

# Causal Graph Discovery with Retrieval-Augmented Generation based Large Language Models

Anonymous ACL submission

## Abstract

Causal graph recovery is traditionally done using statistical estimation-based methods or based on individual’s knowledge about variables of interests. They often suffer from data collection biases and limitations of individuals’ knowledge. The advance of large language models (LLMs) provides opportunities to address these problems. We propose a novel method that leverages LLMs to deduce causal relationships in general causal graph recovery tasks. This method leverages knowledge compressed in LLMs and knowledge LLMs extracted from scientific publication database as well as experiment data about factors of interest to achieve this goal. Our method gives a prompting strategy to extract associational relationships among those factors and a mechanism to perform causality verification for these associations. Comparing to other LLM-based methods that directly instruct LLMs to do the highly complex causal reasoning, our method shows clear advantage on causal graph quality on benchmark datasets. More importantly, as causality among some factors may change as new research results emerge, our method show sensitivity to new evidence in the literature and can provide useful information for updating causal graphs accordingly.

## 1 Introduction

Estimating causal effect between variables from observational data is a fundamental problem to many domains including medical science (Höfler, 2005), social science (Angrist et al., 1996), and economics (Imbens and Rubin, 2015; Yao et al., 2021). It enables reliable decision-making from complex data with entangled associations.

While it is usually expensive and infeasible to investigate causal effects by the golden standard—randomized experiments—researchers employ causal inference (Pearl, 2010) to estimate causal effects from observational data. There are

two main frameworks for causal inference: the potential outcome framework (Rubin, 1974) and the structural causal model (SCM) (Pearl, 1995). Prior causal structures, usually represented as Directed Graphical Causal Models (DGCMS) (Pearl, 2000; Spirtes et al., 2001), are often used to represent and analyze the causal relationships. These causal graphs help disentangle the complex interdependencies and facilitate the analysis of causal effects. Recovering causal graphs often relies on experts’ knowledge or statistical estimation on experimental data (Spirtes and Glymour, 1991). Causal Discovery (CD) algorithms (Spirtes and Glymour, 1991) are the main statistical estimation-based methods that use conditional independence tests to assess associational relationships (called associational reasoning) for inferring causal connections (Spirtes et al., 2001; Chickering, 2002; Shimizu et al., 2006; Sanchez-Romero et al., 2018).

Consequently, the reliability of these algorithms is affected by the quality of data, which can be compromised by issues such as measurement error (Zhang et al., 2017), selection bias (Bareinboim et al., 2014) and unmeasured confounders (Bhattacharya et al., 2021) (See Example A.1 in Appendix A.1). Additionally, CD algorithms often assume certain distribution, such as Gaussian about data, which may fail to accurately reflect the complexity of real-world scenarios.

To mitigate the limitations, Large Language Models (LLMs) (Zhao et al., 2023) have recently been employed for causal graph recovery (Zhou et al., 2023). There are two main streams of these work: 1) directly outputting causal graphs (Choi et al., 2022; Long et al., 2022; Kıcıman et al., 2023); 2) assisting in refining causal graphs generated by statistical estimation-based methods Vashishtha et al. (2023); Ban et al. (2023). Most work have a straightforward way of using LLMs. They directly query the causal relationship between each pair of variables (Choi et al., 2022; Long et al.,

042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082

2022; Kıcıman et al., 2023) by prompting LLMs with the definition of causality, task details and description of the variables of interest. They require LLMs to have extensive domain knowledge and capabilities to perform complex causal reasoning. Whether LLMs have sufficient knowledge in specific domains or whether they have causal reasoning capabilities are questionable (Kandpal et al., 2023; Zečević et al., 2023).

An alternative approach is to exploit LLMs’ capabilities on associational reasoning, e.g., querying the conditional independences (CIs) and recover causal graphs based on extracted associations using CD algorithms (Cohrs et al., 2023). However, it remains difficult for LLMs to understand the CIs between variables, especially the independence conditioned on a large set of variables. (Jiralerspong et al., 2024) tries to inject statistical CI results into LLMs to improve direct causal relationship query results, but the efficacy varies among datasets.

We propose the LLM Assisted Causal Recovery (LACR) method to address these challenges. LACR enhances the knowledge base of LLMs with Retrieval Augmented Generation (RAG) (Lewis et al., 2020; Borgeaud et al., 2022) for reliable associational reasoning. We retrieve highly related knowledge base from a large scientific corpus that contains valuable insight hidden in datasets about associational/causal relationships among variables. We further enhance the accuracy of LACR’s causal recovery results by aggregating the collective extracted information from related literature according to the Wisdom of the Crowd principle (Grofman et al., 1983). LACR also uses an associational reasoning-based causal recovery prompt strategy which elaborately instructs the LLMs the mathematical intuitions behind conditional independence, and builds a surjection from conditional independences extracted by LLMs to causal relationships between variables. LACR is data-driven and does not rely on task-specific knowledge for document retrieval or prompt design. It can serve as a causal graph recovery tool for generic tasks.

Our methodology provides a structured and systematic approach to inferring causal relationships, as it is grounded in a broader evidentiary base and subject to systematic validation. As LACR conducts associational reasoning on a reliable knowledge base, most of which provide evidences based on experimental data analysis, LACR largely overcomes the collection bias problem in statistical estimation-based CD algorithms. We discuss this

in detail in Section 4 by pointing out the causal conflict between the well-known causal discovery results and recent research results extracted by LACR.

### Our Contributions:

- We introduce a novel RAG-based causal graph recovery method that achieves better associational reasoning. The method shows its potential in accurate causal graph construction and overcoming data collection bias issues in traditional methods.
- We design an associational reasoning-based prompting strategy that reduce LLMs’ task complexity to simple associational reasoning to improve the reliability of LLMs’ results. The reliability gain further improve the quality of recovered causal graphs.
- We conduct experiments in several well-known real-world causal graphs and demonstrate the efficacy of LACR. More importantly, based on the scientific evidence returned by our method, we show bias exists in the validation datasets widely used in the CD community, and suggest ways to improve.

## 2 Background

In this section, we introduce the preliminaries of the *directed graphical causal models* (DGCM) and the *causal graph recovery* problem.

### 2.1 Directed Graphical Causal Models

A *Directed Graphical Causal Model* (DGCM) is a tuple  $M = \langle G, P \rangle$ . In the model,  $G = \langle V, E \rangle$  is a Directed Acyclic Graph (DAG), also known as a *causal graph*, where the set of nodes  $V = \{v_1, \dots, v_n\}$  represents random variables (with  $|V| = n$ ), and  $E \subseteq \{(v_i, v_j) \mid v_i, v_j \in V, v_i \neq v_j\}$  is a set of directed edges, also called *causal edges*, that encode *causal relationships*. Let  $\tilde{G} = \langle V, \tilde{E} \rangle$  be the *skeleton* of DAG  $G$ , where each  $(v_i, v_j) \in \tilde{E}$  is an undirected edge, and it indicates that one of  $(v_i, v_j)$  and  $(v_j, v_i)$  is in  $E$ . Let a sequence of distinct nodes  $\ell = (v_{j_1}, v_{j_2}, \dots, v_{j_m})$  denote a *path*, such that for each  $i \in \{1, 2, \dots, m - 1\}$ ,  $(v_{j_i}, v_{j_{i+1}}) \in \tilde{E}$ . A path is a *causal path* from  $v_{j_1}$  to  $v_{j_m}$  if for each  $i \in \{1, 2, \dots, m - 1\}$ ,  $(v_{j_i}, v_{j_{i+1}}) \in E$ . The joint probability distribution of all variables is denoted by  $P$ . Note that we do not consider any variable other than those in  $V$ , that is, we assume there is no so-called latent or exogenous variable.

**Constraints of causal graphs.** A causal graph is subject to a series of constraints on variables’

184 *associational relationships*. Especially, the causal  
185 edges specify the causal relationships between vari-  
186 ables. Given  $(v_i, v_j) \in E$ ,  $v_i$  is a *direct cause* of  $v_j$ .  
187 That is, when holding the other variables constant,  
188 varying the value of  $v_i$  triggers a corresponding  
189 change in the value of  $v_j$ , but not vice versa. This  
190 causal relationship thus entails the associational rela-  
191 tionship between the variables, i.e., their marginal  
192 probability distributions  $P(v_i)$  and  $P(v_j)$  are as-  
193 sociated (or correlated), which does not have the  
194 direction attribute. Notice that two variables can  
195 be associated even though they do not have a direct  
196 causal relationship between each other. Typical  
197 examples are that two variables linked by a causal  
198 path, and two variables pointed to by two causal  
199 paths that have the same starting node (which is  
200 usually called a covariate). The precise constraints  
201 follow an assumption of the causal graph called the  
202 Causal Markov Assumption.

203 **Assumption 2.1** (Causal Markov Assumption). *In*  
204 *any causal graph, each variable is independent of*  
205 *its non-descendants conditioned on its parents in*  
206 *the causal graph.*

207 Therefore, the structure of a causal graph implies  
208 graphical constraints called *d-separation* (Pearl,  
209 2000) that specify a conditional associational re-  
210 lationship between variables. In the rest of this  
211 paper, for any given variable pair  $v_i, v_j \in V$ , we  
212 constantly use  $V'$  to denote an arbitrary subset of  
213  $V \setminus \{v_i, v_j\}$ , unless otherwise specified.

214 **Definition 2.2** (d-separation). *A variable set  $V'$*   
215 *blocks a path  $\ell$  if (i)  $\ell$  contains at least one arrow-*  
216 *emitting variable belonging to  $V'$ , or (ii)  $\ell$  contains*  
217 *at least one collider (variable  $v_i$  is a collider if*  
218  *$(v_{j_{i-1}}, v_{j_i}), (v_{j_{i+1}}, v_{j_i}) \in E$ ) that does not belong*  
219 *to  $V'$  and has no descendant belonging to  $V'$ . If*  
220  *$V'$  blocks all paths from  $v_i$  to  $v_j$ ,  $V'$  is said to*  
221 *d-separate  $v_i$  and  $v_j$ .*

222 If  $V'$  d-separates  $v_i$  and  $v_j$ , then the joint proba-  
223 bility distribution  $P$  encodes that the two variables  
224 are independent conditioned on  $V'$ .

225 Assumption 2.1 is a necessary condition for  
226 the encoding of the associational relationship con-  
227 straints in  $P$ . On the other hand, the following  
228 *faithfulness assumption* is a sufficient condition  
229 that  $P$  encodes such constraints.

230 **Assumption 2.3** (Causal Faithfulness Assumption).  
231 *A joint distribution  $P$  does not encode additional*  
232 *conditional associational relationships other than*  
233 *those consistent with  $G$ 's d-separation information.*  
234 *We call such  $P$  is faithful to  $G$ .*

235 We now formally define the constraints that fol-  
236 low distribution  $P$  faithful to causal graph  $G$ . Let  
237  $\alpha(ij | V') \in \{0, 1\}$  be the *conditional associational*  
238 *relationship* between variables  $v_i, v_j \in V$   
239 conditioned on variable set  $V'$ .  $\alpha(ij | V') = 0$   
240 denotes that  $v_i$  and  $v_j$  are independent conditioned  
241 on  $V'$  according to  $P$ , and  $\alpha(ij | V') = 1$  denotes  
242 associated. We write  $\alpha(ij)$  when  $V' = \emptyset$ .

243 Then, by Assumptions 2.1 and 2.3 and Defini-  
244 tion 2.2, we have that for  $v_i, v_j \in V$ :

- 245 1.  $V'$  d-separates  $v_i$  and  $v_j \implies \alpha(ij | V') = 0$ ;
- 246 2.  $\alpha(ij) = 1$  and  $(v_i, v_j) \notin \bar{E} \implies \exists V'$  s.t.  
247  $\alpha(ij | V') = 0$ ;
- 248 3.  $(v_i, v_j) \in \bar{E} \implies \nexists V'$  s.t.  $\alpha(ij | V') = 0$ .

### 249 3 Methodology

250 We now start to introduce our LLM-based method,  
251 called *large language model assisted causal re-*  
252 *covery* (LACR), that uses a prompt strategy elabo-  
253 rately designed following the process of a statistical  
254 estimation-based CD method, called the *constraint-*  
255 *based causal graph construction* (CCGC). We first  
256 show how CCGC works.

#### 257 3.1 Constraint-based Causal Graph

##### 258 Construction: From Data to Causation

259 Based on Assumptions 2.1 and 2.3, we are able  
260 to partially construct the causal graph  $G$  from a  
261 *knowledge base* KB that is faithful to  $G$  by a statisti-  
262 cal estimation-based method. In a nutshell, KB can  
263 be but not limited to data, the LLM's background  
264 knowledge, and external documents. For more de-  
265 tails, see Section 3.2.1. We take data as the KB in  
266 CCGC. A KB is called faithful to  $G$  if it estimates a  
267 joint distribution that is faithful to  $G$ .

268 The process of CCGC can be divided into  
269 two phases: the *edge existence verification* phase,  
270 which first constructs the skeleton, and the *orienta-*  
271 *tion* phase, which determines the direction of each  
272 undirected edge. LACR only uses the CCGC-based  
273 prompt strategy to conduct the edge existence veri-  
274 fication, and therefore, we only introduce the first  
275 phase of CCGC.

276 For each pair of variables  $v_i, v_j \in V$ , we verify  
277 the existence of the undirected edge in between  
278 (i.e., whether  $(v_i, v_j) \in \bar{E}$  or not) by statistically  
279 testing whether  $v_i$  and  $v_j$  can be d-separated by  
280 any variable set  $V'$ . Let  $\hat{\alpha}_{\text{KB}}(ij | V') \in \{0, 1\}$  be  
281 an estimator of  $\alpha(ij | V')$ , based on KB. Next,  
282 based on a given KB that is faithful to  $G$ , we  
283 define  $\zeta_{\text{KB}} : V \times V \rightarrow \{0, 1\}$  as the *causal*

*edge existence* mapping, such that  $\zeta_{\text{KB}}(ij) = 0$  if  $\exists V'$  s.t.  $\hat{\alpha}_{\text{KB}}(ij | V') = 0$ , otherwise  $\zeta_{\text{KB}}(ij) = 1$ .  $\zeta_{\text{KB}}(ij) = 0$  implies that we estimate there is no edge between  $v_i$  and  $v_j$ , since the pair of variables can be d-separated by at least one variable set. See Appendix A.1 for an example of CCGC’s process.

Compared with LACR, most existing LLM-based causal graph construction methods directly query LLMs the causal relationships. With the introduction of CCGC, we next illustrate the limited reliability of such methods.

**Limited Reliability of Direct Causal Prompt.** We name the prompt used in such direct query of causal relationships as the *direct causal prompt*. Examples include “Is A a cause of B?” and “Does the change of A cause the change of B”, which are wildly used in related work (Kıcıman et al., 2023; Choi et al., 2022; Long et al., 2022). Such prompt directly queries the causal edge existence ( $\zeta_{\text{KB}}(\cdot)$ ) and the causal direction. We argue that such direct prompting requires extensive causal reasoning capability from LLMs. The following proposition (see proof in Appendix C.1) shows the high complexity hidden behind a direct causal prompt.

**Proposition 3.1.** *Assuming that estimating  $\hat{\alpha}_{\text{KB}}(ij | V')$  for a given  $V'$  needs  $O(1)$  time, inferring  $\zeta_{\text{KB}}(ij)$  requires  $O(2^{n-2})$ , where  $n = |V|$ .*

*Proof.* The proof is illustrated in Section C.1.  $\square$

We now start formally introducing the *large language model assisted causal recovery* (LACR) method, which first extracts the conditional associational relationships between variables, and determines the causal relationships following the process of CCGC (see Section 3.1). We implement such a process by a series of separated queries using the constraint-based causal prompt. The LACR consists of two steps: the edge existence verification (LACR 1) and orientation (LACR 2).

### 3.2 LACR 1: Edge Existence Verification

In this phase, we construct the skeleton of the causal graph, i.e., verifying the existence of each edge without clarifying its direction. We use LLMs to mine the statistical evidence to verify the conditional associational relationship between each pair of variables and determine the existence of a causal edge (recall Section 3.1), from the retrieved scientific documents (corresponding to document-based query), LLMs’ internal knowledge (corre-

sponding to background-based query), and statistical estimation-based output. To achieve this target, we design a prompt strategy that encodes the statistical principles of CCGC.

#### 3.2.1 Constraint-Based Causal (CC) Prompt

In LACR, KB can be the LLM’s background knowledge, external documents, and datasets. For each variable pair, namely  $v_i$  and  $v_j$ , we clarify their conditional associational relationship by mining the statistical evidence from KB. We conduct a chain of 4 queries to determine the final opinion of each piece of KB, namely, the **background reminder**, the **association verifier**, the **association type verifier**, and the **association rechecker**. However, if any KB does not contain sufficient information to determine the value of  $\hat{\alpha}_{\text{KB}}(ij | V')$ , we ask the LLM to give an answer UNKNOWN. We classify such knowledge bases as *unusable*, and they are discarded during the decision-making phase. See the original prompt in Appendix C.6.

**Background reminder.** This prompt component helps the LLM to understand the full picture of the task, and avoid misinterpretation of variables’ meaning. We aim to provide minimum external information about the task other than the names of the variables to the LLM. Therefore, we only give the full FACTOR list (i.e., the names of the variables), and specify the DOMAINS from which the variables are. For example, in the ASIA experiment dataset (see Section 4), all variables are from the domains of MEDICAL, BIOLOGY, AND SOCIAL SCIENCE. Finally, we ask the LLM to specify the meaning of each variable, as well as the interaction among them.

**Association verifier.** The component utilizes KR to verify the zero-order associational relationship between  $v_i$  and  $v_j$ , i.e.,  $\hat{\alpha}_{\text{KB}}(ij)$ . The LLM is provided with an ASSOCIATION CONTEXT (an instruction of how to determine whether  $v_i$  and  $v_j$  is associated or not) and KB. Then, the LLM determines the relationship as ASSOCIATED (if  $\hat{\alpha}_{\text{KB}}(ij) = 1$ ), INDEPENDENT (if  $\hat{\alpha}_{\text{KB}}(ij) = 0$ ), or UNKNOWN, based on the statistical evidence extracted from KB. If the decision is an association, the LACR goes to the next query.

**Association type verifier.** Upon determining ASSOCIATED between  $v_i$  and  $v_j$ , we further need to determine whether this association is “indirect” or “direct”, i.e., whether there exists  $V'$  that can d-

separate  $v_i$  and  $v_j$ . Based on the given KB and reasoning, the LLM is asked to read an ACCUSATION TYPE CONTEXT (an instruction of how to judge whether the association is indirect or direct based on KB). Intuitively, the ACCUSATION TYPE CONTEXT illustrates that if the association between  $v_i$  and  $v_j$  is mediated by variables from  $V'$ , then, the association is indirect, otherwise it is direct. To align precisely with CCGC, the ACCUSATION TYPE CONTEXT further explains that “the association mediated by third variables” means that the association is eliminated if we control the third variables constantly. The LACR goes to the final query if the decision is an INDIRECTLY ASSOCIATED.

**Association rechecker.** Considering the potential that the LLM can return INDIRECTLY ASSOCIATED because it judges that the association between  $v_i$  and  $v_j$  is mediated by external variables that are not from  $V$ . Since we do not consider external variables, we ask the LLM to verify whether the set of mediating variables includes any from  $V \setminus \{v_i, v_j\}$ . If yes, the association type should be corrected to DIRECTLY ASSOCIATED.

### 3.2.2 The CC Prompt is Deterministic

Using the above prompt strategy, we demonstrate that the LLM’s return can determine the existence of causal edges based on a given KB. We first specify that the LLM’s return based on the above prompt must be one from set {INDEPENDENT, DIRECTLY ASSOCIATED, INDIRECTLY ASSOCIATED, UNKNOWN}. If the return is UNKNOWN, the KB is unusable. Then, for each *usable* KB, we have the following proposition (see proof in Appendix C.2).

**Proposition 3.2.** *For each variable pair  $v_i$  and  $v_j$ , the mapping from the conditional associational relationship space of  $v_i$  and  $v_j$  to the return set of each usable KB is a surjection, and the mapping from the return set of each usable KB to the range of  $\zeta_{ij}$ , i.e.,  $\{0, 1\}$ , is also a surjection.*

*Proof.* The proof is illustrated in Section C.2.  $\square$

### 3.2.3 LACR 1

With the above CC prompt, we are ready to introduce LACR 1 (Algorithm 1). We initialize the algorithm by setting the skeleton graph  $\bar{G}$  as a complete undirected graph  $\bar{G}^c$ , giving each variable pair  $v_i, v_j$  a pre-retrieved set of  $k$  relevant scientific documents as the document-based knowledge base  $\mathbf{DOC} = \{\mathbf{DOC}_{ij} = \{\mathbf{DOC}_{ij}^1, \dots, \mathbf{DOC}_{ij}^k\}\}_{v_i, v_j \in V, \text{ s.t., } v_i \neq v_j}$ . Then, for

---

#### Algorithm 1 LACR 1

---

```

1: Input:  $\bar{G} \leftarrow \bar{G}^c$ ,  $\mathbf{DOC}$ ,  $\mathbf{D}$ 
2: for  $\forall v_i, v_j \in V$ , s.t.,  $v_i \neq v_j$  do
3:    $S = 0$ 
4:   for  $\mathbf{KB} \in \mathbf{DOC}_{ij} \cup \{\mathbf{BG}\}$  do
5:     if  $\hat{\zeta}_{\mathbf{KB}}(ij) = 1$  then
6:        $S+ = 1$ 
7:     else if  $\hat{\zeta}_{\mathbf{KB}}(ij) = 0$  then
8:        $S+ = -1$ 
9:     if  $S \leq 0$  then
10:       $\bar{G} \leftarrow \bar{G} \setminus (v_i, v_j)$ 
11: Return:  $\bar{G}$ 

```

---

each variable pair, we query the LLM by the CC prompt to estimate  $\hat{\zeta}_{\mathbf{KB}}(ij)$  based on each of the given documents provided in  $\mathbf{DOC}_{ij}$  and the LLM’s background knowledge BG. If the decision of the LLM is INDIRECTLY ASSOCIATED or INDEPENDENT, i.e.,  $\hat{\zeta}_{\mathbf{KB}}(ij) = 0$ , by Proposition 3.2, we add -1 point to the score  $S$ , if the decision is DIRECTLY ASSOCIATED (i.e.,  $\hat{\zeta}_{\mathbf{KB}}(ij) = 1$ ), we add 1 point to  $S$ , otherwise, we do not change  $S$  if LLM answers UNKNOWN based on KB. After considering all of the LLM’s decisions for  $v_i$  and  $v_j$ , if the final score  $S > 0$ , we keep the undirected edge  $(v_i, v_j)$ ; otherwise, we remove it from  $\bar{G}$ . Finally, the algorithm returns the skeleton after querying each variable pair based on all KBs.

Note that using the score  $S$  to aggregate each KB’s “opinion” for each variable pair is equivalent to making the collective decision of  $\hat{\zeta}(ij)$  by the simple majority voting rule (Brandt et al., 2016). We slightly bias the decision towards  $\hat{\zeta}_{\mathbf{KB}}(ij) = 0$  by the setting of removing an edge if  $S \leq 0$ , since generally, the LLM’s decision biases towards DIRECTLY ASSOCIATED. We use this biased setting because (1) almost KBs cannot load the evidence showing  $\hat{\alpha}_{\mathbf{KB}}(ij | V')$  for all possible  $V'$ , and (2) if most retrieved documents are unusable (no research report on  $v_i$  and  $v_j$ ’s association), then, it is more possible that  $v_i$  and  $v_j$  are not associated. By the theory of the Wisdom of the Crowd, LACR’s decision tends to be more accurate than querying a single knowledge base, and it can be improved by adding more relevant documents (see a detailed description in Appendix B).

### 3.3 LACR 2: Orientation

Starting at the skeleton output by LACR 1, we continue to determine the direction of each edge in the skeleton. In LACR 2, we simply utilize direct query

to LLM for the orientation task due to LLMs’ high performance on causal orientation tasks (Kıcıman et al., 2023). For each pair of adjacent variables in the skeleton, we use a two-step prompt strategy:

**Background reminder.** Similar to LACR 1, we provide the main variables and the domain information of the task, and ask the LLM to clarify the variables’ meanings, as well as their interaction.

**Orienting.** With the above clarification, we ask the LLM to thoroughly understand the given KB and a CAUSAL DIRECTION CONTEXT, that specifies that if variable  $A$  is the cause of variable  $B$ , then, the change of  $A$ ’s value causes a change of  $B$ ’s value, but not vice versa. Then, we ask the LLM to give its decision based on all of the above information.

## 4 Experiments

In this section, we first introduce the ground truth datasets and how we collect three research literature pools. Then we introduce the settings of our solution and baselines. Finally, we evaluate the pruning and orienting results, respectively.

### 4.1 Experiment Data

**Validation datasets.** We validate our method on four datasets (namely, ASIA, SACHS, and CORONARY). All datasets have reported causal graphs (see Appendix C.4) based on real-world data. It is worth noting that, we only limit the selection of validation datasets to real-world datasets because LACR uses a realistic knowledge base.

**ASIA (Iau, 1988).** The ASIA dataset has 8 nodes (from domains of medical, biology, and social science) and 8 edges, revealing the potential reasons and symptoms of lung diseases.

**SACHS (Sachs et al., 2005).** The SACHS dataset has 11 nodes (from the medical and biological domains) and 16 edges. It uncovers the interaction among proteins related to several human diseases.

**CORONARY (Reinis et al., 1981).** The CORONARY dataset has 6 nodes (from the medical and biological domains) and 9 *undirected* edges, revealing the causal relationship among several potential reasons of coronary heart disease. We only use it to validate LACR 1 because the edges are undirected.

### 4.2 Experimental Settings

We use GPT-4o in the following experiments.

**Research document pool construction.** In our experiment, we automatically build the pre-retrieved document set for each variable pair (**Initialization** in Algorithm 1) in two steps:

(1) Relevant paper search: We search 20 paper titles by querying “name[ $v_i$ ] and name[ $v_j$ ]” to the Google Scholar engine using the SerpApi (SerpApi), and rank the papers by Google Scholar’s default relevance ranking.

(2) Paper download: Based on the aforementioned ranked paper title list, we use the PubMed API<sup>1</sup> to download the papers. For each paper title, we prioritize downloading the full document from the PubMed Central (PMC) database, and only download the abstract document from the PubMed database if the full version is not available in PMC. For each variable pair, we download up to 10 documents from the top of the ranked title list (note that some papers are unavailable in PubMed).

**Statistical causal discovery method.** In the validation of LACR 1, additional to LLM, we also test the impact of injecting statistical estimation-based results into the decision-making phase. That is, adding point 1 (resp.  $-1$ ) to score  $S$  if the statistical estimation-based method determines  $\hat{\zeta}_{KB}(ij) = 1$  (resp.  $\hat{\zeta}_{KB}(ij) = 0$ ) in Algorithm 1, where KB is numerical data. We use the Peter-Clark (PC)(Spirtes et al., 2001) algorithm as the statistical estimation-based method. We import the data from the bnlearn package (Scutari et al., 2019).

**Baseline methods.** We survey recent LLM-based causal graph construction methods, and for each dataset, we select the baseline method with the best performance. For each dataset, we present two types of baseline LLMs: baseline LLM1, which is a pure LLM-based method, and baseline LLM2, which is a hybrid method combining a statistical estimation-based and an LLM-based method. We do not compare LACR to any baseline method on the CORONARY dataset as the dataset’s absence in such methods’ validation.

**Validation metrics.** We measure LACR 1 and LACR 2 by different metrics. For LACR 1, we show the adjacency precision (AP), the adjacency recall (AR), the F1 score, and the Normalized Hamming Distance (NHD), as follows.

First, we count three attributes of each graph: true positive (TP): the number of edges that are successfully recovered, false positive (FP): the number of edges that are recovered but different from the ground truth graph, and false negative (FN): the number of edges that exist in the ground truth but not recovered in our constructed graph. Then, we compute AP:  $\frac{TP}{TP+FP}$ , AR:  $\frac{TP}{TP+FN}$ , F1:  $\frac{2AP*AR}{AP+AR}$ , and

<sup>1</sup><https://www.ncbi.nlm.nih.gov/home/develop/api/>

	Dataset	AP	AR	F1	SHD
ASIA	LACR 1 (BG)	<b>1</b>	<b>1</b>	<b>1</b>	0
	LACR 1 (DOC)	0.571	<b>1</b>	0.727	0.122
	LACR 1 (PC)	<b>1</b>	0.75	0.857	0.041
	Baseline LLM1	<b>1</b>	0.88	0.93	0.016
	Baseline LLM2	0.8	<b>1</b>	0.89	0.031
CORO	LACR 1 (BG)	0.625	0.625	0.625	0.167
	LACR 1 (DOC)	0.667	0.75	0.706	0.139
	LACR 1 (PC)	<b>0.778</b>	<b>0.875</b>	<b>0.824</b>	0.083
SACHS	LACR 1 (BG)	<b>0.8</b>	0.5	0.615	0.083
	LACR 1 (DOC)	0.467	<b>0.875</b>	<b>0.609</b>	0.149
	LACR 1 (PC)	0.421	0.5	0.457	0.157
	Baseline LLM1	N/A	N/A	0.31	0.63
	Baseline LLM2	0.59	N/A	0.56	0.12

Table 1: Performances of our solution LACR 1 with different KB. We test the performance across three datasets, and compare to baseline methods: ASIA: LLM1: (Jiralerspong et al., 2024), LLM2: (Jiralerspong et al., 2024), SACHS: LLM1: (Zhou et al., 2024), LLM2: (Takayama et al., 2024).

NHD:  $\frac{FP+FN}{n^2}$ , where  $n$  is the number of variables. Intuitively, NHD is the number different edges between two graphs, normalized by  $n^2$ .

In the validation of LACR 2, we simply compute the True Edge Accuracy (TEA), i.e., the ratio of correctly oriented edges among all true positive edges in LACR 1’s output skeleton.

### 4.3 Evaluation

We now first present observations based on experimental results for three datasets, which contains Edge Existence Verification (Section 4.3.1) and Orientation (Section 4.3.2), followed by a comprehensive analysis of the overall results (Section 4.3.3).

#### 4.3.1 Observation on Edge Existence Verification

We present the performance of LACR 1 on causation existence verification with different knowledge bases KB in the section. The orienting performance will be shown in the next section. Table 1 lists the performance of all compared methods, where BG denotes only LLM’s background knowledge, DOC denotes both LLM’s background knowledge and the fixed number of documents, and PC denotes DOC plus the results output by the PC algorithm. We have the following observations:

**ASIA.** We have three observations from the experimental results on the ASIA dataset. First, LACR 1 achieves the best performance when relying solely on BG. It successfully recovers the full skeleton and outperforms the high performance of the pure LLM method in (Jiralerspong et al., 2024). Second, adding retrieved documents into KB reduces performance (AP from 1 to 0.57, and F1 score from 1 to

0.73) according to the given ground truth in (lau, 1988). Third, by further aggregating the output of the PC algorithm, the F1 score increases from 0.73 to 0.86 compared to the ground truth in (lau, 1988). **CORONARY (CORO).** The results differ notably from those based on the ASIA dataset, LACR 1 with only the LLM’s background knowledge achieves the worst performance, with values of 0.625 for all of AP, AR, and F1 scores. By adding documents and the PC algorithm into KB, all metrics increase, reaching 0.875.

**SACHS.** We have three observations from the results on the SACHS dataset. First, the best performance of LACR 1 is achieved using only the LLM’s background knowledge, outperforming the baseline method (combined method of LLM and hybrid statistical methods of DirectLINGAM). Second, adding documents as the knowledge base slightly decreases the F1 score by 0.01. Third, with the PC algorithm’s output, the F1 score reduces to 0.46, which is worse than the baseline method.

#### 4.3.2 LACR 2: Orientation

The results (Table 3 in Appendix C.5) of LACR 2 show TEA is 1 (i.e., orienting all TP edges correctly) on both ASIA and SACHS, based on all KB. We can observe that, upon successfully recovering causal edges by LACR 1, the orientation accuracy is high, reaching 1 for all knowledge bases and all datasets. It demonstrates the efficacy of the orientation prompt as well as LLM’s capability for causal orientation reasoning. We conjecture that the success of this task strongly depends on the rich evidence stored in the scientific literature, and the easy understandability of such evidence, compared to the extraction of associational relationship.

#### 4.3.3 Overall Results Analysis

By summarizing the overall performance of LACR, it is worth noticing the following points:

**LACR’s performance tends to monotonically increase by taking more KB with high quality and readability.** Observing LACR’s performance on ASIA and CORONARY, we notice that our methods generally perform better with high-quality and readability of the input documents and statistical results. While we discuss the performance drop of LACR1 (DOC) on the ASIA dataset later, the overall trend coincides with the Condorcet theorem (Grofman et al., 1983) in voting theory, which suggests that aggregating diverse, high-quality inputs leads to better outcomes.

**LACR performs differently on tail and non-tail data.** We observe that LACR performs better on ASIA and CORONARY datasets compared to the SACHS dataset. This is because the terms in ASIA and CORONARY are more common to LLMs during training. In contrast, SACHS mainly contains symbols with specific meanings in a specific area. Despite feeding scientific documents to LACR on all datasets, the lack of prior knowledge of these symbols in the training phase limits LLMs’ understanding of their meanings, resulting in hallucinations. This has been observed in RAG-based legal research tools (Magesh et al., 2024).

**Updating on the current ground truth causal graphs is necessary.** We found, through LACR’s responses, strong evidence from domain research (see details in Appendix C.3) indicating that an update to the ground truth causal graphs is necessary. For example, on the ASIA dataset, the ground truth being outdated led to reduced performance when using additional documents. Similarly, for the CORONARY dataset, the improvement in performance with added documents and the PC algorithm suggests that the ground truth for CORONARY is more current compared to ASIA.

#### 4.4 Refining the Ground Truth (ASIA and CORONARY)

The ground truth causal graph needs refinement because it is outdated and significantly differs from current SOTA domain knowledge. Additionally, we aim to determine if the LLM can identify new causal relationships based on SOTA literature. In this part, we present the observations and analyses on the refined ground truth causal graphs based on the refined ASIA and CORONARY datasets. We will discuss how refined graph truth affects performance from three perspectives: causal inferring based on LLM’S background knowledge BG, external literature DOC, and statistic data PC. Due to the page limitation, the details of refining are illustrated in Section C.3.

**Background Knowledge BG.** The performance relying solely on background knowledge BG varies between the two datasets. In the ASIAN dataset, all results slightly drop down, whereas in the dataset, results improve. A possible reason is that the background knowledge BG related to ASIA is outdated, while the knowledge BG for CORONARY is more current. Consequently, when the ground truth is updated based on new domain-specific knowledge, the outcomes are different significantly between

	Dataset	AP	AR	F1	SHD
ASIA	LACR 1 (BG)	1	0.8	0.889	0.041
	LACR 1 (DOC)	0.714	1	0.833	0.082
	LACR 1 (PC)	1	0.6	0.75	0.082
CORO	LACR 1 (BG)	0.75	0.75	0.75	0.111
	LACR 1 (DOC)	0.778	0.875	0.824	0.083
	LACR 1 (PC)	0.667	0.75	0.706	0.139

Table 2: Performances of LACR 1 on the refined ground truth with different KB, comparing to baseline LLM-powered methods.

the datasets.

**External Literature DOC.** Incorporating DOC significantly increases performance for both datasets, as the refined ground truth better aligns with SOTA research trends, demonstrating the importance of up-to-date and relevant domain knowledge in improving model accuracy.

**Statistic Data PC.** Different from the results of incorporating DOC, adding PC’s results consistently worsens performance, highlighting the PC-based solution is using an outdated dataset compared to updated SOTA knowledge. As the ground truth evolves to reflect current research advancements, the relative performance of PC-based results diminishes. These findings emphasize the importance of up-to-date domain-specific knowledge for accurate causal graph recovery.

## 5 Conclusion

In this paper, we proposed a novel LLM-based causal graph construction method called LACR which uses the constraint-based causal prompt strategy designed according to the constraint-based causal graph construction (CCGC) method. Comparing to most existing LLM-based causal graph construction methods, that use the direct causal prompt to query LLMs to do highly complex causal reasoning, LACR mainly relies on LLMs to do low-complexity associational reasoning, and follows the process of CCGC to determine the causal relationships. For accurate associational reasoning, we utilize LLMs’ RAG feature to extract statistical evidence with high relevance and quality from a large scientific corpus. Lastly, we validate LACR’s efficacy on several well-known datasets and show LACR’s outstanding performance among LLM-based methods. More importantly, LACR’s responses show the conflict between the ground truths and SOTA domain research, which requests a refinement of the validation ground truths.



## 742 Limitations

743 We first address three technical limitations of the  
744 current version of LACR. The first is the paper  
745 search accuracy. The pre-retrieved document set  
746 needs high quality and relevance to provide rele-  
747 vant evidence. Therefore, we conjecture that using  
748 refined queries and other search engines can en-  
749 hance the performance. The second limitation is  
750 LLMs’ understanding on highly professional doc-  
751 uments. Through our experiments, we found that  
752 LLMs’ poor comprehension capability on specific  
753 domains, e.g., the SACHS dataset, limits LACR’s  
754 performance. An optional solution is to fine-tune  
755 LLMs to better understand such documents. The  
756 third is the complexity of LACR. The method  
757 needs to query each variable pair ( $O(n^2)$ ), and for  
758 each variable pair, multiple documents need to be  
759 queried.

760 We then address other practical limitations.  
761 What comes first is the need of up-to-data prac-  
762 tical validation datasets and causal graphs in causal  
763 discovery community. Many validation datasets  
764 are synthesized, which are not usable in such prac-  
765 tical knowledge-based methods. The second practi-  
766 cal limitation is the access of scientific papers. In  
767 our experiment, we focus on biomedical datasets  
768 for the accessibility of research papers in PubMed,  
769 however, the full contexts of most of the papers  
770 are not open accessible. It would open the possi-  
771 bility of overall better understanding of causal re-  
772 lationships if full documents are accessible in more  
773 research domains and broader scientific databases.

## 774 References

775 1988. Local computations with probabilities on graphi-  
776 cal structures and their application to expert systems.  
777 *Journal of the Royal Statistical Society: Series B*  
778 (*Methodological*), 50(2):157–194.

779 Roya Alavi-Naini, Batool Sharifi-Mood, and Maliheh  
780 Metanat. 2012. Association between tuberculosis  
781 and smoking. *International journal of high risk be-*  
782 *haviors & addiction*, 1(2):71.

783 Genet A Amere, Pratibha Nayak, Argita D Salindri,  
784 KM Venkat Narayan, and Matthew J Magee. 2018.  
785 Contribution of smoking to tuberculosis incidence  
786 and mortality in high-tuberculosis-burden countries.  
787 *American journal of epidemiology*, 187(9):1846–  
788 1855.

789 Joshua D Angrist, Guido W Imbens, and Donald B  
790 Rubin. 1996. Identification of causal effects using  
791 instrumental variables. *Journal of the American sta-*  
792 *tistical Association*, 91(434):444–455.

Taiyu Ban, Lyvzhou Chen, Xiangyu Wang, and Huan-  
huan Chen. 2023. [From query tools to causal ar-  
chitects: Harnessing large language models for ad-  
vanced causal discovery from data.](#) 793  
794  
795  
796

Elias Bareinboim, Jin Tian, and Judea Pearl. 2014. [Re-  
covering from selection bias in causal and statistical  
inference.](#) *Proceedings of the AAAI Conference on*  
*Artificial Intelligence*, 28(1). 797  
798  
799  
800

ELIZABETH Barrett-Connor and K-T Khaw. 1984. 801  
Family history of heart attack as an independent pre- 802  
dictor of death due to cardiovascular disease. *Circu-* 803  
*lation*, 69(6):1065–1069. 804

Rohit Bhattacharya, Tushar Nagarajan, Daniel Malinsky, 805  
and Ilya Shpitser. 2021. [Differentiable causal discov-](#) 806  
[ery under unmeasured confounding.](#) In *The 24th* 807  
*International Conference on Artificial Intelligence* 808  
*and Statistics, AISTATS 2021, April 13-15, 2021, Vir-* 809  
*tual Event*, volume 130 of *Proceedings of Machine* 810  
*Learning Research*, pages 2314–2322. PMLR. 811

Sebastian Borgeaud, Arthur Mensch, Jordan Hoff- 812  
mann, Trevor Cai, Eliza Rutherford, Katie Mill- 813  
ican, George Bm Van Den Driessche, Jean-Baptiste 814  
Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. 815  
Improving language models by retrieving from tril- 816  
lions of tokens. In *International conference on ma-* 817  
*chine learning*, pages 2206–2240. PMLR. 818

Felix Brandt, Vincent Conitzer, Ulle Endriss, Jérôme 819  
Lang, and Ariel D Procaccia. 2016. *Handbook of* 820  
*computational social choice.* Cambridge University 821  
Press. 822

Kai-Chi Chen, Hong-Ren Yu, Wei-Shiang Chen, Wei- 823  
Che Lin, Yi-Chen Lee, Hung-Hsun Chen, Jyun-Hong 824  
Jiang, Ting-Yi Su, Chang-Ku Tsai, Ti-An Tsai, et al. 825  
2020. Diagnosis of common pulmonary diseases in 826  
children by x-ray images and deep learning. *Scien-* 827  
*tific Reports*, 10(1):17374. 828

David Maxwell Chickering. 2002. Optimal structure 829  
identification with greedy search. *Journal of machine* 830  
*learning research*, 3(Nov):507–554. 831

Kristy Choi, Chris Cundy, Sanjari Srivastava, and 832  
Stefano Ermon. 2022. Lmpriors: Pre-trained lan- 833  
guage models as task-specific priors. *arXiv preprint* 834  
*arXiv:2210.12530.* 835

Kai-Hendrik Cohrs, Emiliano Diaz, Vasileios Sitokon- 836  
stantinou, Gherardo Varando, and Gustau Camps- 837  
Valls. 2023. Large language models for constrained- 838  
based causal discovery. In *AAAI 2024 Workshop* 839  
*on "Are Large Language Models Simply Causal Par-* 840  
*rots?"*. 841

David Edwards. 2000. *Introduction to graphical mod-* 842  
*elling.* Springer Science & Business Media. 843

Bernard N. Grofman, Guillermo Owen, and Scott L. 844  
Feld. 1983. Thirteen theorems in search of the truth. 845  
*Theory and Decision*, 15:261–278. 846

847	Himanshu Gupta, Sanjeev Mahajan, Mohan Lal, Adarshjot Kaur Toor, Shyam Sunder Deepti, and Naresh Chawla. 2022. Prevalence of tobacco consumption and smoking and its effect on outcome among microbiologically confirmed new pulmonary tuberculosis patients on daily regimen of dots in amritsar city. <i>Journal of Family Medicine and Primary Care</i> , 11(5):2150–2156.	903
848		904
849		
850		
851		
852		
853		
854		
855	T Hintsu, M Shipley, D Gimeno, M Elovainio, T Chandola, M Jokela, L Keltikangas-Järvinen, J Vahtera, MG Marmot, and M Kivimäki. 2010. Do family history of chd, education, paternal social class, number of siblings and height explain the association between psychosocial factors at work and coronary heart disease? the whitehall ii study. <i>Occupational and environmental medicine</i> , 67(5):330.	
856		
857		
858		
859		
860		
861		
862		
863	Marc Höfler. 2005. Causal inference based on counterfactuals. <i>BMC medical research methodology</i> , 5(1):1–12.	
864		
865		
866	David J Horne, Monica Campo, Justin R Ortiz, Eyal Oren, Matthew Arentz, Kristina Crothers, and Masahiro Narita. 2012. Association between smoking and latent tuberculosis in the us population: an analysis of the national health and nutrition examination survey. <i>PLoS one</i> , 7(11):e49050.	
867		
868		
869		
870		
871		
872	Guido W. Imbens and Donald B. Rubin. 2015. <i>Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction</i> . Cambridge University Press.	
873		
874		
875	Xuefeng Jin, Caiyun Zhang, Chao Chen, Xiaoning Wang, Jing Dong, Yuanyuan He, and Peng Zhang. 2023. Tropheryma whipplei-induced plastic bronchitis in children: a case report. <i>Frontiers in Pediatrics</i> , 11:1185519.	
876		
877		
878		
879		
880	Thomas Jiralerspong, Xiaoyin Chen, Yash More, Vedant Shah, and Yoshua Bengio. 2024. <a href="#">Efficient causal graph discovery using large language models</a> .	
881		
882		
883	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In <i>International Conference on Machine Learning</i> , pages 15696–15707. PMLR.	
884		
885		
886		
887		
888	Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. 2023. Causal reasoning and large language models: Opening a new frontier for causality. <i>arXiv preprint arXiv:2305.00050</i> .	
889		
890		
891		
892	Seung Hoon Kim, Yong-Moon Park, Kyungdo Han, Seung Hyun Ko, Shin Young Kim, So Hyang Song, Chi Hong Kim, Kyu Yeon Hur, and Sung Kyoung Kim. 2022. Association of weight change following smoking cessation with the risk of tuberculosis development: A nationwide population-based cohort study. <i>Plos one</i> , 17(4):e0266262.	
893		
894		
895		
896		
897		
898		
899	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation	
900		
901		
902		
	for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33:9459–9474.	
	Ryan P Lindsay, Sanghyuk S Shin, Richard S Garfein, Melanie LA Rusch, and Thomas E Novotny. 2014. The association between active and passive smoking and latent tuberculosis infection in adults and children in the united states: results from nhanes. <i>PLoS one</i> , 9(3):e93137.	905
		906
		907
		908
		909
		910
	Stephanie Long, Tibor Schuster, and Alexandre Piché. 2022. Can large language models build causal graphs? In <i>NeurIPS 2022 Workshop on Causality for Real-world Impact</i> .	911
		912
		913
		914
	Varun Magesh, Faiz Surani, Matthew Dahl, Mirac Suzgun, Christopher D Manning, and Daniel E Ho. 2024. Hallucination-free? assessing the reliability of leading ai legal research tools. <i>arXiv preprint arXiv:2405.20362</i> .	915
		916
		917
		918
		919
	Prince Ntiamoah, Sanjay Mukhopadhyay, Subha Ghosh, and Atul C Mehta. 2021. Recycling plastic: diagnosis and management of plastic bronchitis among adults. <i>European Respiratory Review</i> , 30(161).	920
		921
		922
		923
	Judea Pearl. 1995. Causal diagrams for empirical research. <i>Biometrika</i> , 82(4):669–688.	924
		925
	Judea Pearl. 2000. <i>Causality: Models, Reasoning and Inference</i> . Cambridge University Press, New York.	926
		927
	Judea Pearl. 2010. Causal inference. <i>Causality: objectives and assessment</i> , pages 39–58.	928
		929
	Z Reinis, J Pokorný, V Bazika, J Tiserova, K Goricán, D Horakova, E Stuchlikova, T Havranek, and F Hrabovský. 1981. Prognostic value of the risk profile in the prevention of ischemic heart disease. <i>Bratislavske Lekarske Listy</i> , 76(2):137–150.	930
		931
		932
		933
		934
	Donald B Rubin. 1974. Estimating causal effects of treatments in randomized and nonrandomized studies. <i>Journal of educational Psychology</i> , 66(5):688.	935
		936
		937
	Karen Sachs, Omar Perez, Dana Pe’er, Douglas A Lauffenburger, and Garry P Nolan. 2005. Causal protein-signaling networks derived from multiparameter single-cell data. <i>Science</i> , 308(5721):523–529.	938
		939
		940
		941
	Ruben Sanchez-Romero, Joseph D Ramsey, Kun Zhang, MR K Glymour, Biwei Huang, and Clark Glymour. 2018. Causal discovery of feedback networks with functional magnetic resonance imaging. <i>bioRxiv</i> , page 245936.	942
		943
		944
		945
		946
	Marco Scutari, Maintainer Marco Scutari, and Hiton-PC MMPC. 2019. Package ‘bnlearn’. <i>Bayesian network structure learning, parameter learning and inference, R package version</i> , 4(1).	947
		948
		949
		950
	SerpApi. <a href="#">Google scholar api</a> .	951
	Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, Antti Kerminen, and Michael Jordan. 2006. A linear non-gaussian acyclic model for causal discovery. <i>Journal of Machine Learning Research</i> , 7(10).	952
		953
		954
		955

956	Peter Spirtes and Clark Glymour. 1991. <a href="#">An algorithm for fast recovery of sparse causal graphs</a> . <i>Social Science Computer Review</i> , 9(1):62–72.	Yu Zhou, Xingyu Wu, Beicheng Huang, Jibin Wu, Liang Feng, and Kay Chen Tan. 2024. <a href="#">Causalbench: A comprehensive benchmark for causal learning capability of large language models</a> . <i>arXiv preprint arXiv:2404.06349</i> .	1010
957			1011
958			1012
959	Peter Spirtes, Clark Glymour, and Richard Scheines. 2001. <i>Causation, Prediction, and Search</i> . The MIT Press.		1013
960			1014
961			
962	Masayuki Takayama, Tadahisa Okuda, Thong Pham, Tatsuyoshi Ikenoue, Shingo Fukuma, Shohei Shimizu, and Akiyoshi Sannai. 2024. <a href="#">Integrating large language models in causal discovery: A statistical causal approach</a> . <i>arXiv preprint arXiv:2402.01454</i> .		
963			
964			
965			
966			
967			
968	Aniket Vashishtha, Abbavaram Gowtham Reddy, Abhinav Kumar, Saketh Bachu, Vineeth N Balasubramanian, and Amit Sharma. 2023. <a href="#">Causal inference using llm-guided discovery</a> .		
969			
970			
971			
972	Jianming Wang and Hongbing Shen. 2009. <a href="#">Review of cigarette smoking and tuberculosis in china: intervention is needed for smoking cessation among tuberculosis patients</a> . <i>BMC public health</i> , 9:1–9.		
973			
974			
975			
976	Ming-Gui Wang, Wei-Wei Huang, Yu Wang, Yun-Xia Zhang, Miao-Miao Zhang, Shou-Quan Wu, Andrew J Sandford, and Jian-Qing He. 2018. <a href="#">Association between tobacco smoking and drug-resistant tuberculosis</a> . <i>Infection and drug resistance</i> , pages 873–887.		
977			
978			
979			
980			
981	Caroline E Wright, Katie O’Donnell, Lena Brydon, Jane Wardle, and Andrew Steptoe. 2007. <a href="#">Family history of cardiovascular disease is associated with cardiovascular responses to stress in healthy young men and women</a> . <i>International Journal of Psychophysiology</i> , 63(3):275–282.		
982			
983			
984			
985			
986			
987	Liuyi Yao, Zhixuan Chu, Sheng Li, Yaliang Li, Jing Gao, and Aidong Zhang. 2021. <a href="#">A survey on causal inference</a> . <i>ACM Transactions on Knowledge Discovery from Data (TKDD)</i> , 15(5):1–46.		
988			
989			
990			
991	Matej Zečević, Moritz Willig, Devendra Singh Dhama, and Kristian Kersting. 2023. <a href="#">Causal parrots: Large language models may talk causality but are not causal</a> . <i>Transactions on Machine Learning Research</i> .		
992			
993			
994			
995	Kun Zhang, Mingming Gong, Joseph Ramsey, Kayhan Batmanghelich, Peter Spirtes, and Clark Glymour. 2017. <a href="#">Causal discovery in the presence of measurement error: Identifiability conditions</a> . <i>arXiv preprint arXiv:1706.03768</i> .		
996			
997			
998			
999			
1000	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. <a href="#">A survey of large language models</a> . <i>arXiv preprint arXiv:2303.18223</i> .		
1001			
1002			
1003			
1004			
1005	Guanglin Zhou, Shaoan Xie, Guangyuan Hao, Shiming Chen, Biwei Huang, Xiwei Xu, Chen Wang, Liming Zhu, Lina Yao, and Kun Zhang. 2023. <a href="#">Emerging synergies in causality and deep generative models: A survey</a> .		
1006			
1007			
1008			
1009			

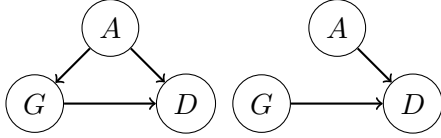


Figure 1: Causal graphs in Example A.1: left-the truth causal graph; right-recovered causal graph by the biased data.

## A Appendix

### A.1 Examples

As follows, we first show an example of statistical estimation-based methods' vulnerability to a type of data bias, the so-called selection bias (Bareinboim et al., 2014).

**Example A.1.** Consider that we would like to investigate the causal relationship of three variables:  $A$  (human age),  $G$  (human gender), and  $D$  (some disease). Assume that the true causal graph is the left figure in Figure 1.

Generally speaking, human age and gender are associated because female has a longer average lifespan. Assuming that this association is only significant for  $A \geq 60$ . However, if each point in a dataset has age under 60, we cannot observe significant difference between the population of male and female. Then, we would recover the causal graph as the right figure in Figure 1.

The second example shows the processing of a well-known constraint-based causal graph discovery algorithm called PC algorithm.

**Example A.2.** Consider a causal discovery task for three variables  $A$ ,  $B$ , and  $C$ , and two different joint probability distributions  $P^1$  and  $P^2$ . We start with a complete undirected graph Figure (a) 2.

Then, by  $P^1$ , we conduct the zero-order independence tests and obtain:  $\hat{\alpha}(AB) = 1$ ,  $\hat{\alpha}(AC) = 1$ , and  $\hat{\alpha}(BC) = 0$ . Then, we keep edges  $(A, B)$  and  $(A, C)$ , and remove  $(B, C)$ , and obtain Figure (b) 2, since  $B$  and  $C$  are not a cause of each other, otherwise they must be associated. Based on the zero-order tests, we can already determine the causal graph as Figure (c) 2, as  $A$  must be a collider since  $B$  and  $C$  are  $d$ -separated by  $\emptyset$ .

On the other hand, if we consider  $P^2$ , we first have zero-order tests showing all pairs are associated, and we cannot remove any edge in Figure (a) 2. We then conduct first-order tests, and obtain:  $\hat{\alpha}(AB | C) = 1$ ,  $\hat{\alpha}(AC | B) = 1$ , and  $\hat{\alpha}(BC | A) = 0$ . Therefore, we can remove the edge  $(B, C)$  from Figure (a) 2, and obtain Figure

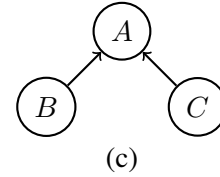
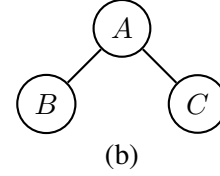
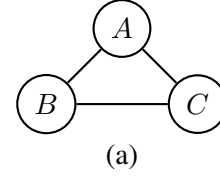


Figure 2: PC algorithm's process.

(b) 2. However, we cannot determine the directions of the edges because all directions of  $A \rightarrow B \rightarrow C$ ,  $A \leftarrow B \leftarrow C$ ,  $A \leftarrow B \rightarrow C$  indicate the conditional independences consistent with  $P^2$ .

## B Enhancing Skeleton Estimation Accuracy by LACR

The theory of Wisdom of the Crowd (Grofman et al., 1983) states that if (1) each individual voter can make the correct decision better than random decision (e.g., by a toss), and (2) voters make their decision independently, then, the accuracy of the collective decision made by simple majority monotonically increases with the number of voters. In LACR, each KB can be seen as a voter. Generally the above conditions tend to be guaranteed because (1) both BG and DOC have high quality and the delivered information is better than random information, and (2) different research papers deliver their results in a relatively independent way because of scientific integrity. Therefore, LACR's decision tends to be more accurate than querying single knowledge base, and it can be improved by adding more relevant documents.

## C Proofs

### C.1 Proof of Proposition 3.1

*Proof.* To verify if  $\zeta_{KB}(ij) = 0$ , by the definition of  $\zeta(ij)$ , we need to check whether there exists a variable set  $V'$  such that  $V'$   $d$ -separates  $v_i$  and

1085  $v_j$ . That is,  $\alpha_{\text{KB}}(ij \mid V') = 0$ . Then, the worst  
 1086 case is that we need to check every combination of  
 1087  $V' \subseteq V \setminus \{v_i, v_j\}$ , which needs  $O(2^{n-2})$  time.  $\square$

## 1088 C.2 Proof of Proposition 3.2

1089 *Proof.* For each usable knowledge base KB, any  
 1090 possible return through the CC prompt must from  
 1091 the set {DIRECTLY ASSOCIATED, INDIRECTLY AS-  
 1092 SOCIATED, INDEPENDENT}.

1093 We first show the first half of the proposi-  
 1094 tion, i.e., the mapping from the conditional as-  
 1095 sociational relationship space between  $v_i$  and  $v_j$ ,  
 1096 i.e.,  $(\hat{\alpha}_{\text{KB}}(ij \mid V'))_{V' \subseteq V \setminus \{v_i, v_j\}}$ , to LLM’s return  
 1097 space based on each usable KB, i.e., {INDEPEN-  
 1098 DENT, DIRECTLY ASSOCIATED, INDIRECTLY AS-  
 1099 SOCIATED} is a surjection. Note that  $(\hat{\alpha}_{\text{KB}}(ij \mid$   
 1100  $V'))_{V' \subseteq V \setminus \{v_i, v_j\}}$  forms a  $2^{|V|-2}$ -dimensional vec-  
 1101 tor, recording  $\alpha_{\text{KB}}(ij \mid V') \in \{0, 1\}$  for all possible  
 1102  $V'$ . We discuss three exclusive cases:

- 1103 1.  $\hat{\alpha}_{\text{KB}}(ij) = 0$ . That is, the zero-order condi-  
 1104 tional associational relationship between  $v_i$   
 1105 and  $v_j$  is independent. This case is mapped to  
 1106 LLM return INDEPENDENT.
- 1107 2. For all possible  $V'$ ,  $\hat{\alpha}_{\text{KB}}(ij \mid V') = 1$ . The  
 1108 case denotes that  $v_i$  and  $v_j$  are always asso-  
 1109 ciated conditioned on any possible  $V'$ , and  
 1110 therefore, this case is mapped to LLM return  
 1111 DIRECTLY ASSOCIATED.
- 1112 3.  $\hat{\alpha}_{\text{KB}}(ij) = 1$  and  $\exists V'$  such that  $|V'| \geq 1$  and  
 1113  $\hat{\alpha}_{\text{KB}}(ij \mid V') = 0$ . In this case, controlling  
 1114 variables in  $V'$ , the statistical association be-  
 1115 tween  $v_i$  and  $v_j$  is eliminated, and then it is  
 1116 mapped to LLM return INDIRECT ASSOCI-  
 1117 ATED.

1118 We then show the second half of the propo-  
 1119 sition, i.e., the mapping from {INDEPENDENT,  
 1120 DIRECTLY ASSOCIATED, INDIRECTLY ASSOCI-  
 1121 ATED} to  $\{\zeta_{\text{KB}}(ij) = 0, \zeta_{\text{KB}}(ij) = 1\}$  is a surjec-  
 1122 tion. If the return is DIRECTLY ASSOCIATED, the  
 1123 KB specifies that  $v_i$  and  $v_j$  cannot be d-separated,  
 1124 and therefore, it is mapped to  $\zeta_{\text{KB}}(ij) = 1$ . On the  
 1125 other hand, if LLM return is INDEPENDENT or IN-  
 1126 DIRECTLY ASSOCIATED, then, it indicates that  $V'$   
 1127 exists that can d-separate  $v_i$  and  $v_j$ , where INDE-  
 1128 PENDENT corresponds to  $V' = \emptyset$ . Therefore, these  
 1129 two last cases correspond to  $\zeta_{\text{KB}}(ij) = 0$ .  $\square$

## C.3 Evidences of Ground Truth Refinement 1130

### C.3.1 ASIA 1131

1132 **Smoking and Tuberculosis** In the documents  
 1133 (Wang and Shen, 2009; Horne et al., 2012; Kim  
 1134 et al., 2022; Wang et al., 2018; Gupta et al., 2022;  
 1135 Lindsay et al., 2014; Amere et al., 2018; Alavi-  
 1136 Naini et al., 2012) fed into LLM as the KB, strong  
 1137 evidence shows that Smoking and Tuberculosis  
 1138 are associated, and the association cannot be elimi-  
 1139 nated by controlling the other variables in the ASIA  
 1140 dataset. This conflicts against the conditional asso-  
 1141 ciational relationship between these two variables  
 1142 in the ground truth causal graph (Appendix C.4),  
 1143 since both of the only two paths have a collider,  
 1144 which indicates that Smoking and Tuberculosis are  
 1145 independent from each other. Based on the sci-  
 1146 entific evidence returned by LACR, a causal link  
 1147 exists between the two factors.

1148 **Bronchitis and X-ray** Documents based on  
 1149 LACR’s response, (Jin et al., 2023; Ntiamoah  
 1150 et al., 2021) show that an association exists be-  
 1151 tween Bronchitis and Positive X-ray report. (Chen  
 1152 et al., 2020) further develops a deep-learning-based  
 1153 method to detect bronchitis directly from X-ray  
 1154 reports for children with age from 1-17 years old.  
 1155 The evidence shows an association between the two  
 1156 variables, and the association is not mediated by  
 1157 the variable “Smoking” as indicated by the causal  
 1158 graph. Therefore, we add a causal link between  
 1159 Bronchitis and X-ray in the ground truth causal  
 1160 graph.

### C.3.2 CORONARY 1161

1162 **Strenuous Mental Work and Family Anamne- 1162**  
 1163 **sis Of Coronary Heart Disease** According to the  
 1164 ground truth causal graph skeleton (Reinis et al.,  
 1165 1981), there is a direct causal relationship between  
 1166 variable Strenuous Mental Work and variable Fam-  
 1167 ily Anamnesis Of Coronary Heart Disease. How-  
 1168 ever, according to the evidence returned by our  
 1169 method and the intuitive description in (Reinis  
 1170 et al., 1981), we observe that this causal linkage  
 1171 should be removed with high probability. By (Ed-  
 1172 wards, 2000) (p.26), this edge is not intuitively  
 1173 expected though it is recovered by the Bayesian  
 1174 learning method. Additionally, non of LACR’s re-  
 1175 sponses suggests the direct association between the  
 1176 two variables, and moreover, it returns evidences  
 1177 showing the positive probability to remove the edge.  
 1178 (Wright et al., 2007) shows that people with Family  
 1179 Anamnesis Of Coronary Heart Disease are easier to 1179

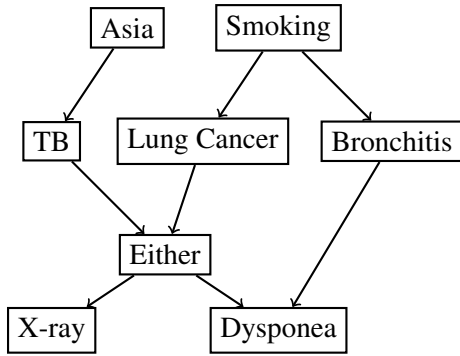


Figure 3: Ground truth causal graph of ASIA in (Iau, 1988).

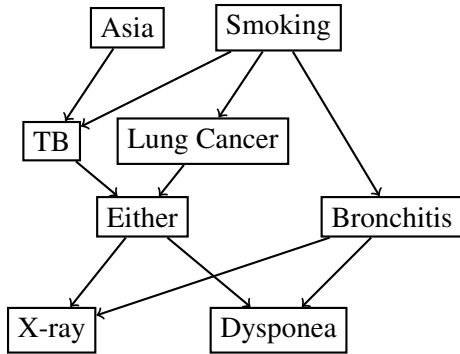


Figure 4: Refined ground truth causal graph of ASIA by LACR.

1180 react to mental stress from work by higher Systolic  
 1181 Blood Pressure. (Hints et al., 2010) shows that the  
 1182 association between psychosocial factors at work  
 1183 and coronary heart disease is largely independent  
 1184 from the Family Anamnesis Of Coronary Heart  
 1185 Disease.

1186 **Systolic Blood Pressure and Family Anamne-**  
 1187 **sis Of Coronary Heart Disease** LACR also re-  
 1188 turns evidence (Barrett-Connor and Khaw, 1984)  
 1189 showing that Family Anamnesis of Coronary Heart  
 1190 Disease and Systolic Blood Pressure are associated  
 1191 even after the adjustment of several variables in-  
 1192 cluding Smoking. We therefore also add this edge  
 1193 between the two variables.

#### 1194 C.4 Additional Experiment Details

1195 The ground truth causal graphs of all datasets in  
 1196 Section 4.

#### 1197 C.5 Additional Experimental Results

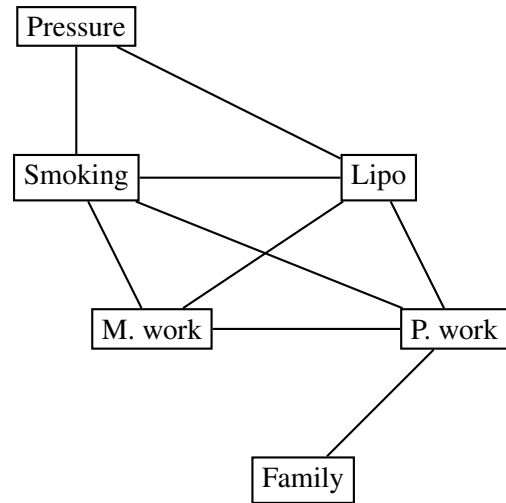


Figure 5: Original ground truth causal graph of CORONARY in (Reinis et al., 1981).

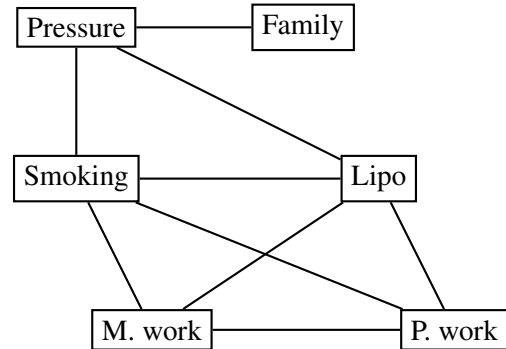


Figure 6: Refined ground truth causal graph of CORONARY by LACR.

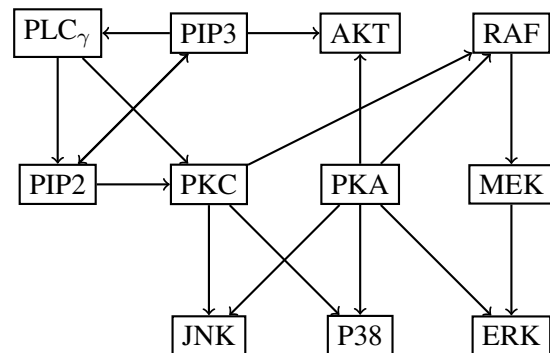


Figure 7: Biological ground truth causal graph in (Sachs et al., 2005).

	ASIA	SACHS
LACR2 (BG)	1	1
LACR2 (DOC)	1	1
LACR2 (PC)	1	1

Table 3: The TEA of LACR 2 on datasets of ASIA, SACHS, based on LACR 1's output skeleton on KB of BG, DOC, and PC, respectively.

1198

1199

1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219

1221

1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251

1253

1254  
1255  
1256  
1257  
1258  
1259  
1260

1262

1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278

## C.6 Prompts

### C.6.1 Association Context

The association relationship between two factors A and B can be associated or independent, and this association relationship can be clarified by the following principles:

1. If A and B are statistically associated or correlated, they are associated, otherwise they are independent.
2. The association relationship can be strongly clarified if there is statistical evidence supporting it.
3. If there is no obvious statistical evidence supporting the association relationship between A and B, it can also be clarified if there is any evidence showing that A and B are likely to be associated or independent statistically.
4. If there is no evidence to clarify the association relationship between A and B, then it is unknown.

### C.6.2 Association Type Context

If two factors A and B are associated, they may be directly associated or indirectly associated with respect to a set of Given Third Factors, and it can be clarified by the following principle:

1. The first principle is to try to find statistical evidence from the given knowledge to clarify the following association types. If you cannot find statistical evidence, at least find evidence that is likely to be able to statistically clarify the association type between A and B. If no obvious evidence can be found, the association type is unknown.
2. If the evidence shows that any factors from the Given Third Factors mediate the association between A and B, then A and B are indirectly associated via these factors.
3. If the evidence shows that by controlling any factors from the Given Third Factors, A and B are not associated any more, then A and B are associated indirectly.
4. If the evidence shows that A and B are still associated even if we control any of the given third factors, then A and B are directly associated.
5. If you think A and B are indirectly associated via any of the given third factors, it must be true that: (1) A and the third factors are directly associated; (2) B and the third factors are directly associated.

### C.6.3 Association Background Reminder

As a scientific researcher in the domains of {domain}, you need to clarify the statistical relationship between some pairs of factors. You first need to get clear of the meanings of the factors in {factors}, which are from your domains, and clarify the interaction between each pair of those factors.

### C.6.4 LLM Association Query (with documents)

Your task is to thoroughly read the given 'Document'. Then, based on the knowledge from the given 'Document', try to find statistical evidence to clarify the association relationship between the pair of 'Main factors' according to the 'Association Context' (delimited by double dollar signs). Consider the given document and the association context. Answer the 'Association Question', write your thoughts, and give the reference in the given document. Respond according to the first expected format (delimited by double backticks).

Document:  
{document}

Main factors:  
{factorA} and {factorB}

Association Context:  
\$\$  
{association\_context}  
\$\$

Association Question:  
Are {factorA} and {factorB} associated?

First Expected Response Format:  
``

Document Identifier: XXX

Thoughts:  
[Write your thoughts on the question]

Answer:  
(A) Associated  
(B) Independent  
(C) Unknown

Reference:  
[Skip this if you chose option C above. Otherwise, provide a supporting sentence from the document for your choice]  
``

### C.6.5 LLM Association Type Query (with documents)

Read and understand the Association Type Context. Consider carefully the role of any of the third factors appearing according to the Association Type Context. Then, based on your thoughts so far, answer the 'Association Type Question' with the 'Given Third Factors', write your thoughts, and give your reference in the given document. Respond according to the expected format (delimited by triple backticks)

Association Type Context:  
\$\$\$  
{association\_type\_context}  
\$\$\$

Given Third Factors:  
{factors} except for {factorA} and {factorB}

Association Type Question: Are {factorA} and {factorB} directly associated or indirectly associated?

Second Expected Response Format:  
````

Thoughts:  
[Write your thoughts on the question]

Answer:  
(D) Directly Associated  
(E) Indirectly Associated  
(C) Unknown

Reference:  
[Skip this if you chose option C above. Otherwise, provide a supporting sentence from the document for your choice]

Intermediary Factors:  
[Skip this if you did not choose D or C above. Otherwise list all factors involved in this indirect association relationship, each separated by a comma]  
````

### C.6.6 LLM Association Query (with background knowledge)

Your task is to thoroughly use the knowledge in your training data to solve a task. Your task is: based on your background knowledge, try to find

1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307

1309  
1310

1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1350  
1351  
1352  
1353  
1354

1356  
1357

1358  
1359  
1360  
1361



1362 statistical evidence to clarify the association  
 1363 relationship between the pair of 'Main factors'  
 1364 according to the 'Association Context' (delimited by  
 1365 double dollar signs).  
 1366 Consider your background knowledge and the  
 1367 association context. Answer the 'Association  
 1368 Question', and write your thoughts. Respond  
 1369 according to the 'First Expected Format' (delimited  
 1370 by double backticks).  
 1371  
 1372 Main factors:  
 1373 {factorA} and {factorB}  
 1374  
 1375 Association Context:  
 1376 \$\$  
 1377 {association\_context}  
 1378 \$\$  
 1379  
 1380 Association Question:  
 1381 Are {factorA} and {factorB} associated?  
 1382  
 1383 First Expected Response Format:  
 1384 ..  
 1385  
 1386 Thoughts:  
 1387 [Write your thoughts on the question]  
 1388  
 1389 Answer:  
 1390 (A) Associated  
 1391 (B) Independent  
 1392 (C) Unknown  
 1393 ..

### 1394 C.6.7 LLM Association Type Query (with 1395 background knowledge)

1396 Read and understand the 'Association Type Context'.  
 1397 Consider carefully the role of any of the third  
 1398 factors appearing according to the Association Type  
 1399 Context. Then, based on your thoughts so far, answer  
 1400 the 'Association Type Question' with the 'Given  
 1401 Third Factors', and write your thoughts. Respond  
 1402 according to the Second Expected Format (delimited  
 1403 by triple backticks)  
 1404  
 1405 Association Type Context:  
 1406 \$\$\$  
 1407 {association\_type\_context}  
 1408 \$\$\$  
 1409  
 1410 Given Third Factors:  
 1411 {factors} except for {factorA} and {factorB}  
 1412  
 1413 Association Type Question: Are {factorA} and {  
 1414 factorB} directly associated or indirectly  
 1415 associated?  
 1416  
 1417 Second Expected Response Format:  
 1418 ...  
 1419  
 1420 Thoughts:  
 1421 [Write your thoughts on the question]  
 1422  
 1423 Answer:  
 1424 (D) Directly Associated  
 1425 (E) Indirectly Associated  
 1426 (C) Unknown  
 1427  
 1428 Intermediary Factors:  
 1429 [Skip this if you did not choose D or C above.  
 1430 Otherwise list all factors involved in this indirect  
 1431 association relationship, each separated by a comma  
 1432 ]  
 1433 ...  
 1434

### 1435 C.6.8 LLM Rethink Query

1436 If none of the Intermediary Factors you found is not  
 1437 in the Given Third Factor list, then, the  
 1438 association type between A and B is direct  
 1439 association.  
 1440 Check your above response, and answer the  
 1441 Association Type Question again. Respond according  
 1442 to the Second Expected Format (delimited by triple  
 1443 backticks).  
 1444  
 1445

1446 Given Third Factors:  
 1447 {factors} except for {factorA} and {factorB}  
 1448  
 1449 Association Type Question: Are {factorA} and {  
 1450 factorB} directly associated or indirectly  
 1451 associated?  
 1452  
 1453 Second Expected Response Format:  
 1454 ...  
 1455  
 1456 Thoughts:  
 1457 [Write your thoughts on the question]  
 1458  
 1459 Answer:  
 1460 (D) Directly Associated  
 1461 (E) Indirectly Associated  
 1462 (C) Unknown  
 1463  
 1464 Intermediary Factors:  
 1465 [Skip this if you did not choose D or C above.  
 1466 Otherwise list all factors involved in this indirect  
 1467 association relationship, each separated by a comma  
 1468 ]  
 1469 ...  
 1470

### 1470 C.6.9 Causal Background Reminder

1471 As a scientific researcher in the domains of {domain  
 1472 }, you need to clarify the statistical relationship  
 1473 between some pairs of factors. You first need to get  
 1474 clear of the meanings of {factorA} and {factorB},  
 1475 which are from your domains, and clarify the  
 1476 interaction between them.  
 1477

### 1478 C.6.10 LLM Causal Direction Query (with 1479 background knowledge)

1480 Your task is to thoroughly use the knowledge in your  
 1481 training data to solve a task. Your task is: based  
 1482 on your background knowledge, try to find  
 1483 statistical evidence to clarify the direction of the  
 1484 causal relationship between the pair of 'Main  
 1485 factors' according to the 'Causal direction context'  
 1486 (delimited by double dollar signs).  
 1487 Consider according to your background knowledge and  
 1488 the 'Causal direction context'. Answer the 'Causal  
 1489 direction question', and write your thoughts.  
 1490 Respond according to the 'Expected Format' (  
 1491 delimited by double backticks).  
 1492  
 1493 Main factors:  
 1494 {factorA} and {factorB}  
 1495  
 1496 Causal direction context:  
 1497 \$\$  
 1498 {causal\_direction\_context}  
 1499 \$\$  
 1500  
 1501 Causal direction question:  
 1502 Is {factorA} the cause of {factorB}, or {factorB}  
 1503 the cause of {factorA}?  
 1504  
 1505 First Expected Response Format:  
 1506 ..  
 1507  
 1508 Thoughts:  
 1509 [Write your thoughts on the question]  
 1510  
 1511 Answer:  
 1512 (A) {factorA} is the cause of {factorB}  
 1513 (B) {factorB} is the cause of {factorA}  
 1514 (C) Unknown  
 1515 ...  
 1516

### 1517 C.6.11 LLM Causal Direction Query (with 1518 documents)

1519 Your task is to thoroughly read the 'Given document'  
 1520 to solve a task. Your task is: based on the 'Given  
 1521 document', try to find statistical evidence to  
 1522 clarify the direction of the causal relationship  
 1523 between the pair of 'Main factors' according to the  
 1524 'Causal direction context' (delimited by double  
 1525 dollar signs).  
 1526

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

First thoroughly read and understand the Given document and the 'Causal direction context'. Then, Answer the 'Causal direction question', and write your thoughts. Respond according to the 'Expected Format' (delimited by double backticks).

Given document:  
{document}

Main factors:  
{factorA} and {factorB}

Causal direction context:  
\$\$  
{causal\_direction\_context}  
\$\$

Causal direction question:  
Is {factorA} the cause of {factorB}, or {factorB} the cause of {factorA}?

First Expected Response Format:  
..

Thoughts:  
[Write your thoughts on the question]

Answer:  
(A) {factorA} is the cause of {factorB}  
(B) {factorB} is the cause of {factorA}  
(C) Unknown

Reference:  
[Skip this if you chose option C above. Otherwise, provide a supporting sentence from the document for your choice]