer		Survey		Asia		Law		Alzheimers		Child		Insurance		Alarm		Avg	
	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J
Zephyr	0.36	0.61	0.34	0.60	0.45	0.66	0.41	0.70	0.35	0.75	0.51	0.70	0.45	0.44	0.46	0.69	0.42	0.63
Mixtral	0.41	0.66	0.39	0.66	0.66	0.75	0.38	0.69	0.31	0.77	0.53	0.77	0.46	0.56	0.50	0.72	0.46	0.70
Neural	0.38	0.77	0.43	0.55	0.53	0.55	0.47	0.72	0.44	0.71	0.48	0.70	0.47	0.43	0.47	0.67	0.45	0.63
Llama	0.40	0.48	0.40	0.54	0.53	0.58	0.67	0.65	0.45	0.61	0.48	0.63	0.42	0.34	0.46	0.65	0.45	0.55
Mistral	0.33	0.67	0.44	0.65	0.60	0.73	0.49	0.67	0.34	0.76	0.48	0.68	0.46	0.47	0.47	0.71	0.44	0.67
GPT-3.5	0.48	0.74	0.42	0.79	0.47	0.61	0.52	0.73	0.39	1.00	0.36	0.60	0.47	0.52	0.48	0.73	0.44	0.71
GPT-4	0.49	0.90	0.51	0.67	0.66	0.76	0.55	0.78	0.47	0.98	0.36	0.53	0.52	0.56	0.49	0.75	0.50	0.73

Table 1: Task 3 Results. Average semantic similarity and LLM-as-Judge metrics to evaluate LLMs in hypothesizing the missing variable in a causal DAG.

5.3 TASK 2: IN-CONTEXT CONTROLLED VARIABLE IDENTIFICATION

We introduce a more complex setting to further challenge the models’ abilities in **missing variable identification**. In this setup, the partial graph has two missing nodes. Alongside out-of-context choices and the ground truth variable, the multiple-choice options also include the second missing node from the partial graph as an in-context distractor. This configuration requires the language model to reason about indirect causal relationships to identify the correct missing variable. To evaluate models’ performance, we present two metrics: accuracy and false node accuracy. The false node accuracy, measures the confusion of LLMs in picking the in-context variable instead of the ground truth.

$$\text{False Node Accuracy (FNA)} \downarrow = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(v_{x_1}^* = v_{x_2})$$

Results. In Figure 4b, we plot both Accuracy and False Node Accuracy across different datasets. Ideally, accuracy should be 1.0, and the FNA should be 0.0. Since there were 5 multiple choices, the random chance is 0.2. We observe that most of the models for larger datasets achieve much higher accuracy than random chance. GPT-3.5 and GPT-4 consistently perform well across all datasets, with high accuracy and low FNA. This suggests that these models are capable of reasoning by identifying the missing nodes in the causal graph and are less likely to be confused by the in-context node variable. On the other hand, open-source models like Mistral, Zephyr, and Mixtral show varying performance across different datasets. For instance, Mistral performs well on the easy Cancer dataset but underperforms in the more complex Alarm dataset. In summary, we observe that most language models can identify causal variables in the presence of multiple missing nodes and an in-context distractor. These results indicate that while most language models can handle **missing variable identification** in the presence of multiple missing nodes and in-context distractors, the robustness of their reasoning abilities varies significantly with dataset complexity and model architecture.

5.4 TASK 3: HYPOTHESIZING IN OPEN WORLD

We recall that the goal is for the language models to be able to complete a causal graph given a partial graph. In realistic scenarios, where scientists provide incomplete graphs without pre-defined answers,

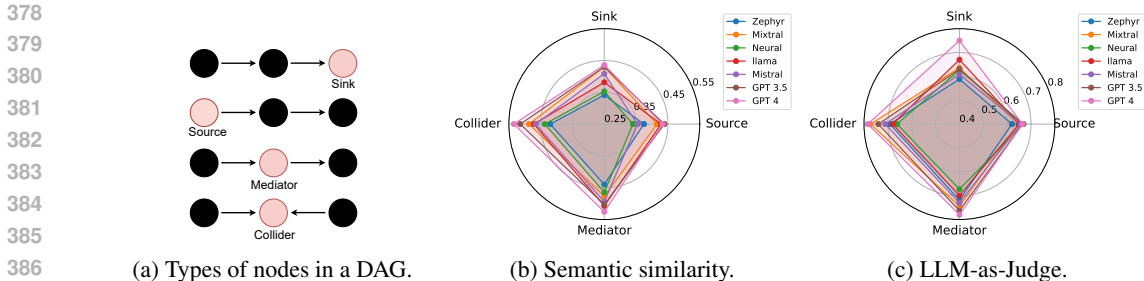


Figure 5: Task 3 Results. Visualizing each model’s performances, averaged across the different datasets, for Sink, Source, Mediator, and Collider nodes.

there is often no single ‘ground truth’ for what the missing variable should be. The correct hypothesis may vary based on domain expertise or available data, making this task fundamentally open-ended. Hence in this test-bed, we aim to leverage LLMs to hypothesize the causal variables. The language model is prompted for $k = 5$ suggestions for the missing node v_x .

We then compare the suggestions against the ground truth, acknowledging that in realistic scenarios, there is often no single ‘ground truth’ for the missing variables. This complexity necessitates careful evaluation, as we suspect that traditional metrics may not fully capture the performance of models, particularly when suggestions must be assessed within the broader context of the entire causal graph (see Appendix C.5 for more details). Hence, for a robust evaluation of this experiment, we use two metrics, semantic similarity, and LLM-as-Judge that incorporate contextual information.

Semantic Similarity: measures the cosine similarity between the embeddings (of another pretrained sentence embedding model) of each suggestion of the model’s predictions, $v_{x_{1:5}}^*$ and the ground truth v_x . The distances of the most similar suggestion are averaged across all nodes $v_x \in \mathbf{V}$. For a detailed explanation of this process, please refer to Appendix A.4.

LLM-Judge: This metric evaluates the quality of the model’s suggestions using a two-step process inspired by Zheng et al. (2023). In particular, LLM-as-Judge compares against ground truth variables to measure contextual semantic similarity beyond semi-exact matching like in semantic similarity metric. In the first step, the language model is prompted to determine which suggestion best fits the partial graph, given the ground truth and the suggestions, $v_{x_{1:5}}^*$. In the second step, the language model is again prompted to rate the selected suggestion on a scale of 1 to 10 in terms of similarity. This is repeated for all nodes, and the ratings are averaged to provide an overall quality measure. Implementation details can be found in Appendix A.5.

Results. We report models’ performances using both semantic similarity and LLM-Judge metrics in Table 1. For brevity, we provided the variances in Appendix C.1. To develop an intuition of LLMs’ performance, we provide a detailed analysis of each metric across different types of node variables (defined in Section 3). We specifically look at sources, sinks, colliders, and mediators for each of the partial causal graphs. The results, fine-grained by node type, are given in Figure 5 that shows each model’s average performance across datasets. For a detailed performance per each dataset individually, see Figure 18.

GPT-4 and Mistral generally achieve higher semantic similarity and LLM-as-Judge scores across most datasets (Figure 18). GPT-3.5 also shows good average performance. We observe that semantic similarity is a stricter metric than LLM-as-judge since it cannot encode contextual information about the causal DAG (see example in Table 7). Despite different scales, semantic similarity and LLM-as-judge metrics both seem to be fairly correlated. In Figure 5, we observe that models display stronger performance for colliders and mediators on average. This suggests that these models are relatively proficient at reasoning about common causes and indirect causal relationships. Sinks are typically the nodes that represent the outcomes or effects of interventions (treatments) applied to other nodes, and the lower performance on these nodes indicates that the models find it challenging to reason about the potential outcomes of the causal graphs. Source nodes represent the causes in a causal graph, and lower performance on these nodes might indicate difficulties in reasoning about the potential treatments from the partial graph.

In Figure 16a, we observe that the model performance increases with k , i.e., with a higher number of suggestions. From Figure 16b, it is also evident that the performance is proportional to the number of total edges, $d_{in} + d_{out}$ (more context about the node). In summary, LLMs show impressive

performance across some of the nodes and can be particularly useful to hypothesize mediators and colliders in a partial causal DAG. It is, hence, potentially beneficial to use LLMs in the real world because, in practice, treatment and outcomes are usually known.

5.4.1 HYPOTHESIZING CONFOUNDER

In causal inference, backdoor paths are alternative causal pathways that confound the estimation of causal effects. They introduce bias when estimating causal effects if not appropriately addressed. Hence hypothesizing and controlling for confounders is an important task in causal inference (Pourhoseingholi et al., 2012). We extract confounder subgraphs from (Sachs et al., 2005), Alarm, and Insurance graphs. From Table 2, with detailed results in Appendix B, we observe that while some confounders were easily hypothesized by LLMs, achieving perfect accuracy, the genomic domain of the SACHS posed challenges for models with potentially less domain-specific knowledge. Similar to the mediator analysis, a large model: GPT-4, does not always perform best across all datasets. This highlights the need for a diverse set of benchmarks, like ours, to fully assess models’ performance. Considering the importance of backdoor paths, we have also benchmarked LLM performance for confounders in addition to colliders. LLM typically performs well when hypothesizing a collider, however, the results for confounders are varied.

	Sachs	Alarm	Ins
Zephyr	0.10 ±0.01	0.45 ±0.05	0.53 ±0.06
Mixtral	0.95 ±0.10	0.85 ±0.09	0.63 ±0.07
Neural	0.30 ±0.03	0.45 ±0.05	0.61 ±0.06
LLama	0.20 ±0.02	0.47 ±0.05	0.63 ±0.06
Mistral	0.20 ±0.02	0.85 ±0.09	0.61 ±0.06
GPT-3.5	0.40 ±0.04	0.49 ±0.05	0.67 ±0.07
GPT-4	0.95 ±0.10	0.73 ±0.07	0.78 ±0.08

Table 2: Hypothesizing Confounders.

5.5 TASK 4: ITERATIVELY HYPOTHESIZING IN OPEN WORLD

In the previous open-world experiment, we observed that LLMs excel at identifying mediators when the treatments and outcomes are given. This observation could be particularly relevant in medical settings, where understanding the mediators can provide insights into causal mechanisms through which a treatment affects a patient’s outcome.

For hypothesizing mediators, we adopted an iterative approach rather than a global (all-at-once) strategy. This interactive process allows the language model to progressively refine its predictions, reducing the search space for subsequent variables. As observed in our empirical results (see Appendix C.7), LLMs underperform when tasked with making multiple simultaneous predictions across different mediators. The iterative approach aligns more closely with human reasoning, as evidenced by Chain-of-Thought (CoT) (Wei et al., 2022) strategies, where sequential decision-making enhances accuracy.

For **unordered mediator evaluation**, the model is prompted iteratively with mediators presented in random order, and the final semantic similarity is averaged across all predictions. In contrast, **ordered mediator evaluation** ranks the mediators using the Mediation Influence Score (MIS), prompting the model in both ascending and descending orders of significance. **We introduce the metric Δ , presenting the difference in performance when mediators are iteratively presented to the LLM in ascending and descending orders of significance defined by the MIS.** Given that some datasets only contain a single mediator, we selected the Asia, Child, Insurance, and Alarm datasets, as they offer a wider range of mediators, ranging from 1 to 10 for the Alarm dataset.

Results. The results of this experiment are in Table 3. Results with variances are provided in Appendix C.1. In a highly complex environment with more than one node missing and with open-world search space, we observe that LLMs can still maintain their performance. Unlike the overall consistent performance of GPT-4 across all of the datasets from the open-world setting, the model showed superior performance in Insurance and Alarm datasets only. As the complexity of the dataset increases, we observe larger differences in hypothesizing the mediators according to the MIS order. Positive Δ values suggest that prompting the LLM based on the MIS metric leads to higher semantic similarity between the mediator hypotheses and the ground truth variables. In summary, we observe that LLMs can be highly effective in iteratively hypothesizing multiple mediators in a DAG, and if present, some domain knowledge about the significance of the mediator can boost the performance.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

	Asia		Child		Insurance		Alarm	
	Sim	Δ	Sim	Δ	Sim	Δ	Sim	Δ
Zephyr	0.61	-0.02	0.54	0.17	0.47	0.19	0.51	0.20
Mixtral	0.87	0.01	0.50	0.18	0.48	0.15	0.52	0.13
Neural	0.65	0.04	0.48	0.21	0.42	0.16	0.46	0.12
Llama	0.80	0.07	0.49	-0.05	0.44	0.21	0.51	0.07
Mistral	0.33	0.02	0.50	0.12	0.48	0.13	0.47	0.11
GPT-3.5	0.48	0.01	0.36	0.25	0.48	0.17	0.51	0.02
GPT-4	0.49	0.04	0.39	0.16	0.52	0.14	0.60	-0.07

Table 3: Sim: semantic similarity for iteratively hypothesizing the mediator nodes (Task 4) when prompted with random order. Δ measures the change in the prediction of each model when repeating the experiment with ordering according to the MIS metric instead of randomly.

5.6 DISCUSSION AND LIMITATIONS

The results show that LLMs effectively hypothesize missing variables, particularly mediators, though performance varies with task complexity. Simple tasks, like identifying missing variables from controlled options, had high success rates. Without unified metrics, we focused on relative rankings across models (Appendix C.2) and observed that no model, including GPT-4, consistently outperformed the others. We hypothesize that the differences in performance across domains may stem from potential biases within the LLMs. These biases may stem from the models’ training data and, therefore, its parametric memory, leading to disparities in how effectively the models handle different tasks introduced in the benchmark. For instance, the models’ ability to hypothesize confounders varied significantly across datasets. In some cases, such as the Sachs dataset (see Appendix B), domain-specific knowledge gaps may have led to lower accuracy.

While this paper aimed to evaluate the ability of current LLMs to identify and hypothesize variables in a partial causal graph, we attempted to improve the performance by fine-tuning the model and few-shot prompting. However, given the limited size of the DAGs used, the resulting datasets were small, leading to mixed results (see Appendix D.1). We suspect that while fine-tuning may help the model to specialize, it can also reduce its ability to leverage the general parametric knowledge (Yang et al., 2024). Future approaches can look at domain-specific fine-tuning.

Given the non-disclosed datasets of models, it is difficult to confirm with absolute certainty that the datasets are not ingested by models during training. However, one of the datasets we used was released recently (Abdulaal et al., 2024) after the announced cut-off date of models. Additionally, our task itself is novel, including the way we verbalize the graphs and prompt the models. **Additionally, in Table 3 we further demonstrate that LLM-generated suggestions are non-verbatim, indicating they generate novel hypotheses rather than retrieving memorized patterns. Finally, we did not observe any direct reconstruction of graphs that would suggest memorization.**

Our setup assumes known edges among missing variables to enable controlled evaluation, which future work can extend. We envision this as a human-LLM collaboration under expert supervision, as LLMs cannot automatically identify the most plausible answer or express confidence in their responses (Zhou et al., 2024). Future work could explore better filtering mechanisms and improve performance on source and sink nodes.

6 CONCLUSION

Most causality literature assumes that the necessary data has been collected, focusing on establishing causal relationships between variables. However, generating hypotheses about which variables to observe is typically done by human experts. LLMs, trained on large-scale datasets, can act as expert proxies for this task. We introduce the novel task of using LLMs to hypothesize missing variables in causal graphs, formalizing it with benchmarks that vary in difficulty and knowledge of the ground truth graph. We evaluate models on identifying missing variables from in-context and out-of-context distractors and hypothesizing variables in an open-world setting. We also explore an iterative approach for populating graphs with up to 10 missing mediator nodes. Our results show that LLMs are particularly effective at hypothesizing mediators, which are often less known than treatments and outcomes. To support further research, we will release our benchmark and codebase.

REFERENCES

- 540
541
542 Ahmed Abdulaal, adamos hadjivasiliou, Nina Montana-Brown, Tiantian He, Ayodeji Ijishakin, Ivana
543 Drobnjak, Daniel C. Castro, and Daniel C. Alexander. Causal modelling agents: Causal graph
544 discovery through synergising metadata- and data-driven reasoning. In *ICLR*, 2024.
- 545 Microsoft Research AI4Science and Microsoft Azure Quantum. The impact of large language models
546 on scientific discovery: a preliminary study using gpt-4. *arXiv*, 2023.
- 547 Taiyu Ban, Lyuzhou Chen, Derui Lyu, Xiangyu Wang, and Huanhuan Chen. Causal structure learning
548 supervised by large language model. *arXiv*, 2023a.
- 549
550 Taiyu Ban, Lyvzhou Chen, Xiangyu Wang, and Huanhuan Chen. From query tools to causal architects:
551 Harnessing large language models for advanced causal discovery from data. *arXiv*, 2023b.
- 552 Ingo A Beinlich, Henri Jacques Suermondt, R Martin Chavez, and Gregory F Cooper. The alarm
553 monitoring system: A case study with two probabilistic inference techniques for belief networks.
554 In *AIME 89: Second European Conference on Artificial Intelligence in Medicine, London, August*
555 *29th–31st 1989. Proceedings*, pages 247–256. Springer, 1989.
- 556 Peter M Bentler. Multivariate analysis with latent variables: Causal modeling. *Annual review of*
557 *psychology*, 31(1):419–456, 1980.
- 558
559 John Binder, Daphne Koller, Stuart Russell, and Keiji Kanazawa. Adaptive probabilistic networks
560 with hidden variables. *Machine Learning*, 29:213–244, 1997.
- 561 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
562 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
563 few-shot learners. In *NeurIPS*, 2020.
- 564
565 Ryan Cory-Wright, Cristina Cornelio, Sanjeeb Dash, Bachir El Khadir, and Lior Horesh. Evolv-
566 ing scientific discovery by unifying data and background knowledge with ai hilbert. *Nature*
567 *Communications*, 15(1):5922, 2024.
- 568 Victor-Alexandru Darvari, Stephen Hailes, and Mirco Musolesi. Large language models are effective
569 priors for causal graph discovery. *arXiv*, 2024.
- 570
571 Nikita Dhawan, Leonardo Cotta, Karen Ullrich, Rahul G Krishnan, and Chris J Maddison. End-to-end
572 causal effect estimation from unstructured natural language data. *arXiv*, 2024.
- 573 Jörg Frohberg and Frank Binder. Crass: A novel data set and benchmark to test counterfactual
574 reasoning of large language models. *arXiv*, 2021.
- 575
576 Gaël Gendron, Qiming Bao, Michael Witbrock, and Gillian Dobbie. Large language models are not
577 abstract reasoners. *arXiv*, 2023.
- 578 Roxana Girju, Dan I Moldovan, et al. Text mining for causal relations. In *FLAIRS conference*, pages
579 360–364, 2002.
- 580
581 Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on
582 graphical models. *Frontiers in genetics*, 10:524, 2019.
- 583 Akash Gupta, Ivaxi Sheth, Vyas Raina, Mark Gales, and Mario Fritz. Llm task interference: An initial
584 study on the impact of task-switch in conversational history. *arXiv preprint arXiv:2402.18216*,
585 2024.
- 586
587 Oktie Hassanzadeh, Debarun Bhattacharjya, Mark Feblowitz, Kavitha Srinivas, Michael Perrone,
588 Shirin Sohrabi, and Michael Katz. Causal knowledge extraction through large-scale text mining.
589 In *AAAI Conference on Artificial Intelligence*, volume 34, pages 13610–13611, 2020.
- 590 Rachael A Hughes, Jon Heron, Jonathan AC Sterne, and Kate Tilling. Accounting for missing data
591 in statistical analyses: multiple imputation is not always the answer. *International journal of*
592 *epidemiology*, 48(4):1294–1304, 2019.
- 593 Intel. Intel neural-chat-7b model achieves top ranking on llm leaderboard! 2023.

- 594 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
595 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
596 Mistral 7b. *arXiv*, 2023.
- 597
598 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
599 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.
600 Mixtral of experts. *arXiv*, 2024.
- 601
602 Zhijing Jin, Yuen Chen, Felix Leeb, Luigi Gresele, Ojasv Kamal, Zhiheng Lyu, Kevin Blin, Fernando
603 Gonzalez, Max Kleiman-Weiner, Mrinmaya Sachan, and Bernhard Schölkopf. CLadder: Assessing
604 causal reasoning in language models. In *NeurIPS*, 2023a.
- 605
606 Zhijing Jin, Jiarui Liu, Zhiheng Lyu, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona Diab,
607 and Bernhard Schölkopf. Can large language models infer causation from correlation? *arXiv*,
608 2023b.
- 609
610 John Kendall. Designing a research project: randomised controlled trials and their principles.
611 *Emergency medicine journal: EMJ*, 20(2):164, 2003.
- 612
613 Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language
614 models: Opening a new frontier for causality. *arXiv*, 2023.
- 615
616 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
617 language models are zero-shot reasoners. *NeurIPS*, 2022.
- 618
619 Kevin B Korb and Ann E Nicholson. *Bayesian artificial intelligence*. CRC press, 2010.
- 620
621 Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei Wang, Liangyou Li, Lifeng Shang, Xin Jiang,
622 Qun Liu, and Kam-Fai Wong. Mt-eval: A multi-turn capabilities evaluation benchmark for large
623 language models. *EMNLP*, 2024.
- 624
625 Steffen L Lauritzen and David J Spiegelhalter. Local computations with probabilities on graphical
626 structures and their application to expert systems. *Journal of the Royal Statistical Society: Series
627 B (Methodological)*, 50(2):157–194, 1988.
- 628
629 Stephanie Long, Tibor Schuster, Alexandre Piché, ServiceNow Research, et al. Can large language
630 models build causal graphs? *arXiv*, 2023.
- 631
632 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist:
633 Towards fully automated open-ended scientific discovery. *arXiv*, 2024.
- 634
635 OpenAI. Gpt-4 technical report. *arXiv*, 2023.
- 636
637 Judea Pearl. *Causality*. Cambridge university press, 2009.
- 638
639 Mohamad Amin Pourhoseingholi, Ahmad Reza Baghestani, and Mohsen Vahedi. How to control
640 confounding effects by statistical analysis. *Gastroenterology and hepatology from bed to bench*, 5
641 (2):79, 2012.
- 642
643 PyWhy-LLM. <https://github.com/py-why/>.
- 644
645 Biqing Qi, Kaiyan Zhang, Haoxiang Li, Kai Tian, Sihang Zeng, Zhang-Ren Chen, and Bowen
646 Zhou. Large language models are zero shot hypothesis proposers. In *NeurIPS 2023 Workshop on
647 Instruction Tuning and Instruction Following*, 2023.
- 648
649 Linlu Qiu, Liwei Jiang, Ximing Lu, Melanie Sclar, Valentina Pyatkin, Chandra Bhagavatula, Bailin
650 Wang, Yoon Kim, Yejin Choi, Nouha Dziri, et al. Phenomenal yet puzzling: Testing inductive
651 reasoning capabilities of language models with hypothesis refinement. In *ICLR*, 2024.
- 652
653 James Robins. A new approach to causal inference in mortality studies with a sustained exposure
654 period—application to control of the healthy worker survivor effect. *Mathematical modelling*, 7
655 (9-12):1393–1512, 1986.

- 648 Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog,
649 M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,
650 Omar Fawzi, et al. Mathematical discoveries from program search with large language models.
651 *Nature*, pages 1–3, 2023.
- 652 Karen Sachs, Omar Perez, Dana Pe’er, Douglas A Lauffenburger, and Garry P Nolan. Causal protein-
653 signaling networks derived from multiparameter single-cell data. *Science*, 308(5721):523–529,
654 2005.
- 655 Marco Scutari and Jean-Baptiste Denis. *Bayesian networks: with examples in R*. CRC press, 2021.
- 656 Shikhar Singh, Nuan Wen, Yu Hou, Pegah Alipoormolabashi, Te-Lin Wu, Xuezhe Ma, and Nanyun
657 Peng. Com2sense: A commonsense reasoning benchmark with complementary sentences. *arXiv*,
658 2021.
- 659 David J Spiegelhalter. Learning in probabilistic expert systems. *Bayesian statistics*, 4:447–465, 1992.
- 660 Fiona Anting Tan, Xinyu Zuo, and See-Kiong Ng. Unicausal: Unified benchmark and repository for
661 causal text mining. In *International Conference on Big Data Analytics and Knowledge Discovery*,
662 pages 248–262. Springer, 2023.
- 663 Jin Tian and Judea Pearl. On the testable implications of causal models with hidden variables. *arXiv*,
664 2012.
- 665 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
666 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
667 efficient foundation language models. *arXiv*, 2023.
- 668 Ruibo Tu, Kun Zhang, Bo Bertilson, Hedvig Kjellstrom, and Cheng Zhang. Neuropathic pain
669 diagnosis simulator for causal discovery algorithm evaluation. *Advances in Neural Information*
670 *Processing Systems*, 32, 2019.
- 671 Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada,
672 Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. Zephyr: Direct
673 distillation of lm alignment. *arXiv*, 2023.
- 674 Tyler J VanderWeele and Nancy Staudt. Causal diagrams for empirical legal research: a methodology
675 for identifying causation, avoiding bias and interpreting results. *Law, Probability & Risk*, 10(4):
676 329–354, 2011.
- 677 Aniket Vashishtha, Abbavaram Gowtham Reddy, Abhinav Kumar, Saketh Bachu, Vineeth N Balasub-
678 ramanian, and Amit Sharma. Causal inference using llm-guided discovery. *arXiv*, 2023.
- 679 Aniket Vashishtha, Abhinav Kumar, Abbavaram Gowtham Reddy, Vineeth N Balasubramanian, and
680 Amit Sharma. Teaching transformers causal reasoning through axiomatic training. *arXiv*, 2024.
- 681 Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak,
682 Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. Scientific discovery in the age of artificial
683 intelligence. *Nature*, 620(7972):47–60, 2023a.
- 684 Ruocheng Wang, Eric Zelikman, Gabriel Poesia, Yewen Pu, Nick Haber, and Noah D Goodman.
685 Hypothesis search: Inductive reasoning with language models. *arXiv*, 2023b.
- 686 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
687 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS*,
688 2022.
- 689 Fangzhi Xu, Qika Lin, Jiawei Han, Tianzhe Zhao, Jun Liu, and Erik Cambria. Are large language
690 models really good logical reasoners? a comprehensive evaluation from deductive, inductive and
691 abductive views. *arXiv*, 2023a.
- 692 Yudong Xu, Wenhao Li, Pashootan Vaezipoor, Scott Sanner, and Elias B Khalil. Llms and the abstrac-
693 tion and reasoning corpus: Successes, failures, and the importance of object-based representations.
694 *arXiv*, 2023b.

702 Haoran Yang, Yumeng Zhang, Jiaqi Xu, Hongyuan Lu, Pheng-Ann Heng, and Wai Lam. Unveiling
703 the generalization power of fine-tuned large language models. In Kevin Duh, Helena Gomez, and
704 Steven Bethard, editors, *NAACL:HLT*, 2024.

705 Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural networks.
706 In *ICML*, pages 7154–7163. PMLR, 2019.

707
708 Matej Zečević, Moritz Willig, Devendra Singh Dhama, and Kristian Kersting. Causal parrots: Large
709 language models may talk causality but are not causal. *TMLR*, 2023.

710
711 Jiaqi Zhang, Joel Jennings, Agrin Hilmkil, Nick Pawlowski, Cheng Zhang, and Chao Ma. Towards
712 causal foundation model: on duality between optimal balancing and attention. In *ICML*, 2024.

713
714 Jiayao Zhang, Hongming Zhang, Weijie Su, and Dan Roth. Rock: Causal inference principles for
715 reasoning about commonsense causality. In *ICML*, pages 26750–26771. PMLR, 2022.

716
717 Li Zhang, Qing Lyu, and Chris Callison-Burch. Reasoning about goals, steps, and temporal ordering
718 with wikipediqa. *arXiv*, 2020.

719
720 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
721 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
722 chatbot arena. *arXiv*, 2023.

723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A IMPLEMENTATION

A.1 DATASETS

We use 7 real-world based datasets. These datasets span different domain knowledge topics. These datasets have ground truth graphs along with their observational data. The simplest dataset used is the cancer dataset with 4 edges and 5 node variables. In addition to the semi-synthetic datasets from the BNLearn library, we also evaluate our approach on a realistic Alzheimer’s Disease dataset Abdulaal et al. (2024), which was developed by five domain experts. Given that each expert created a different causal graph, the final causal DAG comprises only those edges that were agreed upon by consensus.

Dataset	V	E	Description
Cancer	5	4	Factors around lung cancer
Survey	6	6	Factors for choosing transportation
Asia	8	8	Factors affecting dyspnea
Law	8	20	factors around legal system
Alzheimer	9	16	Factors around Alzheimer’s Disease
Child	20	25	Lung related illness for a child
Insurance	27	52	Factors affecting car accident insurance
Alarm	37	46	Patient monitoring system

Table 4: Dataset description.

A.2 REPRODUCIBILITY

For reproducibility, we used temperature 0 and top-p value as 1 across all of the models. We also mentioned the snapshot of the model used. We have also included the prompts and examples below. Our code can be anonymously found here - <https://anonymous.4open.science/r/causal-llm-env-6C8E/README.md>. The datasets are under CC BY-SA 3.0 which allows us to freely modify the datasets for benchmarking. Our benchmark will be released under the CC BY-SA License.

GPT-3.5 GPT-4 were accessed via API. Rest of the models were run on 1 A100 GPU. Since we used off-the-shelf LLM, there was no training to be performed. Since many of the models were run by API, it is difficult to calculate the entire compute, however, all of the experiments for each model took ≈ 6 hours.

A.3 CONTROLLED VARIABLE IDENTIFICATION

For variable identification, we generate multiple choices that remain consistent across all missing nodes and all of the datasets. The words were randomly chosen to be far enough from the nodes. The options chosen were weather, book sales, and movie ratings. We wanted to make sure that the options were not from one specific domain such that the LLM could do the process of elimination.

A.4 SEMANTIC SIMILARITY

Given the task of hypothesizing missing nodes in a partial graph \mathcal{G}^* in the absence of multiple-choices, we evaluate the semantic similarity between the model’s predictions and the ground truth node variable. We leverage an open model namely ‘all-mpnet-base-v2’ to transform the textual representations of the model’s predictions and the ground truth into high-dimensional vector space embeddings. Post transforming textual representations into embeddings and normalizing them, we calculate the cosine similarity. Scores closer to 1 indicate a high semantic similarity, suggesting the model’s predictions align well with the ground truth. This metric gives a score of similarity without the contextual knowledge of the causal graph. We perform our experiments to consider every node of the ground truth as a missing node iteratively. For all the suggestions for a node variable, we calculate the semantic similarity. The average similarity reported is the highest semantic similarity for each of the variable suggestions.

Algorithm 1 Evaluating Semantic Similarity for Hypothesized Missing Nodes

```

1: Input: Partial graph  $\mathcal{G}^*$ , Ground truth node variables  $V_{GT}$ , Language model  $LM =$ 
   'all-mpnet-base-v2'
2: Output: Average highest semantic similarity score
3: procedure SEMANTICSIMILARITY( $\mathcal{G}^*$ ,  $V_{GT}$ ,  $LM$ )
4:   Initialize similarityScores as an empty list
5:   for each node  $v_{GT}$  in  $\mathbf{v}$  do
6:      $predictions \leftarrow \text{GeneratePredictions}(\mathcal{G}^*, LM)$ 
7:     Initialize nodeScores as an empty list
8:     for each prediction  $p$  in predictions do
9:        $embedding_{GT} \leftarrow \text{Embed}(v_{GT}, LM)$ 
10:       $embedding_p \leftarrow \text{Embed}(p, LM)$ 
11:      Normalize  $embedding_{GT}$  and  $embedding_p$ 
12:       $score \leftarrow \text{CosineSimilarity}(embedding_{GT}, embedding_p)$ 
13:      Append  $score$  to nodeScores
14:     end for
15:      $maxScore \leftarrow \text{Max}(nodeScores)$ 
16:     Append  $maxScore$  to similarityScores
17:   end for
18:    $averageScore \leftarrow \text{Average}(similarityScores)$ 
19:   return  $averageScore$ 
20: end procedure

```

Ground Truth:	Smoking status	Alcohol Consumption	Exposure to Radiation	Poor Diet	Genetic Predisposition
<i>LLM Suggestions:</i>	Smoking				
Semantic similarity :	0.72	0.38	0.22	0.22	0.17
Ground Truth:	Employee or self-employed	Job Location	Environmental Awareness	Lifestyle Preferences	Health Consciousness
<i>LLM Suggestions:</i>	Income Level				
Semantic similarity :	0.30	0.25	0.17	0.15	0.10
Ground Truth:	Dyspnea laboured breathing	Chest Pain	Coughing	Fatigue	Weight Loss
<i>LLM Suggestions:</i>	Shortness of breath				
Semantic similarity :	0.57	0.41	0.36	0.29	0.11
Ground Truth:	Montreal Cognitive Assessment score	Neurological Function	Mental Health Status	Risk of Alzheimer's Disease	Memory Performance
<i>LLM Suggestions:</i>	Cognitive Function				
Semantic similarity :	0.60	0.47	0.38	0.36	0.16
Ground Truth:	Grunting in infants	Asthma	Pneumonia	Pulmonary infection	Bronchopulmonary dysplasia
<i>LLM Suggestions:</i>	Respiratory distress				
Semantic similarity :	0.22	0.18	0.17	0.11	0.01
Ground Truth:	Driving history	Distance driven daily	Type of car insurance	Frequency of car maintenance	Location of parking
<i>LLM Suggestions:</i>	Previous accidents				
Semantic similarity :	0.55	0.42	0.27	0.26	0.18
Ground Truth:	Heart rate blood pressure	Blood Pressure	Respiratory Rate	EKG Reading	Blood Oxygen Level
<i>LLM Suggestions:</i>	Pulse Rate				
Semantic similarity :	0.78	0.78	0.57	0.49	0.42

Table 5: Examples of model suggestions from and the corresponding semantic similarity score for a missing node variable from each of the datasets.

A.5 LLM-AS-JUDGE

To capture the domain knowledge of the expert that selects the most relevant causal variable, we use LLM-as-Judge as a proxy expert. This also allows for evaluation based on contextual DAG knowledge as well. Given the impressive results of GPT-4 in (Zheng et al., 2023), we use GPT-4 as a judge for all of the experiments.

Algorithm 2 Evaluating Model Suggestions with LLM as Judge

```

1: Input: Partial graph  $\mathcal{G}^*$ , Ground truth node variables  $V_{GT}$ , Predictions  $P$ , Language model LLM = GPT-4
2: Output: Average quality rating of model’s suggestions
3: procedure LLMASJUDGE( $\mathcal{G}^*$ ,  $V_{GT}$ ,  $P$ , LLM)
4:   Initialize qualityRatings as an empty list
5:   for each node  $v_{GT}$  in  $\mathbf{V}$  do
6:     suggestions  $\leftarrow$  GenerateSuggestions( $\mathcal{G}^*$ ,  $P$ , LLM)
7:     bestSuggestion  $\leftarrow$  SelectBestSuggestion(suggestions,  $v_{GT}$ , LLM)
8:     rating  $\leftarrow$  RateSuggestion(bestSuggestion, LLM)
9:     Append rating to qualityRatings
10:  end for
11:  averageRating  $\leftarrow$  Average(qualityRatings)
12:  return averageRating
13: end procedure
14: function GENERATESUGGESTIONS( $\mathcal{G}^*$ ,  $P$ , LLM)
15:  return A set of suggestions for missing nodes based on  $P$ 
16: end function
17: function SELECTBESTSUGGESTION(suggestions,  $v_{GT}$ , LLM)
18:  Prompt LLM with  $\mathcal{G}^*$ ,  $v_{GT}$ , and suggestions
19:  return LLM’s choice of the best fitting suggestion
20: end function
21: function RATESUGGESTION(suggestion,  $LM$ )
22:  Prompt LLM to rate suggestion on a scale of 1 to 10
23:  return LLM’s rating
24: end function

```

Ground Truth:	Education up to high school or university degree
<i>Top ranked suggestion:</i>	Education level
Rating :	9.5
Ground Truth:	Pollution
<i>Top ranked suggestion:</i>	Smoking history
Rating :	2.0
Ground Truth:	Bonchitis
<i>Top ranked suggestion:</i>	smoking behavior
Rating :	2.0
Ground Truth:	Lung XRay report
<i>Top ranked suggestion:</i>	Lung Damage
Rating :	8.0
Ground Truth:	Socioeconomic status
<i>Top ranked suggestion:</i>	Driver’s lifestyle
Rating :	7.0

Table 6: Examples of model suggestions from and the corresponding LLM-as-judge score for a missing node variable.

Shortcomings of LLM-as-judge. LLM-as-judge uses GPT-4 as a judge model which could be biased towards some data. Since the training datasets are not public for this model, it would be hard

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Ground Truth: Dyspnea laboured breathing
LLM Suggestion: Shortness of breath
Semantic similarity to GT: 0.57
LLM-as-Judge score: 9.5

Table 7: Example comparing the semantic similarity and LLM-as-Judge metrics. Dyspnea is a medical term for shortness of breath. In this example, the contextual information, beyond exact matching, is better captured by LLM-as-Judge.

to judge how these biases might affect the final score. Hence for robust evaluation we also evaluate using the semantic similarity.

A.6 ITERATIVELY HYPOTHESIZING IN OPEN WORLD

For each order, the algorithm prompts the LLM to generate mediator suggestions, selects the suggestion with the highest semantic similarity to the context, and iteratively updates the partial graph with these mediators. Δ , quantifies the impact of mediator ordering by comparing the average highest semantic similarity scores obtained from both descending and ascending orders. This methodical evaluation sheds light on how the sequence in which mediators are considered might affect the LLM’s ability to generate contextually relevant and accurate predictions.

Algorithm 3 Random Order Mediator Hypothesis

```

1: Input: Partial graph  $\mathcal{G}^*$  (where  $\mathcal{G}^* = \mathcal{G} - H$ ), Treatment  $v_t$ , Outcome  $v_y$ , Number of mediators
    $H$ , Number of suggestions  $k$ 
2: Output: Updated graph  $\mathcal{G}^*$  with selected mediators
3: procedure GENERATEMEDIATORSRANDOM( $\mathcal{G}^*$ ,  $v_t$ ,  $v_y$ ,  $H$ ,  $k$ )
4:   for  $i \leftarrow 1$  to  $H$  do
5:      $suggestions \leftarrow$  Generate  $k$  suggestions for  $v_{m_i}$  using  $P_{LLM}(\mathcal{G}^*)$ 
6:     Initialize  $highestSimilarity \leftarrow 0$ 
7:     Initialize  $selectedMediator \leftarrow$  null
8:     for each  $suggestion$  in  $suggestions$  do
9:        $similarityScore \leftarrow$  Calculate semantic similarity for  $suggestion$ 
10:      if  $similarityScore > highestSimilarity$  then
11:         $highestSimilarity \leftarrow similarityScore$ 
12:         $selectedMediator \leftarrow suggestion$ 
13:      end if
14:    end for
15:    Update  $\mathcal{G}^* \leftarrow \mathcal{G}^* \cup \{selectedMediator\}$ 
16:  end for
17:  return  $\mathcal{G}^*$ 
18: end procedure

```

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Algorithm 4 Ordered Mediator Generation and Evaluation Based on MIS

```

1: Input: Partial graph  $\mathcal{G}^*$ , Treatment  $v_t$ , Outcome  $v_y$ , Set of potential mediators  $M$ , Number of
   suggestions  $k$ 
2: Output:  $\Delta$  - measure of the influence of mediator ordering
3: procedure CALCULATEMIS( $v_t, v_y, M$ )
4:   Initialize MISList as an empty list
5:   for each mediator  $v_{m_i}$  in  $M$  do
6:     Calculate NIE( $v_{m_i}$ ) and NDE( $v_{m_i}$ )
7:      $MIS(v_{m_i}) \leftarrow \frac{NIE(v_{m_i})}{NDE(v_{m_i})}$ 
8:     Append MIS( $v_{m_i}$ ) to MISList
9:   end for
10:  return MISList
11: end procedure
12: procedure GENERATEMEDIATORSORDERED( $\mathcal{G}^*, v_t, v_y, M, k$ )
13:  MISList  $\leftarrow$  CALCULATEMIS( $v_t, v_y, M$ )
14:  Sort  $M$  in descending order of MISList to get  $M_{desc}$ 
15:  Sort  $M$  in ascending order of MISList to get  $M_{asc}$ 
16:   $averageDesc \leftarrow$  GENERATEANDEVALUATE( $\mathcal{G}^*, M_{desc}, k$ )
17:   $averageAsc \leftarrow$  GENERATEANDEVALUATE( $\mathcal{G}^*, M_{asc}, k$ )
18:   $\Delta \leftarrow \frac{|averageDesc - averageAsc|}{averageDesc}$ 
19:  return  $\Delta$ 
20: end procedure
21: function GENERATEANDEVALUATE( $\mathcal{G}^*, M_{order}, k$ )
22:  Initialize similarityScores as an empty list
23:  for each mediator  $v_{m_i}$  in  $M_{order}$  do
24:    Perform the same steps as in the refined random order mediator generation
25:    (Generate  $k$  suggestions, select the most similar, update  $\mathcal{G}^*$ )
26:    Append the highest similarity score to similarityScores
27:  end for
28:  return Average of similarityScores
29: end function

```

B CONFOUNDERS

	Sachs	Alarm1	Alarm2	Ins1	Ins2	Ins3	Ins4	Ins5	Ins6	Ins7
Zephyr	0.12	0.37	0.29	0.45	0.49	0.37	0.29	0.33	0.46	0.73
Mixtral	0.89	0.54	0.57	0.57	1.0	0.32	0.23	0.38	0.28	1.0
Neural	0.34	0.27	0.28	0.42	0.47	0.34	0.48	0.48	0.38	0.48
LLama	0.27	0.39	0.44	0.55	1.0	0.29	0.22	0.57	0.45	1.0
Mistral	0.23	0.62	0.46	0.58	1.0	0.28	0.28	0.28	0.28	1.0
GPT-3.5	0.34	0.39	0.48	0.48	1.0	0.58	0.20	0.48	0.47	1.0
GPT-4	0.91	0.49	0.44	0.62	0.39	0.58	0.44	0.58	0.52	1.0

Table 8: Semantic similarity

	Sachs	Alarm1	Alarm2	Ins1	Ins2	Ins3	Ins4	Ins5	Ins6	Ins7
Zephyr	0.10	0.40	0.30	0.45	0.60	0.40	0.40	0.30	0.70	0.80
Mixtral	0.95	0.70	1.0	0.75	1.0	0.80	0.20	0.20	0.20	1.0
Neural	0.30	0.60	0.30	1.0	0.60	0.30	0.80	0.30	0.40	0.60
LLama	0.20	0.50	0.44	0.40	1.0	0.50	0.20	0.70	0.45	1.0
Mistral	0.20	0.90	0.80	0.55	1.0	0.30	0.20	0.70	0.30	1.0
GPT-3.5	0.40	0.50	0.48	0.30	1.0	0.75	0.40	0.75	0.60	1.0
GPT-4	0.95	0.65	0.80	0.60	0.70	0.80	0.85	0.80	0.75	1.0

Table 9: LLM judge

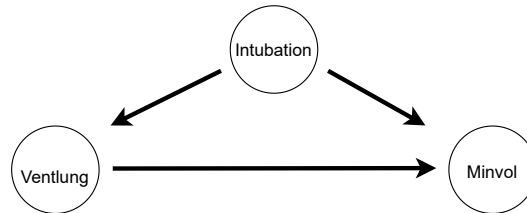


Figure 6: Alarm 1

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

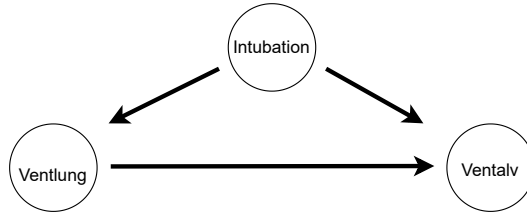


Figure 7: Alarm 2

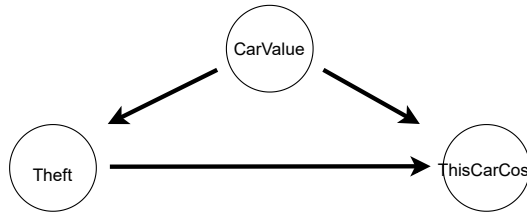


Figure 8: Insurance 1

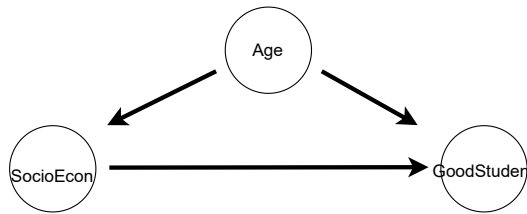


Figure 9: Insurance 2

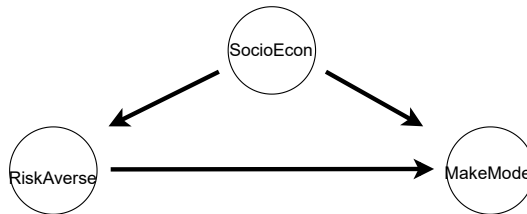


Figure 10: Insurance 3

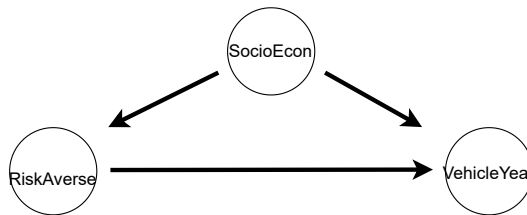


Figure 11: Insurance 4

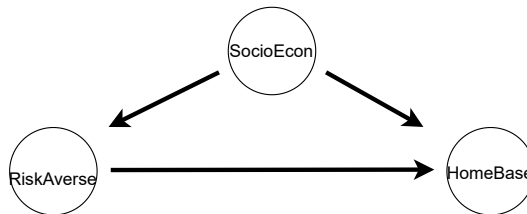


Figure 12: Insurance 5

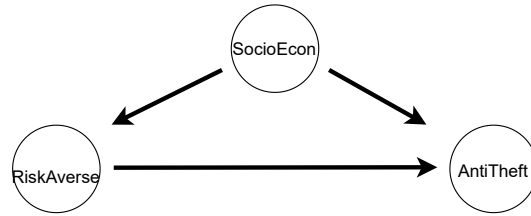


Figure 13: Insurance 6

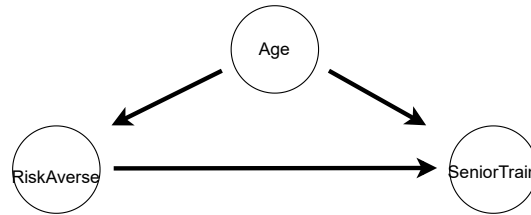


Figure 14: Insurance 7

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

C FURTHER RESULTS

C.1 VARIANCES

For brevity we didnt add variance in the main text, the following results have variances:

	Cancer		Survey		Asia		Alzheimers		Child		Insurance		Alarm		Avg	
	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J	Sim	LLM-J
Zephyr	0.36 ±0.04	0.61 ±0.06	0.34 ±0.07	0.60 ±0.05	0.45 ±0.05	0.66 ±0.04	0.35 ±0.03	0.75 ±0.03	0.51 ±0.02	0.70 ±0.04	0.45 ±0.04	0.44 ±0.05	0.46 ±0.03	0.69 ±0.02	0.42 ±0.04	0.63 ±0.04
Mixtral	0.41 ±0.03	0.66 ±0.04	0.39 ±0.05	0.66 ±0.06	0.66 ±0.02	0.75 ±0.03	0.31 ±0.04	0.77 ±0.02	0.53 ±0.03	0.77 ±0.02	0.46 ±0.03	0.56 ±0.04	0.50 ±0.03	0.72 ±0.06	0.46 ±0.03	0.70 ±0.05
Neural	0.38 ±0.02	0.77 ±0.05	0.43 ±0.02	0.55 ±0.03	0.53 ±0.03	0.55 ±0.04	0.44 ±0.05	0.71 ±0.03	0.48 ±0.04	0.70 ±0.03	0.47 ±0.04	0.43 ±0.05	0.47 ±0.02	0.67 ±0.03	0.45 ±0.03	0.63 ±0.04
Llama	0.40 ±0.03	0.48 ±0.05	0.40 ±0.04	0.54 ±0.05	0.53 ±0.03	0.58 ±0.06	0.45 ±0.05	0.61 ±0.03	0.48 ±0.04	0.63 ±0.03	0.42 ±0.01	0.34 ±0.05	0.46 ±0.02	0.65 ±0.03	0.45 ±0.03	0.55 ±0.04
Mistral	0.33 ±0.01	0.67 ±0.05	0.44 ±0.05	0.65 ±0.04	0.60 ±0.03	0.73 ±0.04	0.34 ±0.04	0.76 ±0.02	0.48 ±0.04	0.68 ±0.03	0.46 ±0.03	0.47 ±0.01	0.47 ±0.03	0.71 ±0.03	0.44 ±0.03	0.67 ±0.03
GPT-3.5	0.48 ±0.03	0.74 ±0.04	0.42 ±0.00	0.79 ±0.03	0.47 ±0.04	0.61 ±0.04	0.39 ±0.05	1.00 ±0.00	0.36 ±0.05	0.60 ±0.05	0.47 ±0.07	0.52 ±0.02	0.48 ±0.04	0.73 ±0.05	0.44 ±0.04	0.71 ±0.03
GPT-4	0.49 ±0.02	0.90 ±0.03	0.51 ±0.06	0.67 ±0.04	0.66 ±0.02	0.76 ±0.03	0.47 ±0.02	0.98 ±0.02	0.36 ±0.05	0.53 ±0.04	0.52 ±0.03	0.56 ±0.03	0.49 ±0.06	0.75 ±0.02	0.50 ±0.04	0.73 ±0.03

Table 10: Average semantic similarity and LLM-as-Judge metrics to evaluate LLMs in hypothesizing the missing variable in a causal DAG.

C.2 ANALYSIS OF DIFFERENCE ACROSS TASKS

Since the metrics are different to evaluate each task, it is not meaningful or straightforward to compare the raw results. It must also be noted that the tasks are not linear. To address this, we rank the model performances across all models and datasets and present these rankings in Figure 15. This allows us to compare the relative performance of the models across different tasks.

As we observe from the graph, GPT-4 model shows consistently top performances in Tasks 1-3, however, it has one of the lowest performances for Task 4. GPT-3.5 shows a strong performance in Task 2 and 4, ranking 2nd, but drops in Tasks 1 and 3. We observe that Zephyr, Neural and Mistral show consistently average performances. These observations motivate the significance of the tasks proposed in our benchmark. They highlight the variability in model performance across different tasks and emphasize the need for comprehensive and diverse benchmarks to fully assess the capabilities of these models.

	Asia		Child		Insurance		Alarm	
	Sim	Δ	Sim	Δ	Sim	Δ	Sim	Δ
Zephyr	0.61 ±0.03	-0.02 ±0.01	0.54 ±0.04	0.17 ±0.02	0.47 ±0.05	0.19 ±0.02	0.51 ±0.05	0.20 ±0.02
Mixtral	0.87 ±0.02	0.01 ±0.01	0.50 ±0.05	0.18 ±0.02	0.48 ±0.05	0.15 ±0.02	0.52 ±0.05	0.13 ±0.01
Neural	0.65 ±0.06	0.04 ±0.02	0.48 ±0.05	0.21 ±0.02	0.42 ±0.04	0.16 ±0.02	0.46 ±0.04	0.12 ±0.01
Llama	0.80 ±0.08	0.07 ±0.02	0.49 ±0.05	-0.05 ±0.01	0.44 ±0.06	0.21 ±0.02	0.51 ±0.05	0.07 ±0.01
Mistral	0.33 ±0.03	0.02 ±0.01	0.50 ±0.05	0.12 ±0.01	0.48 ±0.05	0.13 ±0.02	0.47 ±0.04	0.11 ±0.01
GPT-3.5	0.48 ±0.05	0.01 ±0.01	0.36 ±0.04	0.25 ±0.03	0.48 ±0.05	0.17 ±0.02	0.51 ±0.05	0.02 ±0.01
GPT-4	0.49 ±0.07	0.04 ±0.01	0.39 ±0.05	0.16 ±0.02	0.52 ±0.05	0.14 ±0.02	0.60 ±0.06	-0.07 ±0.01

Table 11: Sim: semantic similarity for iteratively hypothesizing the mediator nodes when prompted with random order. Δ measures the change in the prediction of each model according to the MIS.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

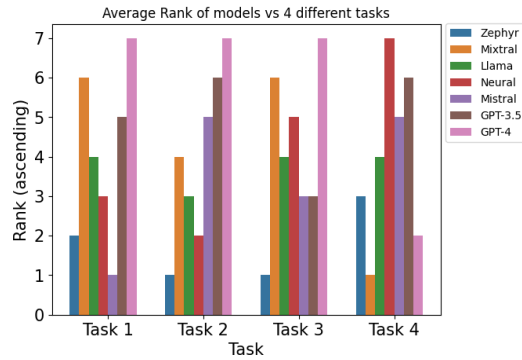


Figure 15: Average Rank of each model against the different tasks. We ranked the mode since the metrics are different to evaluate each task averaged across datasets

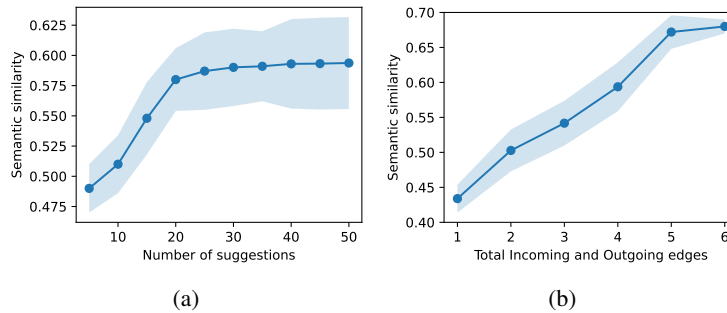


Figure 16: L: Plot of semantic similarity with an increasing number of suggestions for GPT-4 on the Alarm dataset. R: Plot of semantic similarity against the total number of incoming and outgoing edges for GPT-4 on the Alarm dataset.

C.3 BREAKING DOWN THE PERFORMANCE

C.4 EFFECT OF CONTEXT

We observed notable differences in the accuracy of LLM predictions for missing nodes within causal graphs when context was provided versus when it was absent. Specifically, the inclusion of contextual information about the causal graph significantly enhanced the LMs’ ability to generate accurate and relevant predictions. In realistic settings, when this setup is being used by a scientist, they would provide the context of the task along with the partial graph. When context was not provided, the models often struggled to identify the most appropriate variables, leading to a decrease in prediction accuracy, especially for smaller models. Unsurprisingly, providing context was more important for smaller graphs than larger graphs. LLMs were able to understand the context of the graph via multiple other nodes in the graph for larger graphs.

	Cancer		Survey		Asia		Insurance		Alarm	
	<i>X</i>	✓	<i>X</i>	✓	<i>X</i>	✓	<i>X</i>	✓	<i>X</i>	✓
In-Context	0.75	1.00	0.67	1.00	0.68	0.88	0.85	0.90	0.96	0.96
Out-of-Context	0.00	0.25	0.33	0.33	0.53	0.61	0.58	0.58	0.60	0.57
Open world Hypothesis	0.39	0.41	0.40	0.39	0.63	0.66	0.49	0.50	0.44	0.46

Table 12: Model-Mixtral to evaluate the effect of context given in the prompt.

C.5 USING EXPLANATIONS

While using LLMs for hypothesizing the missing nodes within the causal graph for the open world setting, introduced an additional question to prompt the model to provide explanations for each of

their predictions. This was motivated by the fact that incorporating a rationale behind each prediction might enhance the model’s semantic similarity. We present the results in the Table below: We observe that evaluating semantic similarity with explanations leads to a decrease in performance as compared to the earlier setting where the language model returned phrases. This is because semantic similarity, as a metric, evaluates the closeness of the model’s predictions to the ground truth in a high-dimensional vector space, focusing on the semantic content encapsulated within the embeddings. It is a metric that leaves little room for interpretative flexibility, focusing strictly on the degree of semantic congruence between the predicted and actual variables. The introduction of explanations, while enriching the model’s outputs with contextual insights, did not translate into improved semantic alignment with the ground truth.

	Cancer		Survey		Asia		Insurance		Alarm	
	X	✓	X	✓	X	✓	X	✓	X	✓
Sim	0.49	0.38	0.51	0.44	0.66	0.57	0.52	0.40	0.49	0.40
	± 0.02	± 0.07	± 0.06	± 0.10	± 0.02	± 0.09	± 0.03	± 0.07	± 0.06	± 0.06
LLM-Judge	0.90	0.91	0.67	0.69	0.76	0.76	0.56	0.55	0.75	0.75
	± 0.03	± 0.02	± 0.04	± 0.02	± 0.03	± 0.04	± 0.03	± 0.03	± 0.02	± 0.02

Table 13: Model-GPT 4. Evaluating the effect of explanations on different metrics from Task 3.

Ambiguous predictions which semantically represent the same variable. An important linguistic concern that could be missed by semantic similarity is ambiguous hypothesis by the LLM that may have same semantics, which again breaks the semantic similarity metric. This further motivates LLM-judge metric whose input is - the context of the causal graph, the partial causal graph, the ground truth variable, and the model predictions. Given the rich context of the LLM-judge metric we suspect it would be able to overcome the ambiguity. We prompted the model to justify its hypothesis variables using explanations. We observe that evaluating semantic similarity with explanations leads to a decrease in performance as compared to the earlier setting where the language model returned just phrases. In Table 13 we observed a drop in performance for semantic similarity. In contrast, we observe a similar or slight improvement in the LLM-judge metric when the explanation of the model hypothesis is given.

C.6 CHAIN OF THOUGHT

In recent times, Chain-of-Thought prompting has gained popularity due to its impressive performance in proving the quality of LLMs’ output Kojima et al. (2022) also in metadata-based causal reasoning Vashishtha et al. (2023). We also incorporated COT prompting for our prompts. We perform ablation studies in Table. We observe that COT particularly improves the performance of the identification experiments.

	Cancer		Survey		Asia		Insurance		Alarm	
	X	✓	X	✓	X	✓	X	✓	X	✓
In-Context	1.00	1.00	0.83	1.00	0.75	0.88	0.74	0.90	0.91	0.96
Out-of-Context	0.50	0.25	0.18	0.33	0.57	0.61	0.56	0.58	0.54	0.57

Table 14: Model-Mixtral to evaluate the effect of COT given in the prompt.

C.7 ITERATIVE MEDIATOR SEARCH VS ALL AT ONCE

For Task 4, we iteratively hypothesize the missing variables (mediators). Our choice was primarily driven by the complexity of Task 4, which involves predicting multiple missing mediators, ranging from 1 to 10. For a Task with 10 missing mediators, the model would have to predict 50 suggestions at once. We initially hypothesized that LLMs might struggle with making multiple predictions across different variables simultaneously. This was indeed reflected in our results and GPT-4 outputs from Table X. The iterative approach allows the model’s prediction to narrow the search space, which would not be possible in a non-iterative approach. This method is more aligned with the scientific discovery process, where hypotheses are often refined iteratively based on new findings. Furthermore,

our approach simulates a human-in-the-loop scenario, where the most plausible answer is selected and used to guide the next prediction.

	Asia	Child	Insurance	Alarm
Non-iterative	0.42 +- 0.07	0.33 +- 0.06	0.45 +- 0.09	0.54 +- 0.05
Iterative	0.49 +- 0.05	0.39 +- 0.03	0.52 +- 0.02	0.60 +- 0.04

C.8 RESULTS ON NEUROPATHIC DATASET

We added a new dataset, the neuropathic pain dataset Tu et al. (2019), which is not part of common LLM training corpora as one needs to use a python script to download it. The dataset consists of 221 nodes and 770 edges, but for feasibility, we selected a subset of the graph for evaluation. We ran experiments for Task 1, Task 2, and Task 3.

Model	Task 1	Task 2 Result	Task 2 FNA	Task 3 Sim	Task 3 LLM-J
Mistral	0.64	0.51	0.32	0.38	0.53
Mixtral	0.83	0.55	0.34	0.45	0.69
Llama	0.78	0.49	0.27	0.44	0.63
GPT-3.5	0.82	0.53	0.31	0.47	0.72
GPT-4	0.94	0.68	0.24	0.51	0.76

Table 15: Comparison of model performances across tasks on Neuropathic dataset.

C.9 FINE GRAINED MODEL PERFORMANCE

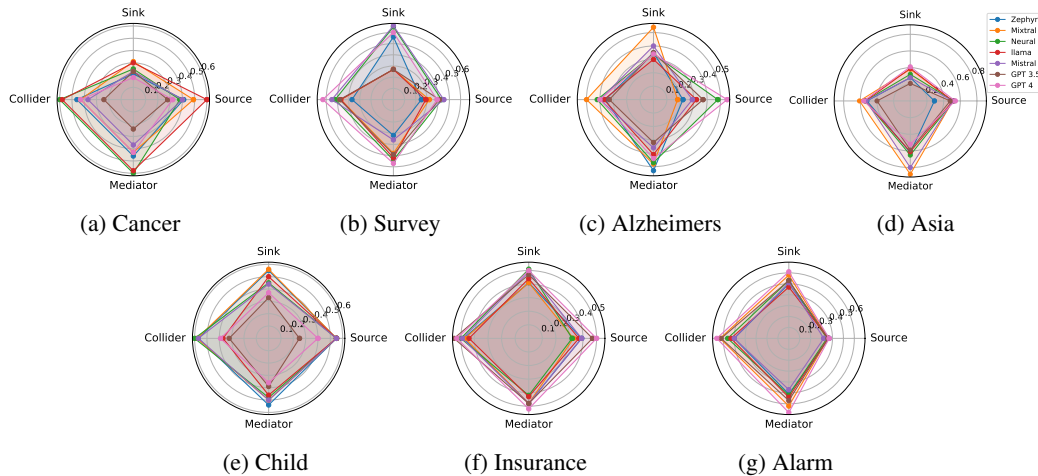


Figure 17: Detailed spider plots for Semantic similarity

D FINETUNING AND FEW-SHOT PROMPTING

D.1 FINETUNING

we aim to assess the LLM’s causal reasoning via prompting. Following are the reasons why fine-tuning is not the most practical solution:

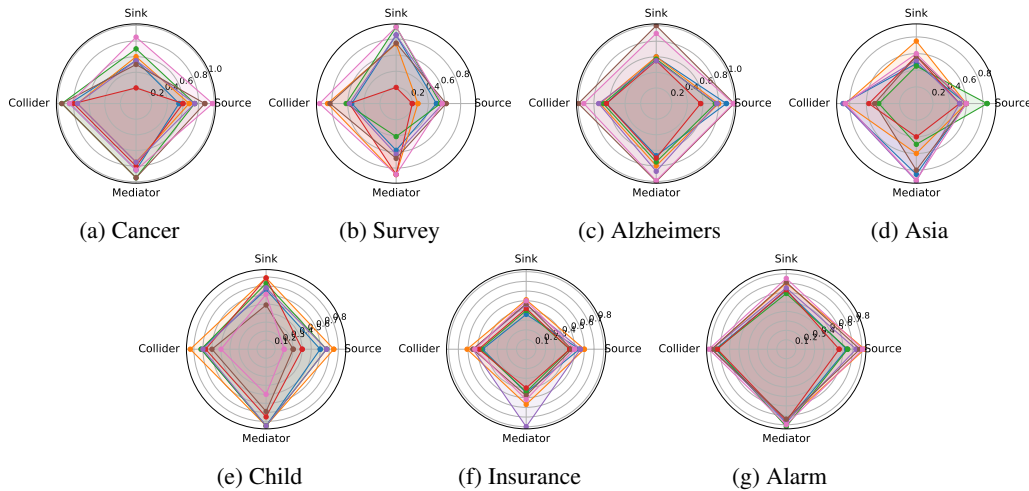


Figure 18: Detailed spider plots for LLM-as-judge metric

- Pretrained models come with a wealth of general knowledge, which we aim to leverage. Fine-tuning these models could potentially limit their ability to draw on this broad knowledge base. We aim to understand the utility of pretrained models, as fine-tuning large models like GPT-4 is not always feasible.
- The training dataset is too small for fine-tuning. Despite considering a large 52-edged graph: Insurance, we would have just 27 datapoints or Alarm with 37 datapoint. Additionally:
 1. Using the same graph as part of train and test would unfortunately lead to training data leakage.
 2. If we consider different graphs for train and test, there would exist a domain shift in the two graphs and the model may be overfitted to the domain of the train graph.

However, to illustrate our hypothesis and alleviate the reviewer’s concern, we performed Supervised Fine-Tuning using QLoRA on the Mistral-7b-Instruct model for hypothesizing in the open world task. The train set here is all of the graphs minus the respective graph it was tested on. We tested on Survey, Insurance and Alzheimers graphs. The model was trained to give one best-fit suggestion for the missing variable.

	Insurance	Survey	Alzheimers
No fine-tuning	0.42 +- 0.03	0.44 +- 0.05	0.34 +- 0.04
Fine-tuned	0.39 +- 0.04	0.39 +- 0.03	0.36 +- 0.07

Table 16: Finetuning results.

From the above results, it is evident that finetuning does not significantly improve over the prompting results. This is because during training the LLM gets biased towards the domains of training datasets which are contextually distant from the test domain, given the diversity of datasets chosen. One may think that training might help the LLM to understand the task, but from prompt-based model output, it was evident that the LLM can instruction-follow. In summary, we were able to extract the LLM knowledge via prompting and domain-specific fine-tuning could be closely looked at in the future works.

D.2 FEWSHOT PROMPTING

Similar to fine-tuning, few-shot learning’s success depends on balancing domain specificity and generality. To avoid test examples becoming part of the shots, we have to use different domains as

1458 examples. Given the complexity of the Alarm graph, we decided to use them as a prior. We performed
1459 experiments with 1-shot and 5-shots for the Mixtral 8x7b model. We would like to remind you that
1460

1461 Dataset	0-shot	1-shot	5-shot
1462 Cancer	0.41	0.43	0.46
1463 Survey	0.39	0.38	0.36
1464 Asia	0.66	0.70	0.72
1465 Alzheimer's	0.31	0.33	0.34
1466 Child	0.53	0.55	0.56
1467 Insurance	0.46	0.42	0.45

1468
1469 Table 17: Fewshot prompting results.
1470

1471 Alarm was a medical dataset which means that providing more examples in a different domain might
1472 hinder the model performance. Drop in performance when changing domain for in-context learning
1473 has been discussed in Kwan et al. (2024) and Gupta et al. (2024).
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

1512 E CAUSAL GRAPHS

1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

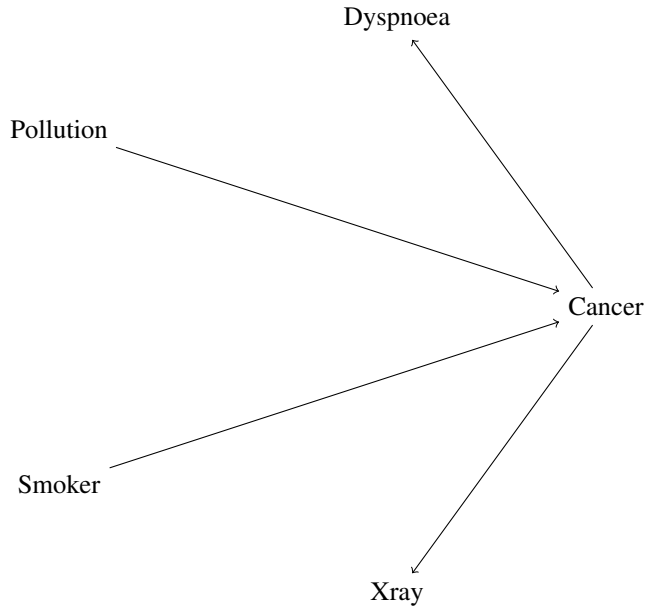


Figure 19: Cancer DAG

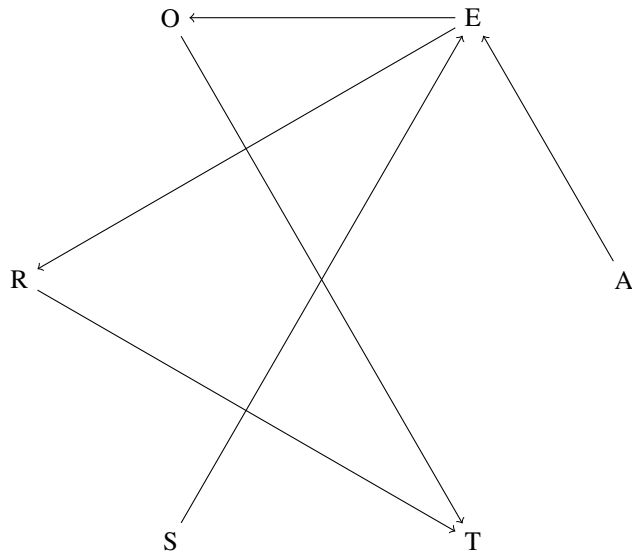


Figure 20: Survey DAG

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

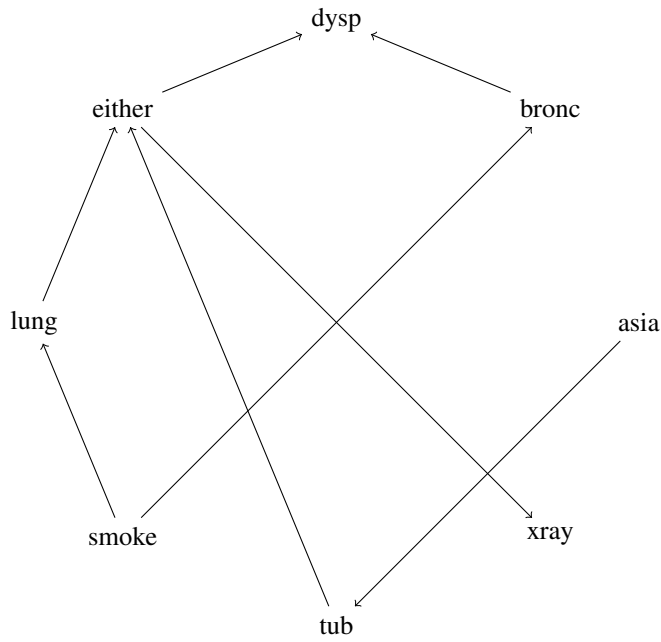


Figure 21: Asia DAG

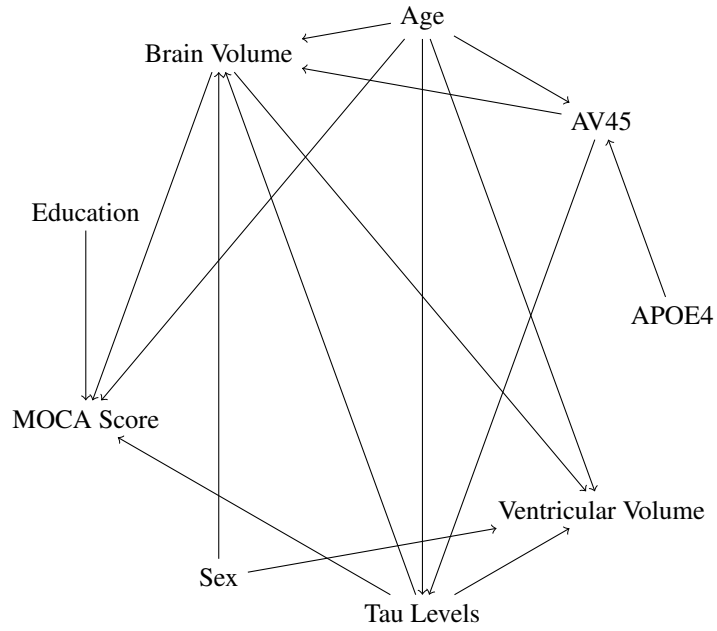


Figure 22: Alzheimer's DAG

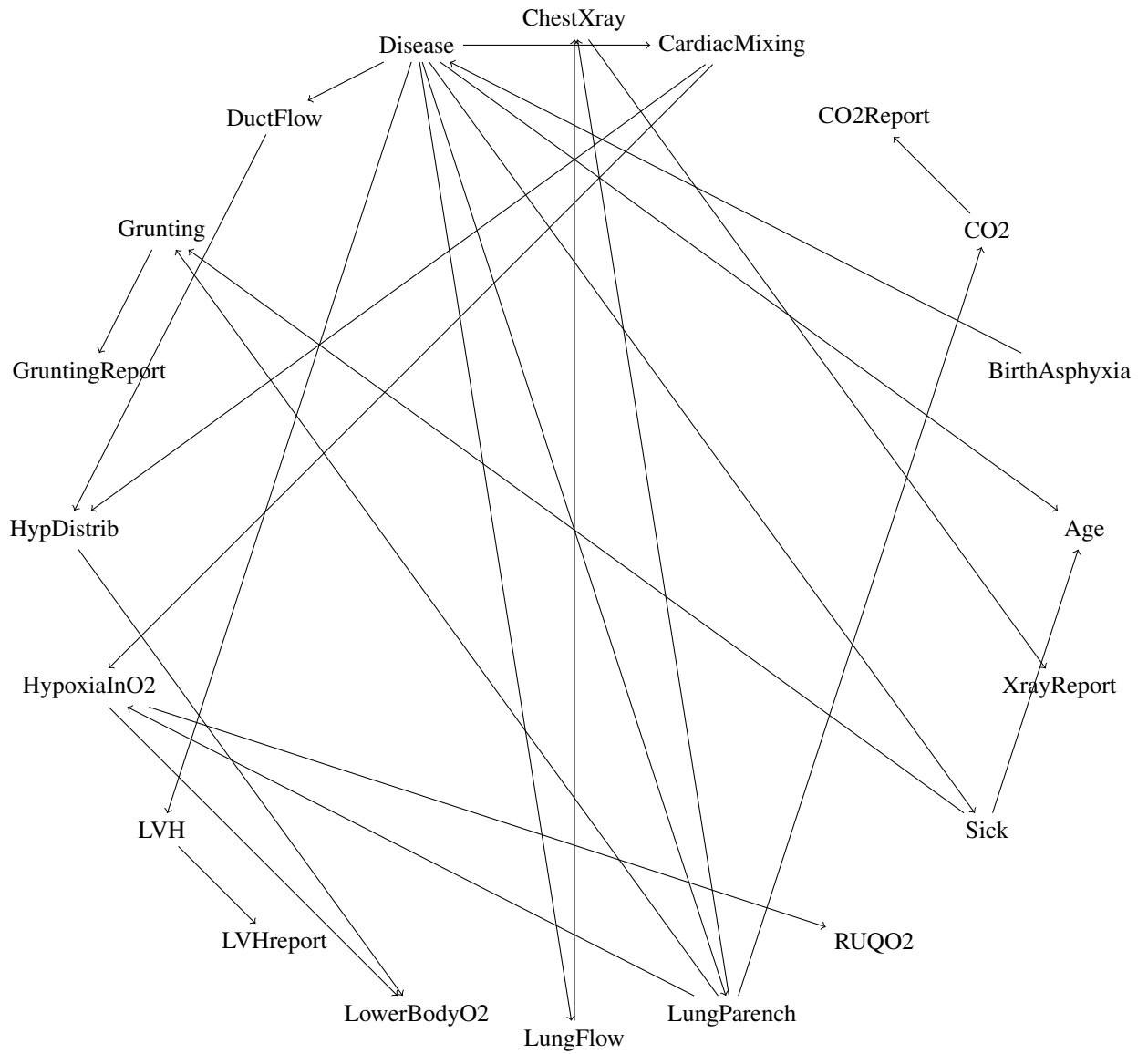


Figure 23: Child DAG

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

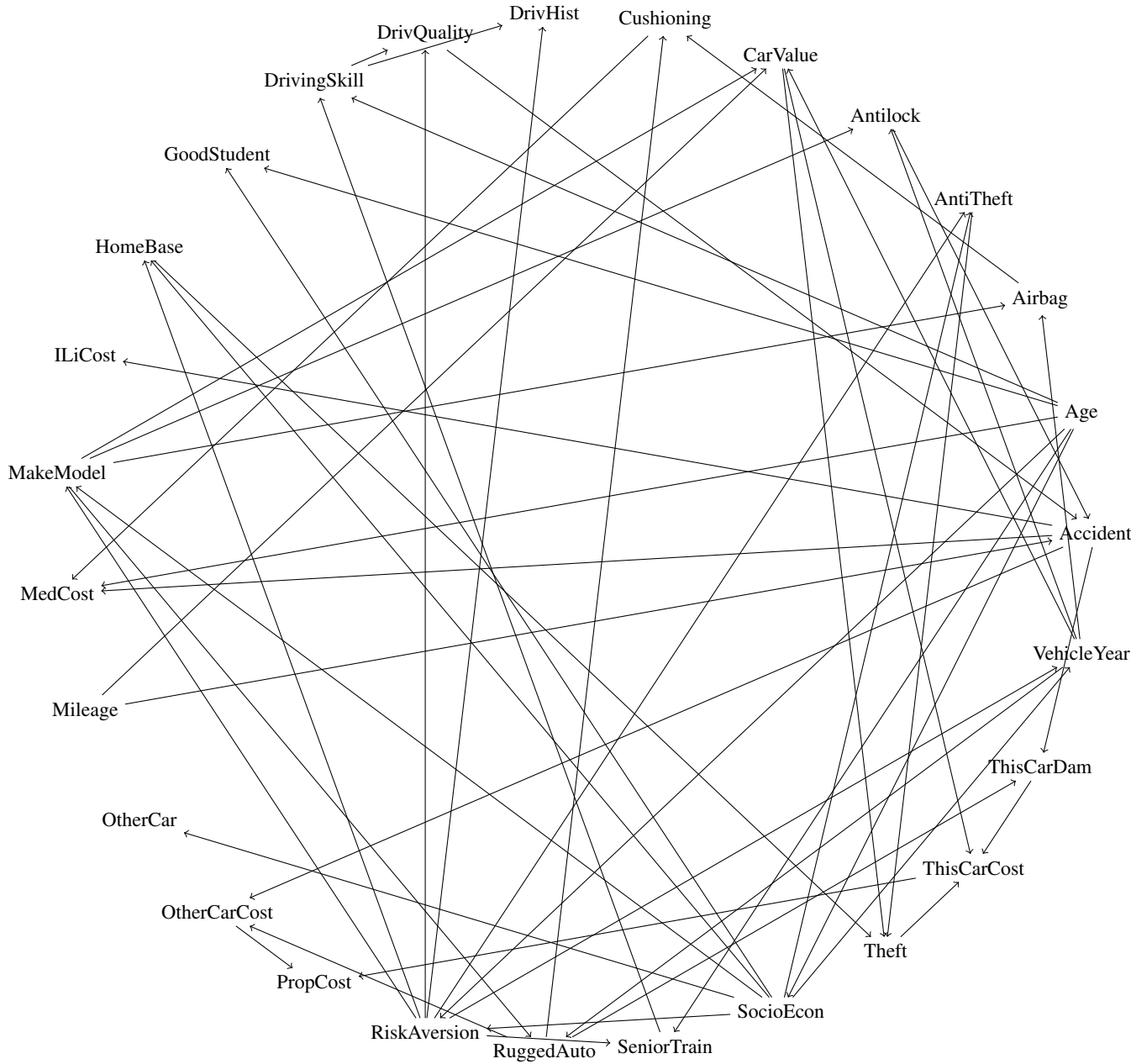


Figure 24: Insurance DAG

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

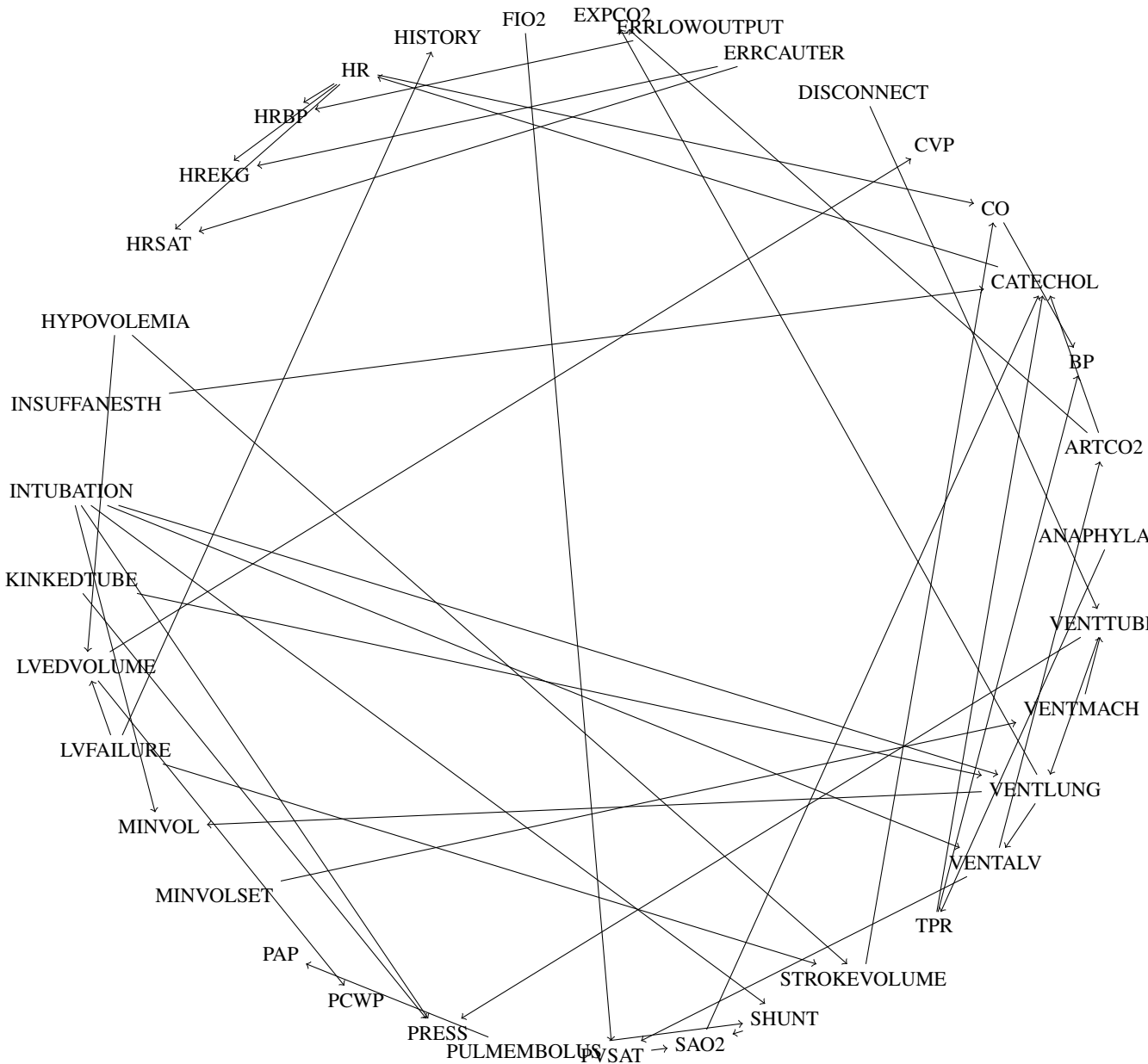


Figure 25: Alarm DAG

F PROMPT TEMPLATE

Hello. You will be given a causal graph. The context of the graph [CONTEXT]. Please understand the causal relationships between the variables - [VERBALISED DAG].

Prompt 1: Base prompt to describe the causal graph

Hello. You will be given a causal graph. The context of the graph is hypothetical patient monitoring system in an intensive care unit (ICU). Please understand the causal relationships between the variables - < anaphylaxis > causes < total peripheral resistance >. < arterial co2 > causes < expelled co2 >. < arterial co2 > causes < catecholamine >. < catecholamine > causes < heart rate >. < cardiac output > causes < blood pressure >. < disconnection > causes < breathing tube >. < error cauter > causes < heart rate displayed on ekg monitor >. < error cauter > causes < oxygen saturation >. < error low output > causes < heart rate blood pressure >. < high concentration of oxygen in the gas mixture > causes < pulmonary artery oxygen saturation >. < heart rate > causes < heart rate blood pressure >. < heart rate > causes < heart rate displayed on ekg monitor >. < heart rate > causes < oxygen saturation >. < heart rate > causes < cardiac output >. < hypovolemia > causes < left ventricular end-diastolic volume >. < hypovolemia > causes < stroke volume >. < insufficient anesthesia > causes < catecholamine >. < intubation > causes < lung ventilation >. < intubation > causes < minute volume >. < intubation > causes < alveolar ventilation >. < intubation > causes < shunt - normal and high >. < intubation > causes < breathing pressure >. < kinked chest tube > causes < lung ventilation >. < kinked chest tube > causes < breathing pressure >. < left ventricular end-diastolic volume > causes < central venous pressure >. < left ventricular end-diastolic volume > causes < pulmonary capillary wedge pressure >. < left ventricular failure > causes < previous medical history >. < left ventricular failure > causes < left ventricular end-diastolic volume >. < left ventricular failure > causes < stroke volume >. < the amount of time using a breathing machine > causes < the intensity level of a breathing machine >. < sudden blockage in the pulmonary arteries > causes < shunt - normal and high >. < sudden blockage in the pulmonary arteries > causes < pulmonary artery pressure >. < pulmonary artery oxygen saturation > causes < oxygen saturation >. < oxygen saturation > causes < catecholamine >. < shunt - normal and high > causes < oxygen saturation >. < stroke volume > causes < cardiac output >. < total peripheral resistance > causes < catecholamine >. < total peripheral resistance > causes < blood pressure >. < alveolar ventilation > causes < arterial co2 >. < alveolar ventilation > causes < pulmonary artery oxygen saturation >. < lung ventilation > causes < expelled co2 >. < lung ventilation > causes < minute volume >. < lung ventilation > causes < alveolar ventilation >. < the intensity level of a breathing machine > causes < breathing tube >. < breathing tube > causes < lung ventilation >. < breathing tube > causes < breathing pressure >.

Prompt 2: An example of the base prompt for Alarm dataset. Each relationship is enclosed in pointed brackets,<> followed by a full stop.

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

Prompt: Hello. You will be given a causal graph. The context of the graph is different aspects of a patient's health related to lung diseases. Please understand the causal relationships between the variables - < x > causes < tuberculosis >. < bronchitis > causes < dyspnoea-laboured breathing >. < individual has either tuberculosis or lung cancer > causes < positive xray >. < individual has either tuberculosis or lung cancer > causes < dyspnoea-laboured breathing >. < lung cancer > causes < individual has either tuberculosis or lung cancer >. < smoking cigarettes > causes < lung cancer >. < smoking cigarettes > causes < bronchitis >. < tuberculosis > causes < individual has either tuberculosis or lung cancer >.Based on your understanding and analytical skills, please identify what the variable X could most likely be from the given multiple choices. It's important to note that you MUST select an option from the provided choices, as one of them is the correct answer. The choices are: [visited Asia; weather; book sales; movie rating]. Please think step by step. Your answer should be formatted as follows - Answer: X = choice. The choice must be one of the provided options, and it should be formatted exactly as it appears in the choices. Generating your own options is not allowed in this task.

Answer: X = visited Asia.My thought process is as follows:The causal graph suggests that X causes tuberculosis, and tuberculosis is related to lung diseases. Visiting Asia is a possible explanation for why an individual may have been exposed to tuberculosis, as the disease is more prevalent in some Asian countries. Therefore, visited Asia is the most likely option for X.

Prompt 3: Out-of-context controlled variable identification, Ground truth variable: visited Asia

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

Prompt: Hello. You will be given a causal graph. The context of the graph is different aspects of a patient's health related to lung diseases Please understand the causal relationships between the variables - $\langle x \rangle$ causes \langle tuberculosis \rangle . \langle bronchitis \rangle causes \langle dyspnoea-laboured breathing \rangle . \langle individual has either tuberculosis or lung cancer \rangle causes $\langle y \rangle$. \langle individual has either tuberculosis or lung cancer \rangle causes \langle dyspnoea-laboured breathing \rangle . \langle lung cancer \rangle causes \langle individual has either tuberculosis or lung cancer \rangle . \langle smoking cigarettes \rangle causes \langle lung cancer \rangle . \langle smoking cigarettes \rangle causes \langle bronchitis \rangle . \langle tuberculosis \rangle causes \langle individual has either tuberculosis or lung cancer \rangle . Based on your understanding and analytical skills, please identify what the variable X could most likely be from the given multiple choices. It's important to note that you MUST select an option from the provided choices, as one of them is the correct answer. The choices are: [weather; book sales; movie rating; positive xray; visited Asia]. Please think step by step. Your answer should be formatted as follows - Answer: X = choice. The choice must be one of the provided options, and it should be formatted exactly as it appears in the choices. Generating your own options is not allowed in this task.

Answer: Answer: X = **visited Asia**. My thought process is as follows:1. The graph shows that X causes tuberculosis, which suggests that X is a variable related to healthcare or medicine. 2. The graph also shows that bronchitis causes dyspnoea-laboured breathing, which further supports the idea that X is a variable related to healthcare or medicine. 3. The graph shows that individual has either tuberculosis or lung cancer causes y, which suggests that X could be a variable related to diagnosing or treating lung diseases. 4. The graph shows that smoking cigarettes causes lung cancer and bronchitis, which suggests that X could be a variable related to smoking or lung health. 5. Based on the above observations, it seems unlikely that X would be weather, book sales, or movie rating.6. Visited Asia is a variable that could potentially relate to healthcare or medicine, particularly if the patient has traveled to a region with a high prevalence of tuberculosis or other lung diseases.7. Therefore, based on the information provided in the graph, it seems most likely that X = **visited Asia**.

Prompt 4: In-context controlled variable identification, Ground truth variable: visited asia

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

Prompt: Hello. You will be given a causal graph. The context of the graph is different aspects of a patient's health related to lung diseases. Please understand the causal relationships between the variables - < x > causes < tuberculosis >. < bronchitis > causes < dyspnoea-laboured breathing >. < individual has either tuberculosis or lung cancer > causes < positive xray >. < individual has either tuberculosis or lung cancer > causes < dyspnoea-laboured breathing >. < lung cancer > causes < individual has either tuberculosis or lung cancer >. < smoking cigarettes > causes < lung cancer >. < smoking cigarettes > causes < bronchitis >. < tuberculosis > causes < individual has either tuberculosis or lung cancer >. Using your causal knowledge, can you suggest what variable X most likely? Give 5 suggestions for what variable X could be. Give reason for your suggestions. Finally, return your answer (without reasoning) in the following format: <Answer> [first suggestion, second suggestion, third suggestion .. and so on] </Answer> No explanations, no reasoning, just answers as a list as described earlier. For example, if the suggestions are A,B,C,D,E then - <Answer> [A, B, C, D, E] </Answer>

Answer: <Answer> [HIV/AIDS, Chronic respiratory conditions, **Pneumonia**, Exposure to asbestos, Silicosis] </Answer>

Reasons:

1. HIV/AIDS: It is known to weaken the immune system, making the individual more susceptible to infections like tuberculosis.
2. Chronic respiratory conditions: These can make an individual more susceptible to lung diseases like tuberculosis.
3. **Pneumonia**: It is a lung infection that can lead to serious conditions like tuberculosis if not treated properly.
4. Exposure to asbestos: Long-term exposure to asbestos can lead to lung diseases including lung cancer and tuberculosis.
5. Silicosis: It is a lung disease caused by inhalation of silica dust, leading to inflammation and scarring of the lung tissue, making the individual more susceptible to tuberculosis.

Prompt 5: Hypothesizing missing variable in open world, Ground truth variable: Visted Asia

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

Prompt: Hello. You will be given a causal graph. The context of the graph is different aspects of a patient's health related to lung diseases. Please understand the causal relationships between the variables - < visited asia > causes < tuberculosis >. < bronchitis > causes < dyspnoea-laboured breathing >. < individual has either tuberculosis or lung cancer > causes < positive xray >. < individual has either tuberculosis or lung cancer > causes < dyspnoea-laboured breathing >. < x > causes < individual has either tuberculosis or lung cancer >. < smoking cigarettes > causes < x >. < smoking cigarettes > causes < bronchitis >. < tuberculosis > causes < individual has either tuberculosis or lung cancer >. Using your causal knowledge, can you suggest what variable X most likely? Give 5 suggestions for what variable X could be. Give reason for your suggestions. Finally, return your answer (without reasoning) in the following format: <Answer> [first suggestion, second suggestion, third suggestion .. and so on] </Answer> No explanations, no reasoning, just answers as a list as described earlier. For example, if the suggestions are A,B,C,D,E then - <Answer> [A, B, C, D, E] </Answer>

Answer: <Answer> [**Lung Cancer**, Chronic Obstructive Pulmonary Disease, Emphysema, Pneumonia, Asthma] </Answer>

Reasons:

1. **Lung Cancer**: Smoking is a major cause of lung cancer and lung cancer can lead to tuberculosis or other lung diseases.
2. Chronic Obstructive Pulmonary Disease: Smoking is a leading cause of COPD and it can lead to other lung diseases including tuberculosis and lung cancer.
3. Emphysema: It is a type of COPD often caused by smoking and can lead to other lung diseases.
4. Pneumonia: Smoking increases the risk of pneumonia and pneumonia can lead to other lung diseases.
5. Asthma: Smoking can trigger asthma attacks and severe asthma can lead to other lung diseases.

Prompt 6: Hypothesizing missing variable in open world, Ground truth variable: Lung cancer

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

Prompt: Hello. You will be given a causal graph. The context of the graph is hypothetical patient monitoring system in an intensive care unit (ICU) Please understand the causal relationships between the variables - < anaphylaxis > causes < total peripheral resistance >. < x > causes < expelled co2 >. < x > causes < x3 >. < x3 > causes < x2 >. < x1 > causes < blood pressure >. < disconnection > causes < x7 >. < error cauter > causes < heart rate displayed on ekg monitor >. < error cauter > causes < oxygen saturation >. < error low output > causes < heart rate blood pressure >. < high concentration of oxygen in the gas mixture > causes < x9 >. < x2 > causes < heart rate blood pressure >. < x2 > causes < heart rate displayed on ekg monitor >. < x2 > causes < oxygen saturation >. < x2 > causes < x1 >. < hypovolemia > causes < left ventricular end-diastolic volume >. < hypovolemia > causes < stroke volume >. < insufficient anesthesia > causes < x3 >. < intubation > causes < x5 >. < intubation > causes < minute volume >. < intubation > causes < x4 >. < intubation > causes < shunt - normal and high >. < intubation > causes < breathing pressure >. < kinked chest tube > causes < x5 >. < kinked chest tube > causes < breathing pressure >. < left ventricular end-diastolic volume > causes < central venous pressure >. < left ventricular end-diastolic volume > causes < pulmonary capillary wedge pressure >. < left ventricular failure > causes < previous medical history >. < left ventricular failure > causes < left ventricular end-diastolic volume >. < left ventricular failure > causes < stroke volume >. < the amount of time using a breathing machine > causes < x6 >. < sudden blockage in the pulmonary arteries > causes < shunt - normal and high >. < sudden blockage in the pulmonary arteries > causes < pulmonary artery pressure >. < x9 > causes < x8 >. < x8 > causes < x3 >. < shunt - normal and high > causes < x8 >. < stroke volume > causes < x1 >. < total peripheral resistance > causes < x3 >. < total peripheral resistance > causes < blood pressure >. < x4 > causes < x >. < x4 > causes < x9 >. < x5 > causes < expelled co2 >. < x5 > causes < minute volume >. < x5 > causes < x4 >. < x6 > causes < x7 >. < x7 > causes < x5 >. < x7 > causes < breathing pressure >. Using your causal knowledge, can you suggest what variable X most likely? Give 5 suggestions for what variable X could be. Give reason for your suggestions. Finally, return your answer (without reasoning) in the following format: <Answer> [first suggestion, second suggestion, third suggestion .. and so on] </Answer> No explanations, no reasoning, just answers as a list as described earlier. For example, if the suggestions are A,B,C,D,E then - <Answer> [A, B, C, D, E] </Answer>

Answer: <Answer> [ventilation rate, respiratory effort, airway resistance, lung compliance, alveolar gas exchange] </Answer>

Reasons for the answer:

1. **Ventilation Rate**: The variable X causes both expelled CO2 and another variable (X3), which indirectly influences heart rate, blood pressure, and oxygen saturation. The rate of ventilation directly affects the amount of CO2 expelled from the body, making it a plausible candidate for X.
2. **Respiratory Effort**: This suggestion is based on the relationship between X and expelled CO2, as well as the influence of X on another variable (X3) that affects cardiovascular and respiratory parameters. Respiratory effort, which encompasses the work of breathing, could affect CO2 expulsion and have downstream effects on the cardiovascular system.
3. **Airway Resistance**: Airway resistance could be a candidate for X because it directly influences the amount of CO2 that can be expelled (a function of airflow) and could also affect other physiological variables through its impact on ventilation efficiency.
4. **Lung Compliance**: Lung compliance, or the ease with which the lungs can expand, directly affects ventilation and, consequently, the expulsion of CO2. It could also indirectly influence other variables in the graph through its impact on gas exchange efficiency.
5. **Alveolar Gas Exchange**: This process directly determines the amount of CO2 expelled and could influence the variable X3, affecting cardiovascular parameters. Efficient gas exchange is crucial for maintaining appropriate oxygen and CO2 levels in the blood, impacting heart rate, blood pressure, and oxygen saturation.

These suggestions are not explicitly mentioned in the input graph but are closely related to the physiological processes and outcomes described. They offer plausible explanations for the causal relationships involving the variable X, considering the context of patient monitoring in an ICU setting.

Prompt 7: Hypothesizing missing variable in open world #1 Ground truth variable: arterial CO2

2106 G ASSUMPTIONS

2107

2108 The **causal sufficiency** of \mathcal{G} , by definition, implies that for every pair of variables within \mathbf{V} , all
 2109 common causes are also included within \mathbf{V} . Extending this assumption to \mathcal{G}^* , we assume that the
 2110 partial graph inherits causal sufficiency for its given that all edges among these variables are preserved
 2111 as in \mathcal{G} . This preservation ensures that the observed relationships within V^* are not confounded
 2112 by omitted common causes. Since the faithfulness of \mathcal{G} ensures that the observed conditional
 2113 independencies among variables in \mathbf{V} are accurately reflected by the causal structure represented
 2114 by \mathbf{E} . By maintaining the same set of edges \mathbf{E} in \mathcal{G}^* for the subset V^* , we uphold the faithfulness
 2115 assumption within the partial graph.

2116

2117 H NDE AND NIE

2118

2119 Average Treatment Effect (ATE) quantifies the expected change in the outcome v_y caused by the unit
 2120 change of the treatment v_t . ATE is part of the causal do-calculus introduced by Pearl (2009). We
 2121 consider binary causal DAGs, i.e., each variable can either take 0 or 1 as values.

2122

$$2123 \text{ATE} = \mathbb{E}[v_y | \text{do}(v_t = 1)] - \mathbb{E}[v_y | \text{do}(v_t = 0)]$$

2124

2125 where the $\text{do}(\cdot)$ operator, represents an intervention. The $E[v_y | \text{do}(v_t = 1)]$ represents the expected
 2126 value of the outcome variable v_y when we intervene to set the treatment variable v_t to 1 (i.e., apply
 2127 the treatment), and $E[v_y | \text{do}(v_t = 0)]$ represents the expected value of v_y when we set v_t to 0 (i.e., do
 not apply the treatment).

2128

2129 H.1 MEDIATION ANALYSIS

2130

2131 Mediation analysis is implemented to quantify the effect of a treatment on the outcome via a third
 2132 variable, the mediator. The total mediation effect can be decomposed into the Natural Direct Effect
 2133 (NDE) and the Natural Indirect Effect (NIE). The Natural Direct Effect (NDE) is the effect of the
 2134 treatment on the outcome variable when not mediated by the mediator variable. The Natural Indirect
 2135 Effect (NIE) is the effect of the treatment variable on the outcome variable when mediated by the
 mediator variable.

2136

$$2137 \text{NDE} = \mathbb{E}[v_{t=1, v_m=0} - v_{t=0, v_m=0}]$$

2138

2139 Here, NDE is calculated by comparing the expected outcome when the treatment variable is set to
 2140 1 and the mediator is fixed at the level it would take under the control treatment $v_t = 0$, with the
 2141 expected outcome when both the treatment and the mediator are set to the control level.

2142

$$2143 \text{NIE} = \mathbb{E}[v_{t=0, v_m=1} - v_{t=0, v_m=0}]$$

2144

2145 Here, NIE is calculated by comparing the expected outcome when the treatment variable is set to 1
 2146 and the mediator is allowed to change as it would under the treatment, with the expected outcome
 when the treatment variable is set to 1 but the mediator is fixed at the control level.

2147

2148

2149

2150

2151

2152

2153

2154

2155

2156

2157

2158

2159

2160 **Prompt:** Strictly follow the format mentioned otherwise you will be disqualified.', 'ello. You will
 2161 be given a causal graph. The context of the graph is hypothetical patient monitoring system in
 2162 an intensive care unit (ICU) Please understand the causal relationships between the variables - <
 2163 anaphylaxis > causes < total peripheral resistance >. < Alveolar Gas Exchange > causes <
 2164 expelled co2 >. < Alveolar Gas Exchange > causes < x2 >. < x2 > causes < x1 >. < x >
 2165 causes < blood pressure >. < disconnection > causes < x6 >. < error cauter > causes < heart rate
 2166 displayed on ekg monitor >. < error cauter > causes < oxygen saturation >. < error low output >
 2167 causes < heart rate blood pressure >. < high concentration of oxygen in the gas mixture > causes
 2168 < x8 >. < x1 > causes < heart rate blood pressure >. < x1 > causes < heart rate displayed on
 2169 ekg monitor >. < x1 > causes < oxygen saturation >. < x1 > causes < x >. < hypovolemia >
 2170 causes < left ventricular end-diastolic volume >. < hypovolemia > causes < stroke volume >. <
 2171 insufficient anesthesia > causes < x2 >. < intubation > causes < x4 >. < intubation > causes
 2172 < minute volume >. < intubation > causes < x3 >. < intubation > causes < shunt - normal and
 2173 high >. < intubation > causes < breathing pressure >. < kinked chest tube > causes < x4 >.
 2174 < kinked chest tube > causes < breathing pressure >. < left ventricular end-diastolic volume >
 2175 causes < central venous pressure >. < left ventricular end-diastolic volume > causes < pulmonary
 2176 capillary wedge pressure >. < left ventricular failure > causes < previous medical history >. <
 2177 left ventricular failure > causes < left ventricular end-diastolic volume >. < left ventricular failure
 2178 > causes < stroke volume >. < the amount of time using a breathing machine > causes < x5 >.
 2179 < sudden blockage in the pulmonary arteries > causes < shunt - normal and high >. < sudden
 2180 blockage in the pulmonary arteries > causes < pulmonary artery pressure >. < x8 > causes < x7 >.
 2181 < x7 > causes < x2 >. < shunt - normal and high > causes < x7 >. < stroke volume > causes
 2182 < x >. < total peripheral resistance > causes < x2 >. < total peripheral resistance > causes <
 2183 blood pressure >. < x3 > causes < Alveolar Gas Exchange >. < x3 > causes < x8 >. < x4 >
 2184 causes < expelled co2 >. < x4 > causes < minute volume >. < x4 > causes < x3 >. < x5 >
 2185 causes < x6 >. < x6 > causes < x4 >. < x6 > causes < breathing pressure >. Using your causal
 2186 knowledge, can you suggest what variable X most likely? Give 5 suggestions for what variable X
 2187 could be. Give reason for your suggestions. Finally, return your answer (without reasoning) in the
 2188 following format: <Answer> [first suggestion, second suggestion, third suggestion .. and so on]
 2189 </Answer> No explanations, no reasoning, just answers as a list as described earlier. For example, if
 2190 the suggestions are A,B,C,D,E then - <Answer> [A, B, C, D, E] </Answer>
 2191 **Answer:** <Answer> [cardiac output, heart function, blood flow, vascular resistance, cardiac cycle]
 2192 </Answer>
 2193 Reasons for the answer:
 2194 1. ****Cardiac Output****: The variable X causes blood pressure and is caused by stroke volume. In
 2195 physiological terms, cardiac output is a primary factor influencing blood pressure, as it represents the
 2196 volume of blood the heart pumps per minute. Since stroke volume (the amount of blood pumped by
 2197 the left ventricle of the heart in one contraction) contributes to cardiac output, it's logical to suggest
 2198 "cardiac output" as a potential identity for X.
 2199 2. ****Heart Function****: Given that X is related to blood pressure and influenced by stroke volume,
 2200 it's reasonable to infer that X pertains to the overall function of the heart. Heart function encompasses
 2201 various aspects of how the heart works, including pumping efficiency, which directly affects blood
 2202 pressure.
 2203 3. ****Blood Flow****: Since X has a direct relationship with blood pressure and an indirect relationship
 2204 with factors affecting the heart's pumping ability, "blood flow" could be a plausible suggestion. Blood
 2205 flow is crucial for maintaining blood pressure and is inherently linked to the heart's ability to pump
 2206 blood effectively.
 2207 4. ****Vascular Resistance****: X's impact on blood pressure and its connection to total peripheral
 2208 resistance suggest that it might involve vascular resistance. Vascular resistance is the resistance that
 2209 must be overcome to push blood through the circulatory system and is a key determinant of blood
 2210 pressure.
 2211 5. ****Cardiac Cycle****: The cardiac cycle involves all the events related to the flow of blood through
 2212 the heart during one complete heartbeat, including systole (contraction) and diastole (relaxation).
 2213 Given X's influence on blood pressure and its relationship with stroke volume, considering the cardiac
 cycle as a potential identity for X makes sense, as it encompasses the mechanisms affecting blood
 pressure regulation.

2214 **Prompt:** Hello. You will be given a causal graph. The context of the graph is hypothetical patient
 2215 monitoring system in an intensive care unit (ICU) Please understand the causal relationships between
 2216 the variables - < anaphylaxis > causes < total peripheral resistance >. < < Alveolar Gas Exchange
 2217 > causes < expelled co2 >. < < Alveolar Gas Exchange > causes < x1 >. < x1 > causes < x >.
 2218 < Cardiac Output > causes < blood pressure >. < disconnection > causes < x5 >. < error cauter
 2219 > causes < heart rate displayed on ekg monitor >. < error cauter > causes < oxygen saturation >.
 2220 < error low output > causes < heart rate blood pressure >. < high concentration of oxygen in the
 2221 gas mixture > causes < x7 >. < x > causes < heart rate blood pressure >. < x > causes < heart
 2222 rate displayed on ekg monitor >. < x > causes < oxygen saturation >. < x > causes < Cardiac
 2223 Output >. < hypovolemia > causes < left ventricular end-diastolic volume >. < hypovolemia >
 2224 causes < stroke volume >. < insufficient anesthesia > causes < x1 >. < intubation > causes <
 2225 x3 >. < intubation > causes < minute volume >. < intubation > causes < x2 >. < intubation >
 2226 causes < shunt - normal and high >. < intubation > causes < breathing pressure >. < kinked chest
 2227 tube > causes < x3 >. < kinked chest tube > causes < breathing pressure >. < left ventricular
 2228 end-diastolic volume > causes < central venous pressure >. < left ventricular end-diastolic volume
 2229 > causes < pulmonary capillary wedge pressure >. < left ventricular failure > causes < previous
 2230 medical history >. < left ventricular failure > causes < left ventricular end-diastolic volume >. <
 2231 left ventricular failure > causes < stroke volume >. < the amount of time using a breathing machine
 2232 > causes < x4 >. < sudden blockage in the pulmonary arteries > causes < shunt - normal and high
 2233 >. < sudden blockage in the pulmonary arteries > causes < pulmonary artery pressure >. < x7
 2234 > causes < x6 >. < x6 > causes < x1 >. < shunt - normal and high > causes < x6 >. < stroke
 2235 volume > causes < Cardiac Output >. < total peripheral resistance > causes < x1 >. < total
 2236 peripheral resistance > causes < blood pressure >. < x2 > causes < < Alveolar Gas Exchange >.
 2237 < x2 > causes < x7 >. < x3 > causes < expelled co2 >. < x3 > causes < minute volume >. <
 2238 x3 > causes < x2 >. < x4 > causes < x5 >. < x5 > causes < x3 >. < x5 > causes < breathing
 2239 pressure > Using your causal knowledge, can you suggest what variable X most likely? Give 5
 2240 suggestions for what variable X could be. Give reason for your suggestions. Finally, return your
 2241 answer (without reasoning) in the following format: <Answer> [first suggestion, second suggestion,
 2242 third suggestion .. and so on] </Answer> No explanations, no reasoning, just answers as a list as
 2243 described earlier. For example, if the suggestions are A,B,C,D,E then - <Answer> [A, B, C, D, E]
 2244 </Answer>

2244 **Answer:** <Answer> [patient anxiety level, patient consciousness level, patient pain level, patient
 2245 respiratory rate, patient skin temperature] </Answer>
 2246 Reasons for the answer:
 2247 1. ****Patient Anxiety Level****: The variable X influences heart rate, blood pressure, oxygen satura-
 2248 tion, and cardiac output, which are all physiological parameters that can be affected by a patient's
 2249 psychological state, such as anxiety. Anxiety can lead to increased heart rate and blood pressure,
 2250 changes in oxygen saturation, and alterations in cardiac output due to the body's stress response.
 2251 2. ****Patient Consciousness Level****: Consciousness level can impact heart rate, blood pressure,
 2252 and oxygen saturation due to the brain's control over these vital signs. Changes in consciousness
 2253 level, such as becoming more alert or lapsing into unconsciousness, can significantly affect these
 2254 parameters, which are all influenced by variable X.
 2255 3. ****Patient Pain Level****: Pain can cause physiological responses that include increased heart rate,
 2256 elevated blood pressure, changes in oxygen saturation due to hyperventilation or altered breathing
 2257 patterns, and potentially affect cardiac output as the body reacts to the stress of pain.
 2258 4. ****Patient Respiratory Rate****: Given that X influences expelled CO2 and is affected by alveolar
 2259 gas exchange, it's plausible that X could relate to the patient's respiratory rate. Respiratory rate
 2260 directly impacts alveolar gas exchange, which in turn affects expelled CO2 levels, oxygen saturation,
 2261 and indirectly influences heart rate and blood pressure through changes in blood chemistry.
 2262 5. ****Patient Skin Temperature****: While not directly related to the cardiovascular or respiratory
 2263 systems, skin temperature can be an indirect indicator of changes in these systems. For example,
 2264 changes in cardiac output and peripheral resistance can affect blood flow to the skin, thereby influenc-
 2265 ing skin temperature. Variable X's impact on cardiac output and blood pressure suggests it could also
 2266 indirectly affect skin temperature.
 2267 These suggestions are based on the understanding that variable X has a broad impact on both
 cardiovascular and respiratory parameters, indicating it could be a factor that influences or reflects
 the overall physiological state of the patient.

2268 **Prompt:** Hello. You will be given a causal graph. The context of the graph is hypothetical patient
 2269 monitoring system in an intensive care unit (ICU) Please understand the causal relationships between
 2270 the variables - < anaphylaxis > causes < total peripheral resistance >. < < Alveolar Gas Exchange
 2271 > causes < expelled co2 >. < < Alveolar Gas Exchange > causes < x >. < x > causes < Patient
 2272 Respiratory Rate >. < Cardiac Output > causes < blood pressure >. < disconnection > causes
 2273 < x4 >. < error cauter > causes < heart rate displayed on ekg monitor >. < error cauter >
 2274 causes < oxygen saturation >. < error low output > causes < heart rate blood pressure >. <
 2275 high concentration of oxygen in the gas mixture > causes < x6 >. < Patient Respiratory Rate >
 2276 causes < heart rate blood pressure >. < Patient Respiratory Rate > causes < heart rate displayed on
 2277 ekg monitor >. < Patient Respiratory Rate > causes < oxygen saturation >. < Patient Respiratory
 2278 Rate > causes < Cardiac Output >. < hypovolemia > causes < left ventricular end-diastolic
 2279 volume >. < hypovolemia > causes < stroke volume >. < insufficient anesthesia > causes < x
 2280 >. < intubation > causes < x2 >. < intubation > causes < minute volume >. < intubation >
 2281 causes < x1 >. < intubation > causes < shunt - normal and high >. < intubation > causes <
 2282 breathing pressure >. < kinked chest tube > causes < x2 >. < kinked chest tube > causes <
 2283 breathing pressure >. < left ventricular end-diastolic volume > causes < central venous pressure >.
 2284 < left ventricular end-diastolic volume > causes < pulmonary capillary wedge pressure >. < left
 2285 ventricular failure > causes < previous medical history >. < left ventricular failure > causes < left
 2286 ventricular end-diastolic volume >. < left ventricular failure > causes < stroke volume >. < the
 2287 amount of time using a breathing machine > causes < x3 >. < sudden blockage in the pulmonary
 2288 arteries > causes < shunt - normal and high >. < sudden blockage in the pulmonary arteries >
 2289 causes < pulmonary artery pressure >. < x6 > causes < x5 >. < x5 > causes < x >. < shunt
 2290 - normal and high > causes < x5 >. < stroke volume > causes < Cardiac Output >. < total
 2291 peripheral resistance > causes < x >. < total peripheral resistance > causes < blood pressure >. <
 2292 x1 > causes < < Alveolar Gas Exchange >. < x2 > causes < x6 >. < x3 > causes < expelled
 2293 co2 >. < x3 > causes < minute volume >. < x3 > causes < x1 >. < x3 > causes < x4 >. <
 2294 x4 > causes < x3 >. < x4 > causes < breathing pressure > Using your causal knowledge, can
 2295 you suggest what variable X most likely? Give 5 suggestions for what variable X could be. Give
 2296 reason for your suggestions. Finally, return your answer (without reasoning) in the following format:
 2297 <Answer> [first suggestion, second suggestion, third suggestion .. and so on] </Answer> No
 2298 explanations, no reasoning, just answers as a list as described earlier. For example, if the suggestions
 2299 are A,B,C,D,E then - <Answer> [A, B, C, D, E] </Answer>

2299 **Answer:** <Answer> [oxygen delivery, arterial blood gas, lung compliance, blood oxygen level,
 2300 respiratory efficiency] </Answer>
 2301 Reasons for the answer:
 2302 1. **Oxygen Delivery**: Given that "alveolar gas exchange" causes "X" and "X" in turn affects
 2303 "patient respiratory rate," it's logical to infer that "X" could be related to how effectively oxygen is
 2304 being delivered throughout the body. Oxygen delivery is a critical aspect of respiratory function and
 2305 would directly impact respiratory rate.
 2306 2. **Arterial Blood Gas (ABG)**: ABG tests measure the levels of oxygen and carbon dioxide in the
 2307 blood and can indicate how well oxygen is being transported and carbon dioxide is being expelled,
 2308 which are functions influenced by alveolar gas exchange. Since "X" is caused by "alveolar gas
 2309 exchange" and affects "patient respiratory rate," it's plausible that "X" could involve the measurement
 2310 or state of arterial blood gases.
 2311 3. **Lung Compliance**: This refers to the elasticity or flexibility of the lungs, which directly
 2312 impacts alveolar gas exchange and, consequently, the respiratory rate. Since "X" is linked to both
 2313 "alveolar gas exchange" and "patient respiratory rate," it could suggest a relationship with the
 2314 mechanical properties of the lung, such as compliance.
 2315 4. **Blood Oxygen Level**: This is a direct outcome of effective alveolar gas exchange and would
 2316 naturally influence the patient's respiratory rate as the body adjusts to maintain adequate oxygenation.
 2317 Given the causal links mentioned, "X" could very well represent the level of oxygen present in the
 2318 blood.
 2319 5. **Respiratory Efficiency**: This encompasses how effectively the respiratory system can supply
 2320 oxygen to and remove carbon dioxide from the body. It's influenced by alveolar gas exchange and
 2321 directly affects the respiratory rate, making it a plausible candidate for "X."
 These suggestions are not explicitly mentioned in the input graph but are logically inferred based on
 the described causal relationships and a basic understanding of respiratory physiology.