# IRIS: An Iterative and Integrated Framework for Verifiable Causal Discovery in the Absence of Tabular Data

Anonymous ACL submission

### Abstract

Causal discovery is fundamental to scientific research, yet traditional statistical algorithms face significant challenges, including expensive data collection, redundant computation for known relations, and unrealistic assumptions. While 006 recent LLM-based methods excel at identifying commonly known causal relations, they fall to uncover novel relations. We introduce IRIS (Iterative Retrieval and Integrated System for Real-Time Causal Discovery), a novel framework that addresses these limitations. Starting with a set of initial variables, IRIS automatically collects relevant documents, extracts variables, and uncovers causal relations. Our hybrid causal discovery method combines sta-016 tistical algorithms and LLM-based methods to discover known and novel causal relations. In 017 addition to causal discovery on initial variables, the missing variable proposal component of IRIS identifies and incorporates missing variables to expand the causal graphs. Our approach enables real-time causal discovery from 022 only a set of initial variables without requiring pre-existing datasets.1

### 1 Introduction

034

A fundamental task in various disciplines of science, including biology, economics and healthcare, is to identify and utilize underlying causal relations (Kuhn, 1962). Although interventional experiments are ideal for discovering causal relations, they are often impractical due to ethical, financial, or logistical constraints. Therefore, researchers develop statistical methods to infer causal relations from purely observational tabular data (Pearl, 2009; Spirtes et al., 2000), though such data is often *not* available for a wide range of NLP applications.

Statistical and large language model (LLM)based causal discovery algorithms face distinct challenges that limit their applicability in realworld scenarios. First, traditional statistical algorithms predominantly require high-quality structured tabular data, which is notoriously difficult to obtain. In contrast, LLM-based methods can consistently estimate causal relations explicitly present in their training data without relying on tabular data. However, these models encounter significant limitations when attempting to uncover causal relationships that were not previously documented (Feng et al., 2024; Zečević et al., 2023). Second, statistical causal discovery algorithms require predefined sets of random variables as input, a constraint that significantly limits their flexibility. LLMs, however, demonstrate the capability to reliably extract and identify concepts and entities as variables directly from texts (Zhang et al., 2011; Glymour et al., 2019). Third, most statistical algorithms are theoretically grounded and mathematically verifiable, but operate under assumptions that rarely hold in real-world scenarios, such as the *causal sufficiency* assumption (i.e., the absence of unobservable variables in the causal graph) and acyclicity assumption (*i.e.*, the absence of cycles in the causal graph) (Pearl, 2009; Neal, 2020). In contrast, the verification of LLMs' predictions in causal discovery remains an open challenge.

039

041

043

044

045

047

051

053

054

059

060

061

062

063

064

065

066

067

068

069

070

071

074

075

076

078

079

To address these limitations, we propose IRIS, Iterative **R**etrieval and Integrated **S**ystem for verifiable causal discovery, in the absence of tabular data for statistical methods. To leverage the strengths of both statistical methods and LLMs, our framework takes a *hybrid* causal discovery approach, combining statistical methods with LLM-based causal relation extraction and verification techniques. This hybrid strategy allows us to leverage known causal relations and uncovering novel causal relations. IRIS begins with a set of initial random variables, which are sent as queries to retrieve a collection of relevant documents. Consequently, LLMs are applied to map the unstructured texts into structured tabular

<sup>&</sup>lt;sup>1</sup>Our code and data are available at https://anonymous. 4open.science/r/iris-7378



Figure 1: Illustration of our framework. The input is a set of initial variables. Using the Google Search API with carefully crafted queries and LLMs, we *collect relevant documents and extract variable values* from these documents to construct structured data. For *hybrid causal discovery*, the statistical branch uses the structured data, while the causal relation extraction branch utilizes the retrieved documents. The results from both branches are merged into the final causal graph. The *missing variable proposal* component identifies new variables with causal relations to the initial variables. We then iteratively use the expanded variables as input to our framework to further expand the causal graphs.

data, which is utilized by an appropriate statistical method to perform causal discovery. Its results are further merged with the causal relations predicted and verified by LLMs. This hybrid approach allows cycles in causal graphs, thereby relaxing the *acyclicity* assumption. Additionally, we introduce a variable proposal component to identify new variables that have causal relations with the initial variables. This component allows us to relax the *causal sufficiency* assumption. We then iteratively use the expanded variables as input to our framework, further expanding the causal graphs.

Our experimental results demonstrate that IRIS significantly surpasses strong baselines across all datasets, achieving an average F1 score improvement of 0.14 and a reduction of 0.14 in the average NHD ratio, as detailed in Section 4.1. Evaluations of individual components reveal that each component outperforms its corresponding baselines. Specifically, the evaluation of value extraction component shows that IRIS with GPT-40 exceeds the strong baselines, which also utilizes GPT-40 (Section 4.2). Our hybrid causal discovery method consistently outperforms both statistical algorithms and LLM-based approaches (Section 4.2.3). Lastly, our variable proposal component is more effective compared to prompt-based baselines (Section 4.3). 106

Primary contributions of IRIS are as follows: 1) We introduce an automatic sample collection and value extraction component that significantly reduces the manual labor for data collection in causal discovery tasks. 2) We propose a hybrid causal discovery method that leverages existing causal relations and uncovers novel causal relations. Our method permits cycles in causal graphs, thus relaxing the *acyclicity* assumption. 3) We develop a missing variable proposal component that identifies new variables that may have causal relations with the initial variables, relaxing the *causal sufficiency* assumption. 4) Experimental results demonstrate that IRIS consistently outperforms its baselines, with each component of IRIS also surpassing corresponding baseline methods.

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

# 2 Background

Causal discovery focuses on uncovering causal relations within a set of variables. Given a pair of variables (X, Y), the objective is to determine whether  $X \leftarrow Y, Y \leftarrow X$ , or no causal influence between them, where  $\leftarrow$  denotes causal direction. A key distinction between causal discovery and relation extraction in NLP is that causal discovery can reveal unknown causal relations, whereas relation extraction focuses on transforming relations in free

219

220

221

222

223

224

225

226

227

228

229

230

231

232

185

text into structured relational tuples.

133

134

135

136

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

155

156

157

159

161

162

163

164

165

167

168

169

170

171

172

173

174

175

176

177

178

180

181

184

Although randomized controlled trials and A/B testing are the gold standard for causal discovery (Fisher, 1935), these experimental approaches are often impractical due to ethical or financial limitations. Thus, researchers turn to rely on statistical analysis of observational data to infer causal relations.

Statistical approaches to causal discovery can be broadly classified into: constraint-based methods, such as Peter and Clark (PC) (Spirtes et al., 2000) and inductive causation (IC) (Pearl, 2009); scorebased methods (Heckerman et al., 1995; Chickering, 2002; Koivisto and Sood, 2004; Mooij et al., 2016); and functional methods (Shimizu et al., 2006; Hyvärinen et al., 2010). These methods employ statistical measures from observational data to construct causal graphs but have notable limitations. First, they require resource-intensive and extensive data collection. Second, theoretically, they cannot precisely identify ground-truth causal graphs but instead yield an equivalence class of true causal graphs (Spirtes et al., 2000; Pearl, 2009).

Furthermore, many statistical approaches, such as PC and Greedy Equivalence Search (GES), operate under assumptions. Causal sufficiency assumption posits that all variables are observed and included, neglecting the potential unobserved variables (Neal, 2020). Some algorithms, such as Tetrad condition-based (Silva et al., 2006; Kummerfeld and Ramsey, 2016) and high-order momentsbased approaches (Adams et al., 2021; Chen et al., 2022) focus on only uncover specific types of unobserved variables, such as latent confounders (i.e., common causes). However, our work aims to identify more general unobserved variables, including confounders, mediators, causes, or effects of observed variables. Acyclicity assumption states that causal graphs contain no cycles, which allows causal discovery to align with Bayesian network and simplifies mathematical challenges. However, this assumption often contradicts real-world phenomena. Many causal graphs are known to contain feedback loops, such as the poverty cycle: poverty  $\rightarrow$  limited access to education  $\rightarrow$  lowpaying jobs  $\rightarrow$  poverty, (Banerjee and Duflo, 2012; De Weiss and Sirkin, 2010) and the predator-prey cycle: increase in predator population  $\rightarrow$  decrease in prey population  $\rightarrow$  decrease in predator population (Schmitz, 2017; Abrams, 2001). In contrast to prior work, our causal discovery framework allows for the inclusion of unobserved variables and

permits cycles within causal graphs to align with real-world scenarios.

The advent of LLMs provides new opportunities to address causal discovery (K1c1man et al., 2023; Zečević et al., 2023; Long et al., 2022). These approaches require LLMs to determine the causal relation between a given pair of variable names. However, the reliability of such methods is under scrutiny. Zečević et al. (2023) argue that LLMs may function as "causal parrots", which depend on *memorization* to recall the causal relations present in their training data rather than infer causal relations. This raises concerns about LLMs' generalization to identify causal relations that are rare or absent in pre-training data. Feng et al. (2024) presents empirical evidence that suggests while LLMs excel at reproducing frequent causal relations in pre-training data, they struggle to uncover novel causal relations.

In contrast to approaches that directly employ LLMs for causal discovery, Liu et al. (2024) utilize LLMs to extract variables and their values from collected documents, then apply statistical methods to uncover causal relations among these variables. Our work diverges from this approach by only taking a set of initial variables as input and employing an automated process to collect relevant documents. After variable value extraction, we implement a hybrid causal discovery approach, which integrate both statistical and LLM-based methods. Furthermore, our framework is capable of identifying new variables that exhibit causal relations with the initial set, thereby enabling an iterative process of data collection and causal discovery on an expanded variables set. This iterative method allows for a comprehensive exploration of the causal relations surrounding the initial variables.

# 3 Methodology

We introduce a real-time causal discovery framework, IRIS. Our method differs from prior causal discovery algorithms in three key aspects. First, IRIS does not rely on pre-existing observational data; instead, it automatically collects and extracts observational data related to the initial variables. Second, our hybrid causal discovery component can utilize known causal relations and uncover novel causal relations. Third, our approach relaxes the *acyclicity* and *causal sufficiency* assumptions.

### 3.1 Problem Definition

237

240

241

242

245

246

247

248

249

250

Given a set of initial variables,  $\mathbb{Z} = (\mathbf{z}_1, \mathbf{z}_2, ..., \mathbf{z}_N)$ , where each  $\mathbf{z}_i$  represents one variable, the goal of real-time causal discovery is to automatically collect relevant unstructured data  $\mathbb{D}$  and extract variable values to form structured data  $\mathbb{X}$ , which enables the discovery of causal relations through unstructured and structured data. After identifying causal relations among initial variables, the process involves identifying new variables causally related to the initial variables, resulting in an expanded set of variables  $\mathbb{Z}_m$ . The final output is an expanded causal graph  $\mathcal{G} = (\mathbb{Z}_m, \mathbb{R})$ , where  $\mathbb{R} = (\mathbf{r}_1, ..., \mathbf{r}_l)$ represents the set of causal relations.

## 3.2 Data Collection and Value Extraction

The first step of IRIS comprises two main steps: collection of relevant documents and extraction of variable values. The detailed procedure is outlined in Algorithm 1 in Appendix A.3.

Retrieval of Relevant Documents We retrieve relevant documents using the Google API<sup>2</sup>. To maximize the relevance to initial variables, we create search queries using a stepwise removal approach: 1) Begin with queries containing all variable names (e.g., "smoking" AND "cancer" AND "pollution"). 2) Progressively remove one variable (e.g., "smoking" AND "cancer"). 3) Stop with single-variable queries (e.g., "smoking"). We also use synonyms 260 of variables to enhance coverage. We select the top-261 k retrieved documents for each query. To ensure relevance to most variables, k is higher for queries 263 containing more variables. The retrieval process continues until the total number of collected docu-265 ments reaches a predefined threshold. The resulting document set is denoted as  $\mathbb{D} = (d_1, .., d_T)$ , where 267  $d_i$  represents one document. 268

Extraction of Variable Values We use LLMs to 269 extract variable values from collected documents  $\mathbb{D}$ . Given an LLM M, we design a prompt l including a document  $d_i$  and a description of one variable  $z_i$ . 272 The variable description includes its name and the 273 meaning of its values. We guide the LLM to gen-274 erate responses following multiple thinking steps, 275 simulating human expert reasoning, and provide 276 the final answer in a specific format (Lin et al., 277 2024). This generation process can be denoted as  $o_{ij} = M(l(d_i, z_j))$ , where  $o_{ij}$  is LLM's response 279 regarding the value of variable  $z_i$  in document  $d_i$ .

We then extract the value  $v_{ij}$  from response  $o_{ij}$ . By iterating through all variables and documents, we construct a structured data X where each column represents a variable and each row represents a document. The prompt template for value extraction is presented in Table 4 in Appendix A.4.

## 3.3 Hybrid Causal Discovery

We employ a hybrid causal discovery approach, leveraging both statistical methods and LLM-based relation extraction techniques. The detailed process of our hybrid causal discovery method is outlined in Algorithm 2 in Appendix A.3.

Statistical Causal Discovery For structured data  $\mathbb{X}$ , we employ statistical causal discovery algorithms including PC (Spirtes et al., 2000), GES (Chickering, 2003), and NOTEARS (Zheng et al., 2018). For instance, the PC algorithm performs conditional independence tests between variable pairs, progressively expanding the conditioning sets to determine the presence of causal relations. These algorithms process structured data  $\mathbb{X}$  to produce a causal graph  $\hat{\mathcal{G}}_s$  as the output.

LLM-based Causal Relation Extraction We introduce a novel causal relation extraction method inspired by causal relation verification (Si et al., 2024; Wadden et al., 2022). We treat each potential causal relation as a claim (e.g., "smoking causes lung cancer") and find documents containing both the cause and effect terms (e.g., "smoking" AND "lung cancer"). To ensure the trustworthiness of retrieved documents, we restrict the search domain to reputable academic repositories  $^{3}$ . We then employ LLMs to assess whether each document supports or refutes or not relates with the causal relation using a carefully designed prompt (see Table 5 in Appendix A.4). If a majority of documents support the causal relation, we incorporate it into a causal graph  $\hat{\mathcal{G}}_v$ . Otherwise, it is excluded.

**Graph Merging** The two branches of our hybrid method produce two causal graphs:  $\hat{\mathcal{G}}_s$  from statistical methods and  $\hat{\mathcal{G}}_v$  from the LLM-based approach. To merge them into the final causal graph  $\hat{\mathcal{G}}$ , we post-process the causal graph  $\hat{\mathcal{G}}_s$  by adding highconfidence causal relations from  $\hat{\mathcal{G}}_v$  and removing those strongly refuted by the verification process. This merging strategy is employed for two reasons: (1) the structured data X from the value extraction phase might contain noise; (2) causal relations that

325

327

328

281

<sup>&</sup>lt;sup>2</sup>https://developers.google.com/custom-search/ docs/overview

<sup>&</sup>lt;sup>3</sup>Our search is limited to the following academic website domains: jstor.org, springer.com, ieee.org, ncbi.nlm.nih.gov, sciencedirect.com, scholar.google.com, arxiv.org.

372

396 397

398

399

400

401

402

403

404

405

are widely supported or refuted by trustworthy documents can be treated as known knowledge.

331 3.4 Missing Variable Proposal

This step aims to identify missing variables not included in the initial set but potentially causally related to them, and append these to  $\mathbb{Z}_m$ , as outlined in Algorithm 3 in Appendix A.3.

Variable Abstraction We first use LLMs to abstract missing variables from the retrieved documents D. For each document, LLMs are instructed to analyze the content of each document, identify variables that could influence or be influenced by the initial variables, and then provide the most possible variable in a specified format. The prompt is provided in Table 6 in Appendix A.4.

Variable Selection To select the most promising 344 variables from all abstracted variables, we employ a dual approach combining causal relation verifi-346 cation and statistical measures. Causal Relation Verification: Using the method described in Section 3.3, we verify whether each new variable has a confirmed causal relation with any initial variable. Variables supported by the majority of documents are added to  $\mathbb{Z}_m$ . *Statistical Measure*: We compute the Pointwise Mutual Information (PMI) between each new variable and the initial variables to quantify their dependence, with higher PMI scores indicating stronger potential causal association. The 356 PMI between two variables  $(z_i, z_j)$  is defined as:

$$PMI(\mathbf{z}_i, \mathbf{z}_j) = \log \frac{p(\mathbf{z}_i, \mathbf{z}_j)}{p(\mathbf{z}_i)p(\mathbf{z}_j)} = \log \frac{\frac{o(\mathbf{z}_i, \mathbf{z}_j)}{C}}{\frac{o(\mathbf{z}_i)}{C} \frac{o(\mathbf{z}_j)}{C}}$$
$$= \log \frac{o(\mathbf{z}_i, \mathbf{z}_j)}{o(\mathbf{z}_i)o(\mathbf{z}_j)} + \log \mathcal{C}$$
(1)

where  $o(z_i, z_j)$  is the count of documents where  $(z_i, z_j)$  co-occur,  $o(z_i)$  is the count where  $z_i$  appears, and C is the total number of retrievable documents. Since C is constant, log C is ignored. These counts are obtained by the Google Search API. We compute the PMI score of each abstracted variable with the initial variables and append the top k variables with the highest aggregate PMI scores to  $\mathbb{Z}_m$ .

With the expanded variables  $\mathbb{Z}_m$ , we can iterate the data collection, value extraction, and causal discovery processes to generate an expanded causal graph  $\mathcal{G} = (\mathbb{Z}_m, \mathbb{R})$  that incorporates these missing variables and new causal relations.

## 4 Experiments

# 4.1 Evaluation of the IRIS Framework

# 4.1.1 Experimental Setup

We evaluate the quality of the resulting expanded causal graphs from the complete pipeline of IRIS. **Datasets.** The initial variables are from five datasets: Cancer (Korb and Nicholson, 2010), Respiratory Disease, Diabetes, Obesity (Long et al., 2022), and Alzheimer's Disease Neuroimaging Initiative (ADNI) (Shen et al., 2020).

Our Method and Baselines. We employ GPT-40 as the LLM component, a choice supported by its superior performance across value extraction, causal discovery, and missing variable proposal tasks (see Sections 4.2, 4.2.3, and 4.3). For the statistical causal discovery algorithms in our method, we utilize the Greedy Equivalence Search (GES) algorithm. This selection is based on GES achieving the highest average F1 score and Normalized Hamming Distance (NHD) ratio across all five datasets, as demonstrated in Section 4.2.3. We introduce a baseline method, coined "Prompt", which relies solely on carefully crafted prompts (see Table 7 in Appendix A.4) with LLM to determine causal relations among expanded variables proposed by our missing variable proposal component.

| Method      | Р    | R    | F1↑  | Pred edge | NHD Ratio↓ |  |  |  |  |
|-------------|------|------|------|-----------|------------|--|--|--|--|
| Cancer      |      |      |      |           |            |  |  |  |  |
| Prompt      | 0.64 | 0.32 | 0.43 | 14        | 0.57       |  |  |  |  |
| IRIS        | 0.89 | 0.57 | 0.7  | 18        | 0.3        |  |  |  |  |
| Respiratory |      |      |      |           |            |  |  |  |  |
| Prompt      | 0.67 | 0.36 | 0.47 | 12        | 0.53       |  |  |  |  |
| IRIS        | 0.67 | 0.55 | 0.6  | 18        | 0.4        |  |  |  |  |
| Diabetes    |      |      |      |           |            |  |  |  |  |
| Prompt      | 0.70 | 0.46 | 0.56 | 17        | 0.45       |  |  |  |  |
| IRIS        | 0.76 | 0.5  | 0.6  | 17        | 0.39       |  |  |  |  |
|             |      |      | Obes | ity       |            |  |  |  |  |
| Prompt      | 0.57 | 0.33 | 0.42 | 14        | 0.58       |  |  |  |  |
| IRIS        | 0.67 | 0.58 | 0.62 | 21        | 0.38       |  |  |  |  |
|             |      |      | ADN  | II        |            |  |  |  |  |
| Prompt      | 0.47 | 0.29 | 0.36 | 17        | 0.64       |  |  |  |  |
| IRIS        | 0.5  | 0.36 | 0.42 | 20        | 0.58       |  |  |  |  |

Table 1: Evaluation results of the complete framework. Pred edge indicates the number of predicted edges.

**Evaluation.** To create ground-truth expanded causal graphs, we hire three domain experts to independently annotate each expanded causal graph. Edges are included if at least two annotators agree. With a Krippendorff's alpha of 0.88, inter-annotator agreement is high (Krippendorff, 2011). The detailed annotation instruction is in Table 8 in Appendix A.6. Following Kıcıman et al. (2023);

358

36

366

367

Feng et al. (2024), we evaluate the results of causal discovery using precision, recall, F1 score, and the Ratio of Normalized Hamming Distance (NHD) to baseline NHD. The ratio is defined as ratio  $= \frac{\text{NHD}}{\text{baseline NHD}}$ , where the baseline NHD is derived from the worst-performing causal graph that has the same number of edges as the predicted graph. A lower ratio signifies a more accurate predicted causal graph.

406

407

408

409

410

411

412

413

414

415 416

417

418

419

420

421

422

423 424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

### 4.1.2 Experimental Results and Analysis

As presented in Table 1, IRIS consistently outperforms the Prompt baseline across all datasets, achieving higher F1 scores and lower NHD ratios. The average F1 score improvement is 0.14. The average NHD ratio decreased by 0.14. ADNI exhibits the lowest overall performance for both methods. This may reflect the inherent complexity of Alzheimer's disease causal relations. IRIS predicts more edges than the baseline (averaging 18.8 vs. 14.8 edges), which ensures a higher recall than the baseline (averaging 0.51 vs. 0.35). This indicates that our method's hybrid causal discovery can capture more causal relations effectively. The expanded causal graphs for each dataset are illustrated in Figures 4, 5, 6, 7, 8 in Appendix A.7.

### 4.2 Evaluation of Value Extraction

## 4.2.1 Experimental Setup

**Datasets.** We evaluate the value extraction component of our method using two table-to-text datasets: AppleGastronome and Neuropathic (Liu et al., 2024). These datasets are particularly suitable for our task as they provide tabular data where columns represent variables and rows represent samples. Each row is associated with a corresponding textual description. The datasets are structured as follows: AppleGastronome contains 7 variables and 100 samples. Variable values are -1, 0, or 1. Neuropathic contains 7 variables and 100 samples. Variable values are 0 or 1.

LLMs and Baselines. We utilize state-of-the-445 art LLMs for our method: Llama-3.1-8b-Instruct 446 (Meta, 2024), GPT-3.5-turbo (OpenAI, 2022), GPT-447 40 (OpenAI, 2024). Additionally, we compare our 448 method with COAT, which also utilizes an LLM 449 to extract values of variables from documents (Liu 450 451 et al., 2024). To ensure a fair comparison, we use GPT-40 in both our method and the COAT imple-452 mentation. 453

454 Metrics. Given that variable values are categorical,455 we frame the value extraction task as a classifica-

| AppleGastronome |      |      |      |  |  |  |  |
|-----------------|------|------|------|--|--|--|--|
|                 | Р    | R    | F1   |  |  |  |  |
| COAT - GPT-40   | 0.74 | 0.76 | 0.75 |  |  |  |  |
| IRIS- Llama     | 0.71 | 0.72 | 0.71 |  |  |  |  |
| IRIS- GPT-3.5   | 0.75 | 0.77 | 0.76 |  |  |  |  |
| IRIS- GPT-40    | 0.79 | 0.82 | 0.79 |  |  |  |  |
| Neuropathic     |      |      |      |  |  |  |  |
| COAT - GPT-40   | 0.72 | 0.80 | 0.79 |  |  |  |  |
| IRIS- Llama     | 0.76 | 0.82 | 0.79 |  |  |  |  |
| IRIS- GPT-3.5   | 0.71 | 0.89 | 0.79 |  |  |  |  |
| IRIS- GPT-40    | 0.73 | 1.0  | 0.84 |  |  |  |  |

Table 2: Result of evaluation of value extraction. Llamarepresents Llama-3.1-8b-instruct

tion problem, predicting the value of a variable in a given document. Therefore, we employ standard classification metrics: precision, recall, and F1. 456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

### 4.2.2 Experimental Results and Analysis

Table 2 presents the evaluation results of our value extraction method across different LLMs on the AppleGastronome and Neuropathic datasets. Our method's superior performance with GPT-40, compared to COAT using the same LLM, indicates that our approach is more effective than COAT under identical LLM. In both datasets, we observe a consistent trend of improvement from Llama-3.1-8b-Instruct to GPT-3.5, and further to GPT-40 when using our method. This progression aligns with the general understanding that more advanced LLMs tend to perform better on complex tasks. Overall, the models perform better on the Neuropathic dataset compared to AppleGastronome. This could be attributed to the simpler binary values of the Neuropathic dataset (values 0 or 1) compared to the ternary values in AppleGastronome (-1, 0, 1). The additional complexity in AppleGastronome might introduce more opportunities for misclassification. The high performance of GPT-40 suggests that it could be highly effective for value extraction in our framework.

### 4.2.3 Evaluation of Causal Discovery

### 4.2.4 Experimental Setup

**Datasets.** We evaluate our hybrid causal discovery component to five datasets: Cancer (Korb and Nicholson, 2010), Respiratory Disease, Diabetes, Obesity (Long et al., 2022), and Alzheimer's Disease Neuroimaging Initiative (ADNI) (Shen et al., 2020). These causal graphs are annotated by domain experts. The ground-truth causal graphs are presented in Figure 3 in Appendix A.5.



Figure 2: Evaluation results of causal discovery component on five datasets. A higher F1 score indicates better performance, while a lower NHD ratio reflects better performance. VCR refers to verified causal relations that are extracted from relevant academic documents and validated by LLMs. "Llama" refers to the use of the Llama-3.1-8b-instruct model as a substitute for GPT-40 in our method.

**Baselines.** We compare our method against several 492 baselines: 1) Pairwise-LLM constructs queries for 493 each pair of variables, using LLMs to determine 494 causal relations. The computational complexity of 495 this method is  $O(n^2)$  (Feng et al., 2024). 2) BFS-496 LLM employs a breadth-first search approach with 497 LLMs, achieving linear computational complexity 498 (Jiralerspong et al., 2024). 3) COAT utilizes LLM 499 to extract values from documents, then applies the PC algorithm for causal discovery (Liu et al., 2024). 501 In our hybrid causal discovery approach, for statistical algorithms, we utilize PC (Spirtes et al., 2000), GES (Chickering, 2003), and NOTEARS (Zheng et al., 2018). Among the three statistical methods, 505 we select the one that demonstrates the best performance for hybrid causal discovery. Based on 507 our value extraction results (see Table 2), we use GPT-40, which demonstrated the best performance, 509 as the LLM for both our method and the baseline 510 approaches. To illustrate how different LLMs af-511 fect the performance of our method, we employ the 512 Llama-3.1-8b-instruct model as a counterpart. 513

514Metrics.We evaluate the quality of causal graphs515using precision, recall, F1, and NHD ratio as de-516tailed in Section 4.1.

### 4.2.5 Experimental Results and Analysis

The evaluation results of the causal discovery component across five datasets are presented in Figure 2. More detailed results are presented in Table 9, 10, 11, 12, 13 in Appendix A.8. In these results, our hybrid method consistently outperforms baseline methods across all datasets. This highlights the effectiveness of combining statistical algorithms with LLM-based methods.

We observe that the performance of individual statistical algorithms (GES, NOTEARS, PC) varied across datasets. PC excels in Respiratory Disease and Obesity. GES demonstrates optimal performance on Diabetes and Obesity. NOTEARS performs best on Cancer and ADNI but struggles significantly with Diabetes and Obesity, achieving a 0 F1 score and a 1 NHD ratio. This variation highlights the importance of selecting statistical algorithms based on the characteristics of the observational data, which presents a compelling area for further research. From our experiments, both GES and PC exhibit strong performances; however, GES outperforms PC, with an average F1 score that is 0.09 points higher and an average NHD ratio that is 0.09 points lower. Given these results, GES is recommended as the primary choice when the suitability of the algorithm is uncertain. When comparing the performance of Llama-3.1-8b-instruct and GPT-40 in our method, GPT-40 consistently outperforms Llama-3.1-8b-instruct across all datasets, with a particularly significant gap observed in the ADNI dataset. We believe this discrepancy arises because ADNI involves specialized knowledge that is less commonly represented in the pre-training data of Llama-3.1-8b-instruct.

526

527

528

529

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

LLM-based methods (Pairwise-LLM and BFS-LLM) show competitive performance on simpler datasets. They perform well on the Cancer and Respiratory Disease datasets. However, their performance degrades on more complex datasets like ADNI. This suggests that while LLMs have potential in causal discovery, they may struggle with more complex causal relations, possibly due to the lower occurrence of such domain-specific causal
relations in their training data (Feng et al., 2024).
The COAT method yields results similar to IRISPC because both approaches extract values from
documents and then perform causal discovery using the PC algorithm.

In conclusion, our experimental results consistently demonstrate that integrating the Verified Causal Relations (VCR) component with statistical algorithms significantly enhances causal discovery performance across datasets, thereby validating the effectiveness of our hybrid approach.

572

574

577

579

581

582

584

585

587

570

566

4.3 Evaluation of Missing Variable Proposal

| Method        | Cancer | Respiratory<br>Disease | Diabetes | Obesity | ADNI  |
|---------------|--------|------------------------|----------|---------|-------|
| Prompt        | 0.4    | 0.25                   | 0.5      | 0.25    | 0.25  |
| MVP - NoVCR   | 0.6    | 0.75                   | 0.5      | 0.75    | 0.25  |
| MVP - NoStats | 0.6    | 0.75                   | 0.75     | 1.0     | 0.375 |
| MVP (Llama)   | 0.4    | 0.5                    | 0.25     | 0.5     | 0.125 |
| MVP           | 0.8    | 0.75                   | 1.0      | 1.0     | 0.5   |

Table 3: Evaluation results (success rate) of the missing variable proposal (MVP) component. MVP-NoVCR excludes verified causal relation extraction; MVP-NoStats omits statistical approaches; Llama is the Llama-3.1-8b-instruct. Except MVP (Llama), other methods use GPT-4o as the LLM.

## 4.3.1 Experimental Setup

**Datasets.** Evaluating the missing variable proposal component presents a unique challenge: the ground-truth missing variables are inherently unknown in real-world scenarios. To address this, we simulate missing variables and assess our method's ability to identify them. We start with complete, ground-truth causal graphs and systematically remove variables to create incomplete graphs. We employ five causal graphs: Cancer, Respiratory Disease, Diabetes, Obesity, and ADNI. For each causal graph, we iteratively remove one variable at a time, creating multiple test cases per graph. We then apply our missing variable proposal component to these incomplete graphs, aiming to identify the removed variables.

589 **Our Method and Baselines.** To ensure a com-590 prehensive evaluation, we introduced a baseline 591 method that uses LLMs to directly suggest new 592 variables via a prompt-based approach. For both 593 our missing variable proposal component and the 594 baseline, we use GPT-40 as the primary LLM. To 595 compare the performance of different LLMs, we 596 also replace GPT-40 with Llama-3.1-8b-instruct in 597 our component. **Metrics.** We evaluate the performance using a *success rate* metric, calculated as follows: 1) For each incomplete causal graph, we check if our method successfully proposes the removed variable in its final set of proposed variables  $\mathbb{Z}_m$ . 2) We count a "success" for each correctly proposed variable. 3) The success rate is computed as: Success Rate = Number of Successes / Total Number of Incomplete Graphs. For instance, in a causal graph with five variables, we create five different incomplete graphs by removing each variable. If our method correctly proposes the removed variable in three of these five graphs, the success rate would be 0.6. For the statistical approach, we select the top-5 variables based on their PMI scores.

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

# 4.3.2 Experimental Results and Analysis

The evaluation results of our Missing Variable Proposal (MVP) component are presented in Table 3. The MVP method consistently outperforms other variants across all datasets. This demonstrates the effectiveness of combining VCR with statistical approach in identifying missing variables. Ablation studies indicate that both VCR and statistical approaches play a crucial role in enhancing the success rate of the MVP. The performance gap between MVP and MVP-Llama indicates the superior capability of GPT-40 in understanding and reasoning about causal relations. The prompt-based baseline consistently underperforms compared to our framework, indicating that relying solely on the internal knowledge of LLMs is not reliable for proposing missing variables.

# 5 Conclusion

In this paper, we introduce IRIS, a novel framework that addresses several longstanding challenges in causal discovery. By integrating automated data collection, hybrid causal discovery methods, and a variable proposal components, IRIS significantly advances our ability to uncover causal relations in real-world scenarios. Our approach not only reduces the reliance on extensive manual data collection but also leverages existing knowledge in order to facilitate the discovery of novel causal relations with novel variables. Our experimental results show that IRIS consistently outperforms competitive baselines. Future work could aim to enhancing the scalability of IRIS for larger and more complex causal graphs by integrating causal relations extracted from texts with the ones identified through statistical algorithms.

## **18** Limitations

649Our approach to uncovering causal relations using650retrieved documents and LLMs has certain limi-651tations. A primary challenge is the potential bias652inherent in both the data and the LLMs. Retrieved653documents may contain sampling biases, inaccu-654racies, or incomplete coverage of causal relations.655Likewise, LLMs may inherit biases from their pre-656training data or face limitations in generalization,657potentially affecting their interpretation of causal658relationships. To mitigate these issues, we retrieve659documents from reliable sources, such as academic660websites, and leverage state-of-the-art LLMs like661GPT-4.

## Ethics Statement

662

664

667

671

672

674

675

676

688

693

We acknowledge the importance of ACL Code of Ethics and agree with it. We ensure that our study is compatible with the provided code.

Our work involves uncovering causal relations using retrieved documents and LLMs, and we acknowledge the ethical considerations associated with this approach. The potential biases inherent in both the retrieved data and the LLMs pose a significant challenge. To mitigate these risks, we prioritize retrieving data from credible sources, such as academic publications and verified websites, to ensure the reliability of the input data. Additionally, we employ state-of-the-art LLMs, like GPT-4, which are designed to provide high-quality and robust outputs. However, we recognize that no system is entirely free from bias, and users of this framework should exercise caution in interpreting its results.

The evaluation of our method involves hiring human experts to annotate causal graphs. We have ensured that the annotation process adheres to ethical guidelines, including providing fair compensation for their contributions. Rigorous measures have been taken to thoroughly anonymize the causal graphs, which do not contain any personally identifiable information or sensitive data related to the contributors. The causal graphs were compiled with contributions from PhD students, which may inherently introduce biases influenced by their demographic backgrounds. We advise researchers utilizing this dataset to carefully account for these potential biases, particularly in studies related to AI fairness, bias, and safety.

### References

Peter A Abrams. 2001. Predator-Prey Interactions. In Evolutionary Ecology: Concepts and Case Studies . Oxford University Press. 696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

726

727

728

730

731

732

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

- Jeffrey Adams, Niels Richard Hansen, and Kun Zhang. 2021. Identification of partially observed linear causal models: Graphical conditions for the nongaussian and heterogeneous cases. In Advances in Neural Information Processing Systems.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-RAG: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations.*
- Ananth Balashankar, Sunandan Chakraborty, Samuel Fraiberger, and Lakshminarayanan Subramanian. 2019. Identifying predictive causal factors from news streams. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2338–2348, Hong Kong, China. Association for Computational Linguistics.
- A.V. Banerjee and E. Duflo. 2012. *Poor Economics:* A Radical Rethinking of the Way to Fight Global *Poverty*. PublicAffairs.
- Adrián Bazaga, Pietro Lio, and Gos Micklem. 2024. Unsupervised pretraining for fact verification by language model distillation. In *The Twelfth International Conference on Learning Representations*.
- Giannis Bekoulis, Christina Papagiannopoulou, and Nikos Deligiannis. 2021. A review on fact extraction and verification. *ACM Computing Surveys (CSUR)*, 55(1):1–35.
- Quoc-Chinh Bui, Breanndán Ó Nualláin, Charles A Boucher, and Peter MA Sloot. 2010. Extracting causal relations on hiv drug resistance from literature. *BMC bioinformatics*, 11:1–11.
- Du-Seong Chang and Key-Sun Choi. 2006. Incremental cue phrase learning and bootstrapping method for causality extraction using cue phrase and word pair probabilities. *Information Processing and Management*, 42(3):662–678.
- Zhengming Chen, Feng Xie, Jie Qiao, Zhifeng Hao, Kun Zhang, and Ruichu Cai. 2022. Identification of linear latent variable model with arbitrary distribution. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 6350–6357.
- David Maxwell Chickering. 2002. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554.
- David Maxwell Chickering. 2003. Optimal structure identification with greedy search. *J. Mach. Learn. Res.*, 3(null):507–554.

- 750 751 760 761 767 770 773 776 777 778 779 781 786 791 796 797 800

- 803

- Tirthankar Dasgupta, Rupsa Saha, Lipika Dey, and Abir Naskar. 2018. Automatic extraction of causal relations from text using linguistically informed deep neural networks. In Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue, pages 306–316, Melbourne, Australia. Association for Computational Linguistics.
- Susan Pick De Weiss and Jenna Sirkin. 2010. Breaking the poverty cycle: The human basis for sustainable development. Oxford University Press.
- Tao Feng, Lizhen Qu, Niket Tandon, Zhuang Li, Xiaoxi Kang, and Gholamreza Haffari. 2024. From pretraining corpora to large language models: What factors influence llm performance in causal discovery tasks? Preprint, arXiv:2407.19638.
  - R. A. Fisher. 1935. The Design of Experiments. Oliver and Boyd.
  - Daniela Garcia. 1997. Coatis, an nlp system to locate expressions of actions connected by causality links. In Proceedings of the 10th European Workshop on Knowledge Acquisition, Modeling and Management, EKAW '97, page 347-352, Berlin, Heidelberg. Springer-Verlag.
  - Clark Glymour, Kun Zhang, and Peter Spirtes. 2019. Review of causal discovery methods based on graphical models. Frontiers in Genetics, 10.
  - David Heckerman, Dan Geiger, and David M Chickering. 1995. Learning bayesian networks: The combination of knowledge and statistical data. Machine learning, 20:197-243.
- Miguel A Hernán, Sonia Hernández-Díaz, and James M Robins. 2004. A structural approach to selection bias. *Epidemiology*, 15(5):615–625.
- Aapo Hyvärinen, Kun Zhang, Shohei Shimizu, and Patrik O. Hoyer. 2010. Estimation of a structural vector autoregression model using non-gaussianity. Journal of Machine Learning Research, 11(56):1709–1731.
- Thomas Jiralerspong, Xiaoyin Chen, Yash More, Vedant Shah, and Yoshua Bengio. 2024. Efficient causal graph discovery using large language models. Preprint, arXiv:2402.01207.
- Christopher S. G. Khoo, Syin Chan, and Yun Niu. 2000. Extracting causal knowledge from a medical database using graphical patterns. In Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics, pages 336–343, Hong Kong. Association for Computational Linguistics.
- Kyungha Kim, Sangyun Lee, Kung-Hsiang Huang, Hou Pong Chan, Manling Li, and Heng Ji. 2024. Can llms produce faithful explanations for fact-checking? towards faithful explainable fact-checking via multiagent debate. Preprint, arXiv:2402.07401.
- Mikko Koivisto and Kismat Sood. 2004. Exact bayesian structure discovery in bayesian networks. The Journal of Machine Learning Research, 5:549-573.

Kevin B. Korb and Ann E. Nicholson. 2010. Bayesian Artificial Intelligence, Second Edition, 2nd edition. CRC Press, Inc., USA.

805

806

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

- Klaus Krippendorff. 2011. Computing krippendorff's alpha-reliability.
- Thomas S. Kuhn. 1962. The Structure of Scientific Revolutions. University of Chicago Press, Chicago.
- Erich Kummerfeld and Joseph Ramsey. 2016. Causal clustering for 1-factor measurement models. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, page 1655-1664, New York, NY, USA. Association for Computing Machinery.
- Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. 2023. Causal reasoning and large language models: Opening a new frontier for causality. Preprint, arXiv:2305.00050.
- Yufeng Li, Rrubaa Panchendrarajan, and Arkaitz Zubiaga. 2024. Factfinders at checkthat! 2024: Refining check-worthy statement detection with llms through data pruning. Preprint, arXiv:2406.18297.
- Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Chandu, Chandra Bhagavatula, and Yejin Choi. 2024. The unlocking spell on base LLMs: Rethinking alignment via in-context learning. In The Twelfth International Conference on Learning Representations.
- Chenxi Liu, Yongqiang Chen, Tongliang Liu, Mingming Gong, James Cheng, Bo Han, and Kun Zhang. 2024. Discovery of the hidden world with large language models. Preprint, arXiv:2402.03941.
- Stephanie Long, Tibor Schuster, and Alexandre Piché. 2022. Can large language models build causal graphs? In NeurIPS 2022 Workshop on Causality for Real-world Impact.

Meta. 2024. Meet llama 3.1.

- Isabelle Mohr, Amelie Wührl, and Roman Klinger. 2022. CoVERT: A corpus of fact-checked biomedical COVID-19 tweets. In Proceedings of the Thirteenth Language Resources and Evaluation Conference, pages 244–257, Marseille, France. European Language Resources Association.
- Joris M. Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. 2016. Distinguishing cause from effect using observational data: methods and benchmarks. J. Mach. Learn. Res., 17(1):1103-1204.
- Brady Neal. 2020. Introduction to Causal Inference from a Machine Learning Perspective.
- OpenAI. 2022. Introducing chatgpt.
- OpenAI. 2024. Hello gpt-4o.

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

909

910

911

912

- 857
- 859
- 861
- 863
- 864
- 866 867
- 870
- 871 872

874

- 875
- 878

880

- 881

900 901

902 903

904 905

906

907 908

- Rrubaa Panchendrarajan and Arkaitz Zubiaga. 2024. Claim detection for automated fact-checking: A survey on monolingual, multilingual and cross-lingual research. Natural Language Processing Journal, 7:100066.
- Judea Pearl. 2009. Causality. Cambridge university press.
- Oswald Schmitz. 2017. Predator and prey functional traits: understanding the adaptive machinery driving predator-prey interactions. F1000Research, 6.
- Xinpeng Shen, Sisi Ma, Prashanthi Vemuri, and Gyorgy Simon. 2020. Challenges and opportunities with causal discovery algorithms: application to alzheimer's pathophysiology. Scientific reports, 10(1):2975.
- Shohei Shimizu, Patrik O. Hoyer, Aapo Hyvärinen, and Antti Kerminen. 2006. A linear non-gaussian acyclic model for causal discovery. Journal of Machine Learning Research, 7(72):2003–2030.
- Jiasheng Si, Yibo Zhao, Yingjie Zhu, Haiyang Zhu, Wenpeng Lu, and Deyu Zhou. 2024. CHECKWHY: Causal fact verification via argument structure. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15636–15659, Bangkok, Thailand. Association for Computational Linguistics.
  - Ricardo Silva, Richard Scheine, Clark Glymour, and Peter Spirtes. 2006. Learning the structure of linear latent variable models. Journal of Machine Learning Research, 7(8):191-246.
- Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. 2000. Causation, Prediction, and Search. MIT press.
- David Wadden, Kyle Lo, Bailey Kuehl, Arman Cohan, Iz Beltagy, Lucy Lu Wang, and Hannaneh Hajishirzi. 2022. SciFact-open: Towards open-domain scientific claim verification. In Findings of the Association for Computational Linguistics: EMNLP 2022, pages 4719-4734, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Somin Wadhwa, Silvio Amir, and Byron Wallace. 2023. Revisiting relation extraction in the era of large language models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15566-15589, Toronto, Canada. Association for Computational Linguistics.
- Haoran Wang and Kai Shu. 2023. Explainable claim verification via knowledge-grounded reasoning with large language models. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 6288-6304, Singapore. Association for Computational Linguistics.

- Jie Yang, Soyeon Caren Han, and Josiah Poon. 2022. A survey on extraction of causal relations from natural language text. Knowledge and Information Systems, 64(5):1161–1186.
- Matej Zečević, Moritz Willig, Devendra Singh Dhami, and Kristian Kersting. 2023. Causal parrots: Large language models may talk causality but are not causal. Transactions on Machine Learning Research.
- Kun Zhang, Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. 2011. Kernel-based conditional independence test and application in causal discovery. In Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence, UAI'11, page 804-813, Arlington, Virginia, USA. AUAI Press.
- Xiaoyan Zhao, Yang Deng, Min Yang, Lingzhi Wang, Rui Zhang, Hong Cheng, Wai Lam, Ying Shen, and Ruifeng Xu. 2024. A comprehensive survey on relation extraction: Recent advances and new frontiers. ACM Comput. Surv., 56(11).
- Youwen Zhao, Xiangbo Yuan, Ye Yuan, Shaoxiong Deng, and Jun Quan. 2023. Relation extraction: advancements through deep learning and entity-related Social Network Analysis and Mining, features. 13(1):92.
- Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P. Xing. 2018. Dags with no tears: continuous optimization for structure learning. In Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18, page 9492-9503, Red Hook, NY, USA. Curran Associates Inc.

#### Appendix Α

#### **Related Work** A.1

Causal Discovery Causal discovery aims to uncover causal structures among variables, distinguishing itself from relation extraction in NLP by revealing novel causal relations rather than merely extracting known relations. While experimental approaches such as randomized controlled trials are gold standard methods(Fisher, 1935), practical limitations often necessitate statistical methods using observational data. These include constraintbased and score-based approaches (Spirtes et al., 2000; Pearl, 2009; Heckerman et al., 1995). However, statistical methods face challenges in data collection and theoretical limitations. Recent advancements in LLMs have introduced new possibilities for causal discovery without direct data access (Kıcıman et al., 2023; Zečević et al., 2023; Long et al., 2022). However, concerns about LLMs functioning as "causal parrots" and their ability to generalize to novel relations have been raised

(Zečević et al., 2023; Feng et al., 2024). Alterna-961 tive approaches, such as using LLMs for variable 962 proposer and combining them with statistical meth-963 ods (Liu et al., 2024), have emerged. Our work 964 builds upon these ideas, introducing an automated 965 document collection process, a hybrid causal dis-966 covery method integrating statistical and relation 967 extraction techniques, and a hybrid approach for new variable proposal. 969

**Relation Extraction** Relation extraction aims to 970 transform unstructured textual relations into struc-971 tured relation tuples of the form  $\langle e_1, r, e_2 \rangle$ , where  $e_1$  and  $e_2$  represent entities and r denotes the relation between them (Yang et al., 2022; Dasgupta 974 975 et al., 2018). While relation extraction can identify cause-effect relationships from documents, it fun-976 damentally differs from causal discovery in that it 977 relies on explicitly stated relations in texts, whereas 978 causal discovery can uncover novel causal relation-979 ships from observational data even in the absence of explicit textual mentions. Nevertheless, relation 981 extraction can serve as a complementary method for identifying commonly known causal relations in textual data. Several studies have focused on extracting causal relations from natural language texts 985 (Balashankar et al., 2019; Bui et al., 2010; Chang 986 and Choi, 2006). The methods for causality ex-987 traction can be divided into pattern-based and deep learning-based approaches. Pattern-based methods utilize predefined linguistic patterns to extract relevant text segments, which are then converted into 991 tuples using hand-crafted algorithms (Garcia, 1997; 992 Khoo et al., 2000). However, these methods often suffer from limited coverage of causal relations and require significant effort in pattern design. Deep learning-based methods employ neural networks to learn high-level abstract features and represen-997 tations from sentences, framing relation extraction as a sequence-to-sequence task (Zhao et al., 2023, 2024). While these approaches offer improved per-1000 formance, they typically require large fine-tuning 1001 datasets and may not consistently produce struc-1002 1003 turally correct output tuples.

> A notable limitation of many relation extraction systems is the lack of verification for extracted relations, potentially leading to the extraction of false or unreliable relations from untrustworthy sources (Si et al., 2024; Wadhwa et al., 2023). Our work addresses this issue by adopting a novel approach: instead of directly extracting causal relations from documents, we pre-create textual men-

1004

1005

1006 1007

1008

1009

1011

tions of causal relations (e.g., "smoking causes lung 1012 cancer") and employ LLMs to verify the veracity 1013 of these relations based on relevant documents. We 1014 consider a causal relation to hold if the majority of 1015 documents support its veracity, thereby enhancing 1016 the reliability of our extracted causal relations. 1017

Claim Verification Claim verification aims to 1018 assess the veracity of claims based on relevant doc-1019 uments (Bekoulis et al., 2021). This process typi-1020 cally encompasses several key components: claim 1021 detection, document retrieval, veracity prediction, 1022 and explanation generation. Research in this field 1023 often focuses on specific aspects of the verifica-1024 tion pipeline. For instance, Panchendrarajan and 1025 Zubiaga (2024) and Li et al. (2024) concentrate 1026 on identifying check-worthy statements from large 1027 text corpora. Others, such as Wadden et al. (2022) 1028 and Mohr et al. (2022), prioritize veracity predic-1029 tion, while Wang and Shu (2023) emphasize the 1030 importance of generating explanations for verifica-1031 tion outcomes. The emergence of LLMs has sig-1032 nificantly influenced the field, with numerous stud-1033 ies leveraging LLMs for claim verification through 1034 carefully crafted prompts (Kim et al., 2024; Bazaga 1035 et al., 2024; Asai et al., 2024). Building on these 1036 advancements, one branch of our hybrid causal dis-1037 covery approach reframes causal discovery as a 1038 causal relation verification task. We employ LLMs 1039 to assess the veracity of causal relations based on 1040 retrieved documents, subsequently incorporating 1041 verified relations into a causal graph. This method-1042 ology bridges the gap between traditional claim 1043 verification techniques and causal discovery, offer-1044 ing a novel approach to uncovering and validating 1045 causal relations.

## A.2 Reproducibility Statement

We release our code and scripts at https:// anonymous.4open.science/r/iris-7378. Detailed descriptions of the algorithms used in each component of our framework can be found in Appendix A.3. We provide all prompts utilized throughout our framework in Appendix A.4. The ground-truth causal graphs employed in our evaluation experiments are outlined in Appendix A.5. Additionally, Appendix A.6 presents human annotation instruction and interface for the human annotation tasks involved in evaluating the expanded causal graphs. The annotated expanded causal graphs, alongside the predicted causal graphs, are documented in Appendix A.7.

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1064

1065

1066

1067

1068

1069

# A.3 Algorithms

In this section, we provide detailed descriptions of the algorithms for each component of our method. The data collection and value extraction process is outlined in Algorithm 1. The hybrid causal discovery algorithm can be found in Algorithm 2. Finally, the algorithm for proposing missing variables is detailed in Algorithm 3.

| Algorithm I Document Collection and Value Ex-   |
|---|
| traction  |
| <b>Require:</b> Initial Variables $\mathbb{Z}$ , LLM $M$ , threshold                    |
| T, prompt $l$   |
| Document Collection   |
| $\mathbb{D} \leftarrow \emptyset  \triangleright$ Initialize an empty set for collected |
| documents   |
| while $ \mathbb{D}  < T$ do   |
| $queries \qquad \leftarrow$   |
| $[(z_1, z_2, \ldots, z_n), (z_1, z_2, \ldots, z_{n-1}), \ldots, (z_i)]$                 |
| ▷ queries considering all variables and their   |
| synonyms  |
| for each q in queries do  |

 $n \leftarrow 20 \times \text{len}(q)$  > Determine the number of URLs to collect

 $urls \leftarrow google\_search(q, n) \triangleright Search$ with query q and retrieve top-n URLs

for each *url* in *urls* do

 $D \leftarrow \text{extract text from } url$ 

 $\mathbb{D} \leftarrow \mathbb{D} \cup \{D\} \mathrel{\vartriangleright} \mathsf{Add} \text{ extracted text}$  to the document set

end for end while

# Value Extraction

 $V \leftarrow Matrix of dimensions T \times N \triangleright Initialize$ a matrix with T rows and N columns for each  $d_i$  in D do for each  $z_j$  in Z do  $o_{ij} \leftarrow M(l(d_i, z_j)) \triangleright Determine value$ of  $z_j$  in  $d_i$  by LLM  $v_{ij} \leftarrow extract(o_{ij}) \triangleright Extract value from$ LLM output  $V[i][j] \leftarrow v_{ij} \triangleright Store$  the value  $v_{ij}$  in matrix V at position (i, j)end for end for Output: D, V

# A.4 Prompts

In this section, we show prompts we used in IRIS1071in Table 4, 5, 6. The prompt used in the "Prompt"1072baseline in evaluation of expanded causal graphs is1073shown in Table 7.1074

1070

1075

# A.5 Ground-Truth Causal Graphs

The ground-truth causal graphs for causal discovery1076can be found in Figure 3.1077





Figure 3: The ground-truth causal graphs from original sources (Hernán et al., 2004; Long et al., 2022; Shen et al., 2020; Korb and Nicholson, 2010).



Table 4: The prompt for value extraction, where doc indicates the content of a document, var refers to a variable name.

| Given a document: {doc}   |
|---|
|   |
| Please complete the below task.   |
| We have a claim: '{claim}'. We need to check the veracity of this claim. The value of veracity is True or False or Unknown. |
| True indicates that the given document supports this claim,   |
| False indicates that the given document refutes the claim.  |
| Unknown indicates that the given document has no relation to the claim.   |
| Please form the answer with a logical reasoning chain according to the following format.                                    |
| First, provide an introductory sentence that explains what information will be discussed.                                   |
| Next, list the logical reasoning chain in detail, ensuring clarity and precision.   |
| Finally, conclude the veracity of claim '{claim}' using the following template:   |
| The veracity of claim '{claim}' is  |

Table 5: The prompt for causal relation verification, where doc indicates the content of a document, claim refers to a causal relation (*e.g.*, smoking causes lung cancer).

# A.6 Causal Relation Annotation Task

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

The detailed instructions for the causal relation annotation task are presented in Table 8. This table provides comprehensive guidance to annotators on how to identify and annotate causal relations among the given variables.

## A.7 human-annotated Causal Graphs

The human-annotated causal graphs are demonstrated in Figure 4, 5, 6, 7, 8.

# A.8 Evaluation of Causal Discovery Component

The detailed evaluation results of the causal discovery component are presented in Table 9, 10, 11, 12, 13.





(b) Human

Figure 4: Illustration of expanded causal graphs for Cancer. Squared nodes represent initial variables, while round nodes denote new proposed variables.

|   | Given a document: {doc}  |
|---|--|
| l |  |
| l | Please complete the below task.  |
| l | We have some given variables: '{observed_variables}'.  |
| l | What are the high-level variables in the provided document that have causal relations to variables in the given variable set?                      |
| l | Please form the answer using the following format.   |
| l | First, propose as many variables as possible that have causal relationships with the given variables, based on your understanding of the document. |
| l | Please ensure these proposed variables are different from the ones already provided.   |
| l | Next, refine your list of candidate variables by reducing semantic overlap among them and shortening their names for clarity.                      |
| l | Finally, determine the most reliable variable candidate as the final answer using the template provided below:                                     |
| l | The new abstracted variable is <var></var> .   |
|   |  |
|   | Table 6: The prompt for missing variable abstraction.  |
|   |  |

The task is to determine the cause-effect relation between two variables. The variables are: {variable1} and {variable2}. The answer should be {variable1} ->{variable2} or {variable1} <- {variable2} or no causal relation. Let's provide a step-by-step process to analysis the relation between them, then provide your final answer using the following format: The final answer is \_\_\_\_.

Table 7: The prompt used in the baseline for evaluation of expanded causal graphs.



Figure 5: Illustration of expanded causal graphs for Respiratory Disease. Squared nodes represent initial variables, while round nodes denote new proposed variables.



Figure 6: Illustration of expanded causal graphs for Diabetes. Squared nodes represent initial variables, while round nodes denote new proposed variables.

### **Causal Relation Annotation Task**

### Task overview:

Your task is to identify and annotate causal relations among a set of variables. A causal relation exists when one variable directly influences another.

#### Instructions:

1. Consider each pair of variables and determine if there is a direct causal relationship between them.

- 2. If you believe variable A causes variable B, indicate this as:  $A \rightarrow B$
- 3. Be cautious of confusing correlation with causation. Only mark a relationship if you believe there is a direct causal link.
- 4. Consider the direction of causality carefully. For example, "Obesity  $\rightarrow$  Heart Failure" suggests obesity causes heart failure, not the other way around.
- 5. It's okay to have multiple causes for a single effect, or multiple effects from a single cause.
- 6. Not all variables will necessarily have causal relationships with others.
- 7. Use your best judgment based on available knowledge and logical reasoning.

### **Examples:**

lifestyle -> obesity heart defect -> cardiac output genetic disorder -> heart defect

#### Submission:

Please submit your annotations as a list of causal relations in the format: Variable A -> Variable B Thank you for your careful consideration of this task!

### Task 1: Cancer

### Variables:

pollution smoker cancer x-ray dyspnoea air quality education health issues toxicity chronic illness covid-19 inflammation respiratory issues immunity carcinogens early detection

### **Causal Relations:**

### Table 8: Instructions and interface of causal relation annotation task.

| Cancer (4 edges, 5 nodes)  |           |        |      |                      |            |  |  |  |
|----------------------------|-----------|--------|------|----------------------|------------|--|--|--|
| Method                     | Precision | Recall | F1↑  | # of predicted edges | NHD Ratio↓ |  |  |  |
| Pairwise-LLM               | 0.75      | 0.75   | 0.75 | 4                    | 0.25       |  |  |  |
| BFS-LLM                    | 0.6       | 0.75   | 0.67 | 5                    | 0.33       |  |  |  |
| COAT                       | 0.13      | 0.25   | 0.17 | 8                    | 0.83       |  |  |  |
| IRIS- GES                  | 0.25      | 0.5    | 0.33 | 8                    | 0.67       |  |  |  |
| IRIS- NOTEARS              | 1.0       | 0.25   | 0.4  | 1                    | 0.6        |  |  |  |
| IRIS- PC                   | 0.14      | 0.25   | 0.18 | 7                    | 0.82       |  |  |  |
| IRIS- VCR                  | 1.0       | 0.75   | 0.86 | 3                    | 0.14       |  |  |  |
| IRIS (Llama) - NOTEARS+VCR | 0.375     | 0.75   | 0.5  | 8                    | 0.5        |  |  |  |
| IRIS- NOTEARS+VCR          | 1.0       | 0.75   | 0.86 | 3                    | 0.14       |  |  |  |

Table 9: Evaluation results of causal discovery on cancer graph. VCR refers to verified causal relations that are extracted from and validated by relevant academic documents. "Llama" refers to the use of the Llama-3.1-8b-instruct model as a substitute for GPT-40 in our method.

| Respiratory Disease (5 edges, 4 nodes) |           |        |      |                      |            |  |  |  |
|--|-----------|--------|------|----------------------|------------|--|--|--|
| Method                                 | Precision | Recall | F1↑  | # of predicted edges | NHD Ratio↓ |  |  |  |
| Pairwise-LLM                           | 1.0       | 0.6    | 0.75 | 3                    | 0.25       |  |  |  |
| BFS-LLM                                | 0.67      | 0.4    | 0.5  | 3                    | 0.5        |  |  |  |
| COAT                                   | 1.0       | 0.8    | 0.89 | 4                    | 0.11       |  |  |  |
| IRIS- GES                              | 1.0       | 0.8    | 0.89 | 4                    | 0.11       |  |  |  |
| IRIS- NOTEARS                          | 1.0       | 0.2    | 0.33 | 1                    | 0.67       |  |  |  |
| IRIS- PC                               | 0.83      | 1.0    | 0.91 | 6                    | 0.09       |  |  |  |
| IRIS- VCR                              | 1.0       | 0.8    | 0.89 | 4                    | 0.11       |  |  |  |
| IRIS (Llama) - PC+VCR                  | 1.0       | 0.8    | 0.89 | 4                    | 0.11       |  |  |  |
| IRIS- PC+VCR                           | 0.83      | 1.0    | 0.91 | 6                    | 0.09       |  |  |  |

Table 10: Evaluation results of causal discovery on respiratory disease graph.

| Diabetes (5 edges, 4 nodes) |           |        |      |                      |            |  |  |  |
|-----------------------------|-----------|--------|------|----------------------|------------|--|--|--|
| Method                      | Precision | Recall | F1↑  | # of predicted edges | NHD Ratio↓ |  |  |  |
| Pairwise-LLM                | 0.67      | 0.4    | 0.5  | 3                    | 0.5        |  |  |  |
| BFS-LLM                     | 0.67      | 0.4    | 0.5  | 3                    | 0.5        |  |  |  |
| COAT                        | 0.25      | 0.2    | 0.22 | 4                    | 0.78       |  |  |  |
| IRIS- GES                   | 0.5       | 0.6    | 0.55 | 6                    | 0.45       |  |  |  |
| IRIS- NOTEARS               | 0         | 0      | 0    | 0                    | 1.0        |  |  |  |
| IRIS- PC                    | 0.25      | 0.2    | 0.22 | 4                    | 0.78       |  |  |  |
| IRIS- VCR                   | 1.0       | 0.2    | 0.33 | 1                    | 0.67       |  |  |  |
| IRIS (Llama) - GES+VCR      | 0.67      | 0.4    | 0.5  | 3                    | 0.5        |  |  |  |
| IRIS- GES+VCR               | 1.0       | 0.6    | 0.75 | 3                    | 0.25       |  |  |  |

Table 11: Evaluation results of causal discovery on diabetes graph.

| Obesity (5 edges, 4 nodes) |           |        |      |                      |            |  |  |
|----------------------------|-----------|--------|------|----------------------|------------|--|--|
|                            | Precision | Recall | F1↑  | # of predicted edges | NHD Ratio↓ |  |  |
| Pairwise-LLM               | 0.83      | 1.0    | 0.91 | 6                    | 0.09       |  |  |
| BFS-LLM                    | 0.6       | 0.6    | 0.6  | 5                    | 0.4        |  |  |
| COAT                       | 0.25      | 0.2    | 0.22 | 4                    | 0.78       |  |  |
| IRIS-GES                   | 0.25      | 0.2    | 0.22 | 4                    | 0.78       |  |  |
| IRIS- NOTEARS              | 0         | 0      | 0    | 2                    | 1.0        |  |  |
| IRIS- PC                   | 0.25      | 0.2    | 0.22 | 4                    | 0.78       |  |  |
| IRIS- VCR                  | 1.0       | 1.0    | 1.0  | 5                    | 0          |  |  |
| IRIS (Llama) - PC+VCR      | 0.83      | 1.0    | 0.91 | 6                    | 0.09       |  |  |
| IRIS- PC+VCR               | 1.0       | 1.0    | 1.0  | 5                    | 0          |  |  |

Table 12: Evaluation results of causal discovery on obesity graph.

| ADNI (7 edges, 8 nodes)    |           |        |      |                      |            |  |  |  |  |
|----------------------------|-----------|--------|------|----------------------|------------|--|--|--|--|
| Method                     | Precision | Recall | F1↑  | # of predicted edges | NHD Ratio↓ |  |  |  |  |
| Pairwise-LLM               | 0.5       | 0.14   | 0.22 | 2                    | 0.78       |  |  |  |  |
| BFS-LLM                    | 0.33      | 0.14   | 0.2  | 3                    | 0.8        |  |  |  |  |
| COAT                       | 0.11      | 0.14   | 0.13 | 9                    | 0.87       |  |  |  |  |
| IRIS- GES                  | 0.08      | 0.14   | 0.11 | 12                   | 0.89       |  |  |  |  |
| IRIS- NOTEARS              | 0.33      | 0.14   | 0.2  | 3                    | 0.8        |  |  |  |  |
| IRIS- PC                   | 0.11      | 0.14   | 0.13 | 9                    | 0.87       |  |  |  |  |
| IRIS- VCR                  | 0.4       | 0.29   | 0.33 | 5                    | 0.67       |  |  |  |  |
| IRIS (Llama) - NOTEARS+VCR | 0.08      | 0.14   | 0.11 | 12                   | 0.89       |  |  |  |  |
| IRIS- NOTEARS+VCR          | 0.38      | 0.43   | 0.4  | 8                    | 0.6        |  |  |  |  |

Table 13: Evaluation results of causal discovery on ADNI graph.

Algorithm 2 Hybrid Causal Discovery

**Require:** Initial variables  $\mathbb{Z}$ , LLM M, structured data X, prompt l, hyperparameters  $\alpha, \beta$ **Statistical Causal Discovery**  $\hat{\mathcal{G}}_s \leftarrow \text{causal\_discovery\_alg}(\mathbb{X}) \triangleright \text{Apply causal}$ discovery algorithms (e.g., PC algorithm)

# **Causal Relation Verification**

 $\mathcal{G}_v \leftarrow \text{causal graph with no edges}$ remove\_edges  $\leftarrow \emptyset$ for each  $z_i$  in  $\mathbb{Z}$  do for each  $z_i$  in  $\mathbb{Z}$  do if  $z_i \neq z_j$  then  $r \leftarrow "z_i \text{ causes } z_i"$  $veracity_r \leftarrow \emptyset$ ▷ Initialize the veracity list for relation rfor each d in  $\mathbb{D}_{z_i, z_i}$  do  $\triangleright \mathbb{D}_{Z_i,Z_i}$ denotes documents containing both  $z_i$  and  $z_j$  $ver_d \leftarrow \boldsymbol{M}(l(r,d))$ Determine the veracity of r based on document d $veracity_r$  $\leftarrow veracity_r \cup$  $\{ver_d\}$ end for if  $veracity_r.count(True) > \alpha \times$  $len(veracity_r)$  then  $\hat{\mathcal{G}}_v \leftarrow \hat{\mathcal{G}}_v \cup \{r\} \triangleright \text{Add relation } r$ to the causal graph  $\hat{\mathcal{G}}_v$ else if  $veracity_r.count(False) >$  $\beta \times len(veracity_r)$  then remove\_edges  $\leftarrow$ remove\_edges  $\cup \{r\}$ end if end if end for end for Merge  $\hat{\mathcal{G}}_s$  and  $\hat{\mathcal{G}}_v$ for each edge r in  $\hat{\mathcal{G}}_v$  do  $\hat{\mathcal{G}}_s \leftarrow \hat{\mathcal{G}}_s \cup \{r\}$   $\triangleright$  Add relation r to  $\hat{\mathcal{G}}_s$ end for for each edge r in remove\_edges do  $\hat{\mathcal{G}}_s \leftarrow \hat{\mathcal{G}}_s \setminus \{r\} \triangleright$  Remove relation r from  $\hat{\mathcal{G}}_s$ if it exists end for

 $\hat{\mathcal{G}} \leftarrow \hat{\mathcal{G}}_s$  $\triangleright$  The final merged causal graph **Output:**  $\hat{\mathcal{G}}$ 

Algorithm 3 Missing Variable Proposal

**Require:** Initial variables  $\mathbb{Z}$ , LLM M, collected documents  $\mathbb{D}$ , prompt l, hyperparameter  $\alpha$ Step 1: Abstraction of Missing Variable Candidates  $\mathbb{Z}_c \leftarrow \emptyset$ ▷ Initialize the set of candidates for each document d in  $\mathbb{D}$  do  $\mathbf{z} \leftarrow \boldsymbol{M}(l(\mathbb{Z}, d))$ ▷ Abstract a candidate variable from document d  $\mathbb{Z}_c \leftarrow \mathbb{Z}_c \cup \{z\}$ end for Step 2: Missing Variable Proposal Based on Verified Causal Relations  $\mathbb{Z}_m \leftarrow \emptyset \triangleright$  Initialize the set of missing variables for each variable  $z_i$  in  $\mathbb{Z}_c$  do for each given variable  $z_i$  in  $\mathbb{Z}$  do  $r_1 \leftarrow "z_i \text{ causes } z_i"$  $veracity_{r_1} \leftarrow \emptyset \triangleright$  Initialize the veracity list for relation  $r_1$ for each document d in  $\mathbb{D}_{z_i, z_i}$  do  $\mathbb{D}_{z_i,z_i}$  denotes documents containing both  $z_i$  and  $\mathbf{z}_j$  $ver_d \leftarrow \boldsymbol{M}(l(r_1, d)) \quad \triangleright \text{ Determine}$ the veracity of r1 based on document d $veracity_{r_1} \leftarrow veracity_{r_1} \cup \{ver_d\}$ end for if veracity<sub>r1</sub>.count(True) >  $\alpha \times$  $veracity_{r_1}.count(False)$  then  $\mathbb{Z}_m \leftarrow \mathbb{Z}_m \cup \{z_i\} \quad \triangleright \operatorname{Add} z_i \text{ to the}$ set of proposed variables end if  $r_2 \leftarrow "z_j \text{ causes } z_i"$  $\triangleright$  Repeat the process for the reverse causal relation . . . end for end for Step 3: Missing Variable Proposal Based on **Statistical Methods** ▷ Initialize a set for PMI scores  $\mathbb{S} \leftarrow \emptyset$ 

for each variable  $z_i$  in  $\mathbb{Z}_c$  do  $s_i \leftarrow \emptyset$ for each given variable  $z_i$  in  $\mathbb{Z}$  do  $s_{ij} \leftarrow PMI(z_i, z_j) \triangleright Compute PMI of$  $(z_i, z_j)$  by Equation 1  $s_i \leftarrow s_i \cup \{s_{ij}\}$ end for  $\mathbb{S} \leftarrow \mathbb{S} \cup \{\sum(s_i)\}$ ▷ Aggregate the PMI scores for  $z_i$ end for  $\mathbb{Z}_m \leftarrow \mathbb{Z}_m \cup \text{top-k}(\mathbb{S}, \mathbb{Z}_c)$  $\triangleright$  Select the top-k

variables based on their PMI scores

**Output:**  $\mathbb{Z}_m \triangleright$  Return the final set of proposed missing variables



Figure 7: Illustration of expanded causal graphs for Obesity. Squared nodes represent initial variables, while round nodes denote new proposed variables.



(b) Human

Figure 8: Illustration of expanded causal graphs for ADNI. Squared nodes represent initial variables, while round nodes denote new proposed variables.