

Discovery of the Hidden World with Large Language Models

Anonymous Authors¹

Abstract

Despite progress in the past decades, it has been a major impediment to broader real-world applications of causal methods for relying on high-quality measured variables given by human experts. In this position paper, we argue that the rise of large language models (LLMs) trained on massive observations of the world has great potential to bridge the gap. In particular, we envision a system consisting of two *mutually beneficial* components: LLMs and causal discovery methods (CDs). On the one hand, LLMs that learn rich world knowledge about the world can assist with discovering high-level hidden variables from low-level observational data. CDs, on the other hand, can uncover the underlying causal relations between the discovered high-level variables with guarantees. More importantly, CDs can also provide feedback to improve the identification of the variables. The combination of LLMs and CDs iteratively improves the discovery of the hidden causal world, and hence opens up a broader adoption of various causal methods. To provide justifications for our envision, we also introduce a preliminary implementation of the system **Causal representatiOn AssistanT** (COAT). We present a comprehensive analysis of COAT with four case studies ranging from analysis of human reviews to diagnosis of neuropathic and brain tumors. The effectiveness of COAT demonstrates a promising direction toward foundation models that empower various downstream applications of causal methods and ultimately the general artificial intelligence.

1. Introduction

Despite the progress in the past decades, existing causal discovery algorithms mainly rely on high-quality measured

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

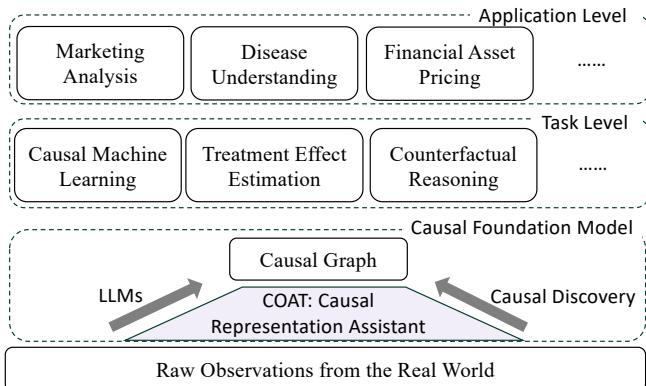


Figure 1. COAT combines both strengths of LLMs that learn the rich knowledge of the world, and causal discovery methods to uncover the hidden causal world. With the uncovered causal knowledge, COAT can empower broader applications of causal methods.

variables given by human experts (Spirtes et al., 2000; 2010; Vowels et al., 2022). However, the causal variables and their measurements are usually available in a wide range of real-world applications. For example, Amazon sellers who want to analyze the factors related to user ratings only have user reviews, generated by the underlying user preferences for certain product characteristics. Therefore, the lack of measured high-quality causal variables has been a major impediment to broader real-world applications of causal or causality-inspired methods (Schölkopf et al., 2021).

In this position paper, we argue that the recent emergence of Large Language Models (LLMs) (Brown et al., 2020; OpenAI, 2022; Touvron et al., 2023a; OpenAI, 2023) has a great potential to mitigate the gap. In particular, as shown in Fig. 1, we envision an entangled system towards causal foundation models, which consists of two *mutually beneficial* components: LLMs and causal discovery methods (CDs). On the one hand, LLMs that learn rich world knowledge about the world can assist with discovering high-level hidden variables from low-level observational data (Bubeck et al., 2023), or even certain commonsense causal knowledge (Zhang et al., 2023a; Kiciman et al., 2023; Anonymous, 2024). on the other hand, it has also been found that the reliability of LLMs in reasoning of causality remains a debate (Zečević et al., 2023; Jin et al., 2023a;b; Zhang et al., 2023a), due to a series of drawbacks of LLMs (Zhang et al.,

055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075

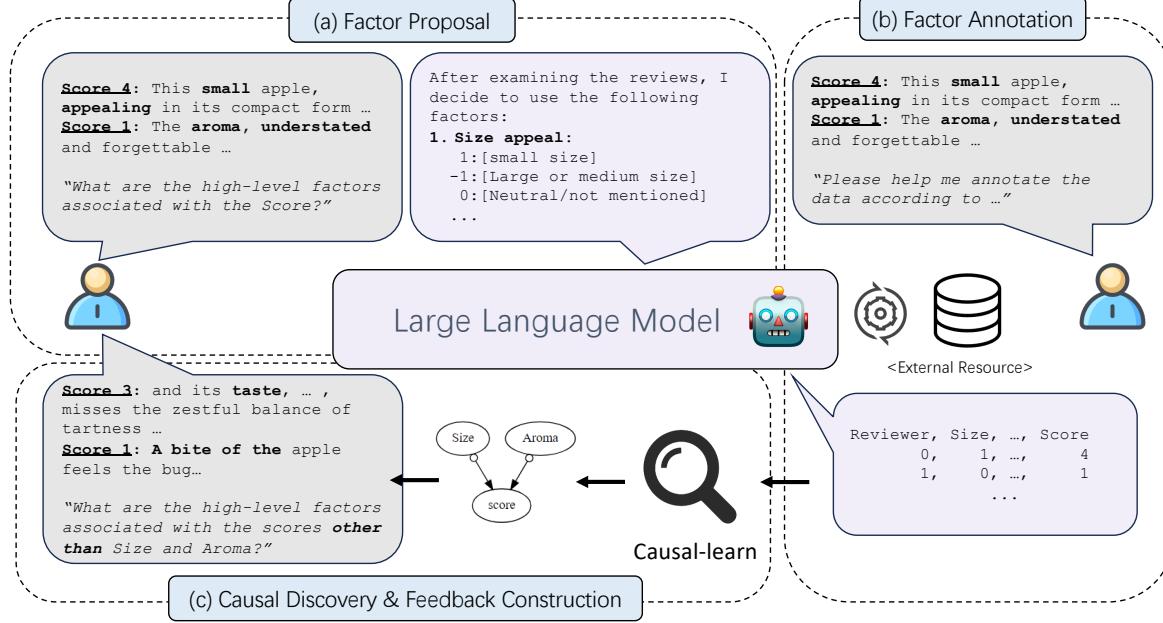


Figure 2. Illustration of COAT framework to analyze the rating scores of AppleGastronome. COAT aims to uncover the underlying Markov Blanket with respect to the given ratings of the apples (i.e., factors that fit the preferences of gastronomes). COAT first (a) adopts an LLM to read, comprehend, and relate the rich knowledge about tasting the apples. The LLM needs to propose a series of candidate factors such as apple sizes and smells, along with some meta-information such as annotation guidelines. Based on the candidate factors, COAT then (b) prompts another LLM to annotate the unstructured review into structured data. (c) The CD algorithm then finds causal relations among the factors, and constructs feedback based on samples where the ratings can not be well explained by the existing factors. By looking into the new samples, the LLM is expected to associate more related knowledge to uncover more desired causal factors.

2023c; Cui et al., 2023; Berglund et al., 2023), as well as the ethical considerations (Lee et al., 2023; Tu et al., 2023).

Nevertheless, CDs can uncover the underlying causal relations between the discovered high-level variables with guarantees. More importantly, CDs can also provide feedback to improve the identification of the variables. The combination of LLMs and CDs iteratively improves the discovery of the hidden causal world, and hence opens up a broader adoption of various causal methods.

To provide justifications for our envision, we also introduce a preliminary implementation of the system **Causal representatiOn AssistanT (COAT)**. Specifically, we will be focusing on discovering the Markov Blanket with respect to a high-value target variable from the unstructured observational data. When given a new task, such as identifying causally related high-level variables about the rating score of an apple from multiple gastronomes, shown as in Fig. 2, COAT first leverages an LLM to relate its acquired knowledge about the task, and comprehend the sample inputs to propose candidate causal factors such as the apple size and smells associated with certain meta-information (e.g., annotation guidelines). Then, COAT can either fetch the data of the candidate factors from external tools or resources when available, or employ another LLM to annotate the raw observational data according to the guidelines. Afterward, the

annotated data is readily input to the suitable CD algorithm to obtain the intermediate causal hypothesis. Throughout the paper, we mainly adopt the FCI algorithm (Spirtes et al., 2000) as it allows the existence of latent confounders and is compatible with our objective to find the Markov Blanket.

Nevertheless, LLMs often require proper prompts to unlock their full capability and elicit high-quality outputs (Wei et al., 2022). Meanwhile, LLMs may also give certain unreliable factors or miss necessary factors (Bubeck et al., 2023; Zhang et al., 2023c). Therefore, COAT also leverages the CD algorithm to audit the candidate causal factors and filter out inappropriate factors. For example, one could filter out factors that are independent of the target variable. Meanwhile, COAT can also construct feedback to further improve the factor discovery by LLMs. For example, one could leverage samples where the current factors can not explain the target variable very well, to prompt LLMs to look into the samples and find better factors.

We present a comprehensive analysis of COAT with four case studies ranging from analysis of human reviews to diagnosis of neuropathic and brain tumors (see more details in Appendix B.7). The experiments demonstrate that with COAT, LLMs can effectively identify the underlying causal variables, and provide reliable annotations for CDs

110 to uncover the causal relations. The strong performances
 111 of COAT not only provide preliminary verification of our
 112 envision, but more importantly, demonstrate a promising
 113 direction toward foundation models that empower various
 114 downstream applications of causal methods and ultimately
 115 the general artificial intelligence.

116

117 2. Related Work

118 We discuss closely related work in this section and leave
 119 more details in Appendix A.

120 **Causal discovery** aims to discover the *unknown* causal relations
 121 from the observational data (Spirtes et al., 2010; 2000),
 122 which is critical to both real-world applications and scientific
 123 discoveries (Pearl & Robins, 1999; Pearl & Mackenzie,
 124 2018). Despite recent theoretical and empirical improvements
 125 (Glymour et al., 2019; Vowels et al., 2022), most
 126 existing causal discovery approaches rely on well-structured
 127 data with human-crafted factors as inputs, and can easily
 128 suffer from the low quality of annotated data (e.g., latent
 129 confounders) (Dong et al., 2023). In this work, we argue
 130 that incorporating LLMs that learn world knowledge from
 131 massive training data could effectively relieve the need for
 132 artificial causal factor annotation. Meanwhile, COAT also
 133 opens up a new line to learn causal representations with
 134 rich pre-trained knowledge as well as feedback from causal
 135 discovery algorithms (Schölkopf et al., 2021).

136 **Reasoning with LLMs** has achieved remarkable performance
 137 across a variety of tasks with few demonstrations of
 138 the samples (Brown et al., 2020; OpenAI, 2022; Ouyang
 139 et al., 2022; OpenAI, 2023). The strong capabilities
 140 of LLMs show that it is evident that LLMs could ac-
 141 quire and understand commonsense knowledge about the
 142 world (Bubeck et al., 2023). The power of LLMs can
 143 be further unlocked with suitable context as inputs (Wei
 144 et al., 2022). Nevertheless, LLMs have also been found to
 145 make mistakes in basic algorithmic reasoning (Cobbe et al.,
 146 2021; Dziri et al., 2023), easy to hallucinate nonfactual re-
 147 sults (Zhang et al., 2023c), tend to learn shortcuts or dataset
 148 biases (Liu et al., 2023a; Berglund et al., 2023; Cui et al.,
 149 2023), and fall short in complex planning and reasoning
 150 tasks (Bubeck et al., 2023). The drawbacks of LLMs ren-
 151 der it *risky* to rely on the reasoning results to derive any
 152 rigorous results. Therefore, we do not directly derive the
 153 results from LLMs. Rather, we merely leverage the learned
 154 world knowledge in LLMs to find useful causal factors by
 155 constructing proper instructions based on causal feedback.

156 **Causal learning with LLMs** has received lots of attention
 157 results from the community (Kiciman et al., 2023; Zhang
 158 et al., 2023a). Kiciman et al. (2023) find that LLMs can re-
 159 cover the pairwise causal relations very well. Lampinen et al.
 160 (2023) show that transformer-based agents can learn causal
 161

162 strategies passively if allowed intervention during tests.
 163 Choi et al. (2022); Long et al. (2023a); Ban et al. (2023);
 164 Anonymous (2024) propose to incorporate the causal dis-
 165 covery results by LLMs as a prior or constraint to improve
 166 the performance of data-driven causal discovery algorithms.
 167 However, Willig et al. (2022); Zečević et al. (2023) find that
 168 LLMs can not understand causality while simply retelling
 169 the causal knowledge contained in the training data. Zhang
 170 et al. (2023a); Jin et al. (2023b,a) find that LLMs can hardly
 171 provide satisfactory answers for discovering new knowledge
 172 or decision-making tasks. Although Long et al. (2023b) find that
 173 LLMs can build causal graphs with 3-4 nodes, Tu et al.
 174 (2023) find that the performance of LLMs in more complex
 175 causal discovery remains limited as LLMs can hardly un-
 176 derstand new concepts and knowledge. The aforementioned
 177 debate implies the limitations in directly adopting the causal
 178 discovery results by LLMs, which motivates us to incorpo-
 179 rate the existing causal discovery algorithms with rigorous
 180 guarantees instead of LLMs to learn the causal relations.

181 Despite certain successes of indirectly or directly leveraging
 182 the LLMs to reason for the causal relations, the unreliability
 183 of LLM reasoning poses risks to trusting the causal
 184 discovery results. Different from existing approaches, we
 185 advocate for another paradigm that merely leverages LLMs
 186 to find high-level causal variables. The CD algorithm in
 187 COAT can provide necessary guarantees of the discovered
 188 causal relations, as well as feedback for maximally eliciting
 189 the rich world knowledge of LLMs. We envision that the
 190 system consists of two mutually beneficial components that
 191 provide a promising direction to develop foundation mod-
 192 els empowering broader applications of various causal or
 193 causality-based methods.

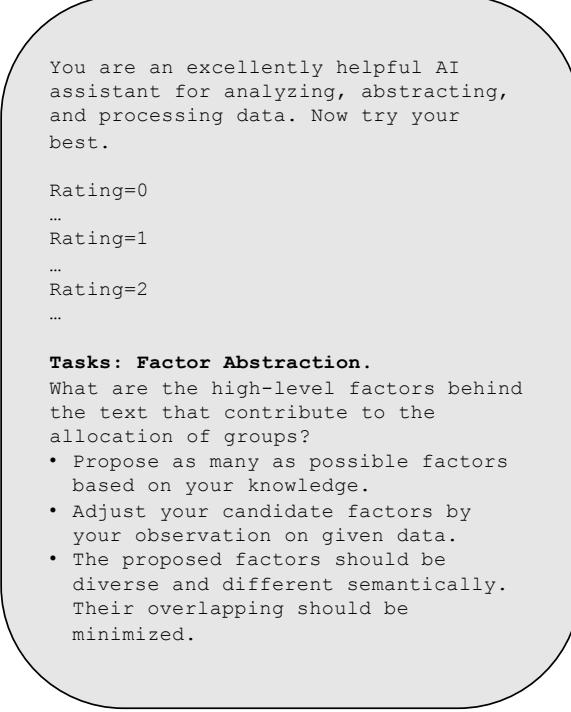
3. Causal Discovery with LLMs

In this work, we focus on discovering causation from text data while our framework could also be extended to other modalities such as images (Radford et al., 2021; Betker et al., 2023; OpenAI, 2023). Different from typical causal discovery settings with well-structured data (Spirtes et al., 2010), we start from the unstructured empirical observations where the causal variables are not available.

3.1. Problem Definition

We consider uncovering the unobserved latent variables $\mathbf{z} = (z_1, z_2, \dots, z_l)$ with respect to an observed high-value target variable y , as well as their causal relations e from the *observed unstructured* data \mathbf{x} (e.g., natural languages, images, etc.). The target variable y can be considered as the trigger to query the underlying causation (e.g., an indicator for getting on a disease, etc.). Nevertheless, our framework can also be generalized to discover the whole causal graph.

165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191



You are an excellently helpful AI assistant for analyzing, abstracting, and processing data. Now try your best.
Rating=0
...
Rating=1
...
Rating=2
...
Tasks: Factor Abstraction.
What are the high-level factors behind the text that contribute to the allocation of groups?

- Propose as many as possible factors based on your knowledge.
- Adjust your candidate factors by your observation on given data.
- The proposed factors should be diverse and different semantically. Their overlapping should be minimized.

Figure 3. Illustration of the prompt for factor proposal.

192
193
194
195
196
197
198
199

Factors behind unstructured data. We consider a broad spectrum of unstructured data, predominantly text, symbolically represented as $x \in \mathcal{X}$. This data is assumed to be generated by an unknown mapping $g : \cup\{\mathcal{Z}_i\}_{i=1}^l \times \mathcal{N} \rightarrow \mathcal{X}$, such that $x = g(z, \epsilon)$. The latent factors $z = (z_1, z_2, \dots, z_l)$ with $z_i \in \mathcal{Z}_i$ are the underlying causal variables that control the generation of x . The randomness is captured by $\epsilon \in \mathcal{N}$, where \mathcal{N} is the noise space.

200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219

Observable target variable. The observable target variable is denoted as $y \in \mathbb{R}$, and is generated as $y = f_y(\text{PA}(y), \epsilon_y)$ where $\text{PA}(y) \subseteq z$ refer to the direct parents of y in the causal graph and ϵ_y captures the noise when generating y via f_y . The observable variable y serves as a crucial element in our analysis, which triggers the causal discovery, as well as provides a measurable outcome associated with the data.

210
211
212
213
214
215
216
217
218
219

Learning latent causal structures. There exists an underlying causal structure $\mathcal{G} = (z, e)$ where z is the set of nodes representing the latent factors and e is the set that contains edges linking the nodes. The objective of COAT is, given n i.i.d. observations of $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$,

- a) to identify a set of latent factors $\hat{z} \subseteq z$ such that all of the underlying useful factors to predict y are identified $\text{MB}(y) \subseteq \hat{z}$, where $\text{MB}(y)$ refers to the Markov Blanket (Pearl, 1991) with respect to y such that $y \perp\!\!\!\perp z_i \setminus \text{MB}(y) | \text{MB}(y)$;

- b) to uncover the causal relations among $\text{MB}(y)$;

The discovery process resembles local causal discovery, which offers meaningful insights about the target variable y (Aliferis et al., 2010; Gupta et al., 2023). Nevertheless, by shifting the target variable to other identified causal factors, we could also recover the whole causal graph.

3.2. The COAT Framework

Under the aforementioned setup, we propose a new framework called COAT that leverages LLMs to propose and annotate factors given the unstructured observation. By additionally incorporating a causal discovery algorithm to audit the proposed factors and provide feedback, COAT could secure reliability when leveraging the learned knowledge by LLMs. Therefore, COAT is proceeded in an iterative manner among the two modules. The algorithm of COAT is given in Algorithm 1, which consists of four major steps.

Factor proposal. Denote an LLM as Ψ that takes k observations $\widehat{\mathcal{D}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^k$ as inputs. Usually, k is a relatively small number (i.e., $k \leq n$) due to the limited context understanding capability of LLMs (Liu et al., 2023b). Ψ will be instructed by s^t to find potential causal factors using the pre-trained knowledge from large observations of the world:

$$(\hat{z}^t, \hat{s}^t) = \Psi(\widehat{\mathcal{D}}^t, s^t), \quad (1)$$

where $\hat{z} = \{\hat{z}_1, \hat{z}_2, \dots, \hat{z}_m\}$ is the set of proposed factors, and $\hat{s} = \{\hat{s}_1, \hat{s}_2, \dots, \hat{s}_m\}$ is the set of corresponding semantic descriptions of the factors, such as the guidelines for curating the specific values of the factors. We will use the superscript t to denote the inputs and outputs at the i -th iteration: $\widehat{\mathcal{D}}^t, s^t, \hat{z}^t, \hat{s}^t$. The factor proposal stage aims to imitate the process of human experts in selecting and defining causal variables (Pearl & Mackenzie, 2018). $\widehat{\mathcal{D}}$ will be grouped according to values of y , and then Ψ will be instructed to explain the potential factors associated with the changes in y , shown as in Fig. 3. In addition, the metadata about the task such as the task description and context can also be incorporated into s if available.

Factor parsing. Once with the candidate factors, we will then collect the values of the factors from the unstructured observations. In previous works, they are usually collected from human experts according to the given factors (Spirtes et al., 2000). When without any additional information, in the factor proposal stage, Ψ will also need to give the corresponding value curation guidelines in \hat{s} . Then, another LLM Ψ_s could be leveraged to read the guidelines and parse the unstructured observations into structured data:

$$\hat{z}^t = \Psi_s(\mathcal{D}, \hat{v}^t, \hat{s}^t, s_p), \quad (2)$$

where s_p refers to the additional instruction to prompt Ψ_s to parse the observed data, and $\hat{v} = \{\hat{v}_1, \hat{v}_2, \dots, \hat{v}_n\}$ refers to the parsed values for the corresponding factors in \hat{z} .

220 **Algorithm 1** The COAT Algorithm

```

1: Required: Dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ; LLM for factor proposal  $\Psi$ ; Model for factor parsing  $\Psi_s$ ; causal discovery algorithm  $\mathcal{A}$ ; Feedback constructor  $\mathcal{F}$ ; Maximal rounds  $T$ ;
2: Random sampling  $\hat{\mathcal{D}}^0$ ;
3: Constructing  $s^0$ ;
4: while not converge and current round  $t < T$  do
5:    $(\hat{\mathbf{z}}^t, \hat{s}^t) \leftarrow \Psi(\hat{\mathcal{D}}^t, s^t)$ ; //factor proposal
6:    $\hat{\mathbf{v}}^t \leftarrow \Psi_s(\mathcal{D}, \hat{\mathbf{z}}^t, \hat{s}^t, s_p^t)$ ; //factor parsing
7:    $\hat{\mathcal{G}}^t \leftarrow \mathcal{A}(\hat{\mathbf{v}}^t \cup \{y\})$ ; //causal discovery
8:    $(\hat{\mathcal{D}}^{t+1}, s^{t+1}) \leftarrow \mathcal{F}(\hat{\mathcal{G}}^t, \mathcal{D}, s^t)$ ; //feedback
9: end while
10: return  $\mathcal{G}^T$ 

```

When the data curation of the proposed factors requires additional domain-specific knowledge/skills (e.g., intervening on the external environments) that the LLMs do not acquire, we could obtain \hat{z} through some external process (Schick et al., 2023; Xi et al., 2023). For example, studying the causes of a disease requires annotating relevant symptoms from diagnosis records and conducting additional medical checks (Tu et al., 2019). In experiments, we examine both strategies and show that COAT can effectively uncover the hidden factors under both factor parsing schemes.

Causal discovery. With the given values $\hat{v} \cup \{y\}$ associated with the candidate factors $\hat{v} \cup \{y\}$, a causal discovery algorithm \mathcal{A} can be allocated to reason about the causal hypothesis from the parsed data:

$$\widehat{\mathcal{G}}^t = \mathcal{A}(\widehat{z}^t \cup \{y\}), \quad (3)$$

where $\hat{\mathcal{G}}^t = (\hat{\mathbf{z}}^t, \hat{\mathbf{e}}^t)$ is the discovered causal hypothesis. In general, the inputs in each round may contain noises as well as the latent confounders, any causal discovery algorithms with suitable theoretical assumptions could be used for \mathcal{A} . The noises injected through LLM-based parsing may be of independent interest to the literature of causal discovery (Glymour et al., 2019; Vowels et al., 2022).

In this work, we demonstrate the idea of COAT via the FCI algorithm (Spirtes et al., 2010) as it is flexible to the functional classes of the underlying generation process, allows for the existence of latent confounders, and aligns well with our objective of discovering $\text{MB}(y)$.

Improving factor proposal with causal feedback. LLMs require proper prompts to fully unlock their capabilities (Wei et al., 2022; Huang et al., 2023; Qiao et al., 2023; Bubeck et al., 2023). In complex tasks, it usually requires a decomposition into intermediate steps (Wei et al., 2022). When it comes to factor proposing, it is also hard for LLMs to propose all factors at once. Nevertheless, from the causal discovery results, we could find useful information and thus

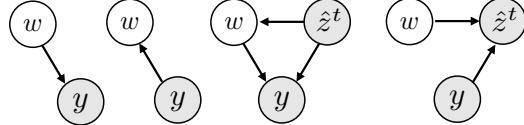


Figure 4. Illustration of variables that could be discovered with COAT. When w is the direct parent or child of y , finding hard-to-explain samples can help uncover it, with a sufficiently powerful LLM. When w is the direct parent and also a child of the discovered variable \hat{z}^t , or the spouse of y with \hat{z}^t as the common child as y , conditioning on \hat{z}^t facilitates the discovery of w .

provide feedback to further improve the factor proposal:

$$(\widehat{\mathcal{D}}^{t+1}, \mathbf{s}^{t+1}) = \mathcal{F}(\widehat{\mathcal{G}}^t, \mathcal{D}, \mathbf{s}^t), \quad (4)$$

where \mathcal{F} samples specific examples from \mathcal{D} and constructs new prompts according to the feedback from \mathcal{A} for the next round of factor proposal. For example, FCI could imply the existence of latent confounders, from which we could refine s to prompt Ψ to focus on the corresponding factors.

3.3. Causal Feedback

After setting up COAT framework, let us elaborate more on the causal feedback used in our demonstration. Since we are interested in discovering all of the relevant factors of the target variable y , let us assume there exists some underlying latent variable $w \in \text{MB}(y)$, but not discovered (i.e., $w \notin \hat{z}^t$) at round t . Then, we know w is not conditional independent with y given \hat{z}^t by the property of $\text{MB}(y)$ (Pearl, 1991) (Note that $H(\cdot)$ is the notation for entropy):

$$H(y|\hat{z}^t, w) < H(y|\hat{z}^t), \quad (5)$$

which also means incorporating w into the causal factors could facilitate the explanation and predictions of y . If the observational data contains sufficiently diverse examples, and the LLM is sufficiently powerful, Eq. 5 can help with progressively expanding the discovered factors and uncovering all the useful factors in $\text{MB}(y)$, as shown in Fig. 4.

In general, LLMs have been demonstrated to be able to resolve a variety of tasks when given a few samples ([Brown et al., 2020](#); [Bubeck et al., 2023](#)). When given the current observations \widehat{D}^t , Ψ is expected to find useful \hat{w} such that

$$H_{\widehat{\mathcal{D}}^t}(y|\hat{z}^t) - H_{\widehat{\mathcal{D}}^t}(y|\hat{z}^t, \hat{w}) \geq \Delta H_\Psi, \quad (6)$$

where $H_{\widehat{\mathcal{D}}^t}(y|\hat{z}^t)$ refers to the conditional entropy measured on $\widehat{\mathcal{D}}^t$, and ΔH_Ψ represents the sensitivity of the LLM Ψ in identifying the associated factors. Intuitively, more powerful LLMs tend to have a lower ΔH_Ψ . For the samples where Ψ can not find the desired w , it means \hat{z}^t is sufficient to capture the variances in y . To this end, we need to select suitable $\widehat{\mathcal{D}}^*$ such that

$$\widehat{\mathcal{D}}^* = \arg \max_{\widehat{\mathcal{D}} \subset \mathcal{D}} H_{\widehat{\mathcal{D}}}(y|\hat{z}^t), \quad (7)$$

which means the currently discovered causal variables \hat{z}^t can not well explain the target variable y . Recall that one of the key metric for measuring the quality of discovered Markov Blanket variables is the predictivity (Pearl, 1991). Therefore, Eq. 7 can be solved by converting it into a classification problem, where $\hat{\mathcal{D}}^*$ are the samples that the fitted classification model yields a large prediction error. In our experiments, we implement the classification as a clustering with respect to \hat{z}^t to perform sample selection of $\hat{\mathcal{D}}^*$. The clustering elicits c groups $\{\hat{\mathcal{D}}_i\}_{i=1}^c$ associated with $\{\hat{z}_i^t\}_{i=1}^c$ in the i -th group.¹

$$\{\hat{\mathcal{D}}_i\}_{i=1}^c, \{\hat{z}_i^t\}_{i=1}^c = \text{K-Means}(\mathcal{D}, \hat{z}^t). \quad (8)$$

With the groups obtained, we take the group with the largest conditional entropy to construct the feedback. The entropy of i -th group is estimated by $\hat{h}_i = -\sum \hat{p}_i(Y = y) \log \hat{p}_i(Y = y)$, where the summation is on the set of Y 's values occurred in the i -th cluster, and $\hat{p}_i(Y = y)$ is the estimated probability mass function evaluated in the $\hat{\mathcal{D}}_i$.

3.4. Practical Discussions

After establishing the main framework of COAT, we then discuss some cases where we need to handle them properly in practice, due to the instability and variety of LLM outputs.

Factor filtering. In some cases, LLMs may output several factors that have similar semantics, or even exhibit multicollinearity in the annotated data, which will hinder the causal discovery process. To mitigate the issue, one could do factor filtering, that adopts PCA or early conditional independence tests given the currently discovered variables in the Markov Blanket to detect and get rid of these variables.

Factor pool. In some cases, LLMs may discover useful factors in early rounds while may be discarded. For example, the underlying spouse variables of the target label y may be independent with y without conditioning on their common children variables. To resolve the issue, we could introduce a factor pool that stores the candidate variables proposed in the past, and replay the variables that have not been passed by conditional independence tests with existing variables in the Markov Blanket for a double check.

4. Analysis And Discussions of Feedback Driven Causal Discovery in COAT

The previous discussion leaves many aspects mysterious that the success of COAT in finding useful causal information still largely depends on the capability of LLMs. In other words, it remains unknown to what extent LLMs can find useful factors and parse the data properly, in order to find the underlying Markov Blanket of the target variable.

¹For the sake of clarity, here we denote the empirical values of \hat{z}^t within each group as \hat{z}_i^t , with a little abuse of notation.

Table 1. Causal discovery results in AppleGastronome. MB, NMB and OT refer to the number of causal factors discovered in the underlying Markov Blanket, in the causal graph but not the Markov Blanket, and the other variables. Recall, precision, and F1 for factor proposal evaluate the discovered causal ancestors.

LLM	METHOD	FACTOR PROPOSAL				
		MB	NMB	OT	RECALL	PRECISION
GPT-4	META	4	1	2	0.67	0.29
	DATA	4	1	0	1.00	0.60
	COAT	3	0	0	1.00	1.00
GPT-3.5	META	5	1	4	1.00	0.30
	DATA	5	1	0	1.00	0.50
	COAT	5	0	0	1.00	0.75
LLAMA2	META	3	0	5	0.67	0.40
	DATA	2	1	0	0.67	0.50
	COAT	2	1	0	0.67	0.80
MISTRAL	META	3	1	0	0.67	0.50
	DATA	3	1	0	1.00	0.75
	COAT	3	1	0	1.00	0.75

Thus, in this section, we construct the first benchmark under our setting, called AppleGastronome for causal discovery from unstructured data. We examine the capabilities of the predominant LLMs such as GPT-3.5-turbo (OpenAI, 2022), GPT-4-turbo (OpenAI, 2023), LLaMA2-chat (Touvron et al., 2023b), as well as Mistral-Medium (Jiang et al., 2024) in realizing COAT. For LLaMA2, we mainly evaluate the 70b variant, as we find the 13b and 7b variants can not follow the instructions well in AppleGastronome.

4.1. The AppleGastronome Benchmark

Benchmark construction. In AppleGastronome benchmark, we consider the target variable as a rating score of the apple by several gastronomes. Each apple has its own attributes, including size, smell, and taste (or sweetness). Each gastronome has a unique preference for some attributes of the apple. They will give and rating as well as write a review according to the matchness of the apple with respect to their preference. We generate the review using GPT-4 by fetching GPT-4 the preferences and the apple attributes. More details about the construction of the AppleGastronome are left in Appendix B.1. Finally, we have the AppleGastronome benchmark \mathcal{D} with (x_i, y_i) as the review and the score of the i -th gastronome, respectively. We aim to leverage COAT to identify the underlying factors that cause the rating score.

Evaluation. In AppleGastronome benchmark, since we do not have an external tool to fetch the data of the proposed factors, we employ two LLMs Ψ and Ψ_s to propose the useful factors along with the annotation criteria, and parse the unstructured review data into the structured tabular data, respectively. We evaluate both factor proposal and recovery of the underlying causal graph based on the annotated data.

Factor proposal baselines. Since most of the existing causal discovery approaches start with structured data, we

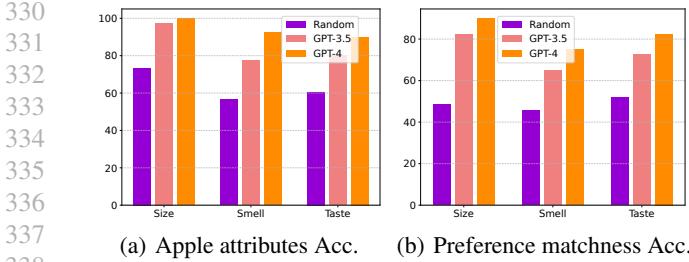


Figure 5. Annotation accuracies of GPT-4 and GPT-3.5 for apple attributes and preference matchness in AppleGastronome.

mainly compared COAT with the alternative use of using LLMs (Kiciman et al., 2023). Specifically, for the comparison of causal factor proposal, we mainly employ two different uses of LLMs, i.e., **Meta** that uses LLMs to directly output the factors given the context; **Data** that provides the same examples as the first round in COAT to the LLM.

Causal relation inference baselines. For the causal relation reasoning, we follow the previous works Kiciman et al. (2023) that prompts LLMs to reason for the causal direction of each pair of the discovered variables by **Data** baseline.

4.2. Analysis with AppleGastronome Benchmark

Can LLMs be an effective factor proposer? The causal discovery results are given in Table 3. It can be found that, compared to other uses of LLMs, COAT obtain significant improvements regardless of which LLM is used. In contrast, directly using LLMs to reason about the causal relations results in a high sensitivity to the capabilities of LLMs. When incorporating a relatively weak backbone such as LLaMA2, both **Meta** and **Data** can not find too much meaningful information than that with a stronger LLM backbone such as GPT-3.5 or GPT-4. Interestingly, when fed the examples to LLMs, all methods can generically obtain improvements. Therefore, we will directly perform the causal relation inference based on the factors proposed by **Data**.

Can LLMs be an effective factor annotator? Moreover, since LLMs are also used to annotate the data according to the annotation guidelines proposed by Ψ , we analyze the capabilities of LLMs in terms of annotation accuracy and whether the annotation will bring additional confounders that hinder the causal discovery. We mainly evaluate two LLMs GPT-3.5 and GPT-4. Fig. 5 shows the annotation accuracies of GPT-3.5 and GPT-4 when given proper guidelines. It can be found that both LLMs are generally good at annotating subjective attributes. When it comes to objective human preferences, the performance of GPT-3.5 will decrease while being relatively high.

Will LLMs introduce additional confounders in annotating factors? In addition, since the annotated results by

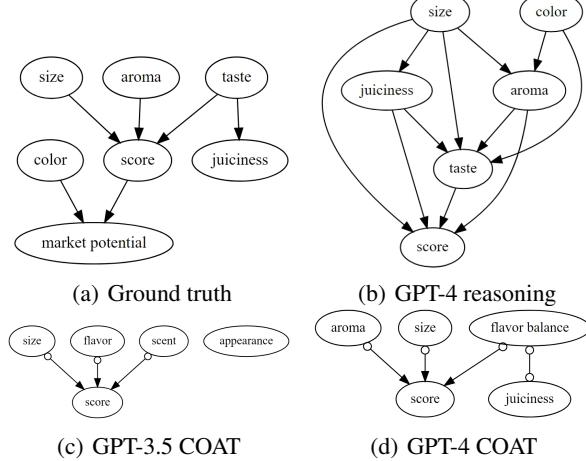


Figure 6. The discovered causal graphs in AppleGastronome. Compared to the ground truth results, directly adopting LLMs to reason about the causal relations can easily elicit lots of false positive causal relations. In contrast, the relations recovered by COAT have a high precision as well as the recall. The directed edge between “taste” and “juiciness” can not be recovered by COAT is because of the limitations of FCI.

LLMs will involve additional noises, or even additional confounders, we also conduct independence tests among the annotation noises and the features. Table 4 presents the independence testing results between the annotation noises and the annotated features, and the noises themselves. It can be found that, the introduced noises are independent of the attributes, therefore, will not introduce much additional interference to the causal discovery procedure.

Can COAT reliably recover the causal relationships? We present quantitative and qualitative results in Table 3 and Fig. 6, respectively. Compared to directly adopting LLMs to reason the causal relations, COAT significantly boosts the causal relation recovery. When considering the F1 score, even the weakest LLM LLaMA2 can outperform all the other LLMs when incorporated into COAT. GPT-3.5 is also boosted by COAT to find better causal relations. The results serve as strong evidence for the effectiveness of COAT.

5. Empirical Study with Realistic Benchmarks

After examining the capabilities of COAT in AppleGastronome, we are further motivated to challenge COAT in a more complex setting from neuronpathic panic diagnosis (Tu et al., 2019). Tu et al. (2023) have conducted tests with GPT-3.5 on the dataset and find the performance of direct reasoning with GPT-3.5 is rather limited. In this study, we construct a more realistic and challenging benchmark based on the dataset.

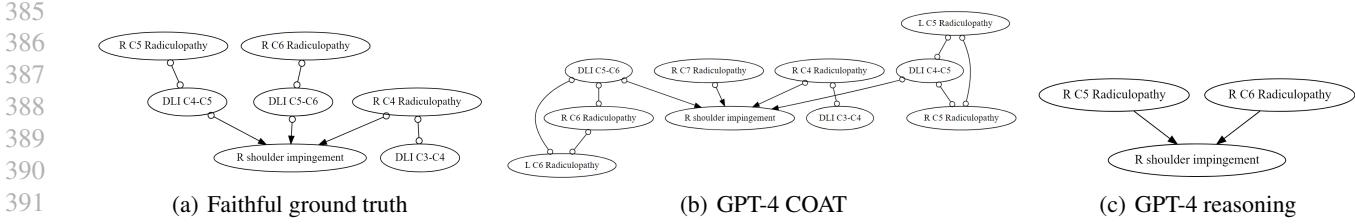


Figure 7. The discovered causal graphs in Neuropathic.

5.1. The Neuropathic Benchmark

Benchmark construction. In the Neuropathic benchmark, we convert the dataset into a clinical diagnosis task. In the original dataset, there are three levels of causal variables, including the symptom-level, radiculopathy-level and the pathophysiology-level. In experiments, we mainly consider the target variable of right shoulder impingement. When generating the clinical diagnosis notes as α using GPT-4, we will avoid any mentioning of variables other than symptoms.

As we intend to leverage the Neuropathic benchmark to simulate the real-life diagnosis, after the factor proposal stage, we directly incorporate external experts that measure the values of the candidate factors. More details about the construction of the Neuropathic are left in Appendix B.5.

Evaluation and baselines. In Neuropathic, we adopt a similar evaluation protocol and the baselines as in AppleGastronome. Nevertheless, due to the faithfulness issue of the original dataset (Tu et al., 2019), for the evaluation of causal relation discovery, we mainly conduct a qualitative comparison between the ground truth that is faithful to the data, against the baselines and COAT.

5.2. Empirical Results on Neuropathic Benchmark

Factor proposal. The quantitative results on Neuropathic benchmark are given in Table 2. Similarly, we can find that COAT consistently outperforms all of the baselines regardless of which LLMs are incorporated. In particular, COAT can boost the weakest backbone LLaMA2-7b to be better than any other LLMs.

Causal relation recovery. Fig. 7(a) shows the causal graph obtained by FCI running on the original data, where we can find that several causal relations cannot hold on the data. As shown in Fig. 7, when using LLMs to perform the reasoning, LLMs cannot identify the faithfulness issues. In contrast, COAT can imply faithful causal insights.

6. Conclusions

In this position paper, we advocate a new paradigm to incorporate LLMs towards building a causal foundation model for

Table 2. Causal discovery results in Neuropathic. PA, AN, and OT refer to the parents, ancestors, and others, respectively. Accuracy and F1 measure the recovery of the causal ancestors.

LLM	METHOD	FACTOR PROPOSAL				
		PA	AN	OT	ACC	F1
GPT-4	META	3	5	6	0.91	0.59
	DATA	2	2	0	0.95	0.50
	COAT	3	6	3	0.96	0.80
GPT-3.5	META	3	5	6	0.91	0.59
	DATA	3	5	4	0.94	0.67
	COAT	3	5	2	0.96	0.77
LLAMA2-70B	META	2	4	5	0.91	0.53
	DATA	3	3	1	0.95	0.60
	COAT	3	6	2	0.97	0.86
LLAMA2-13B	META	1	3	6	0.88	0.40
	DATA	3	6	4	0.95	0.75
	COAT	3	6	2	0.97	0.86
LLAMA2-7B	META	1	1	17	0.72	0.08
	DATA	3	6	3	0.96	0.80
	COAT	3	6	2	0.97	0.86
MISTRAL-MED	META	3	6	3	0.96	0.80
	DATA	3	3	2	0.94	0.66
	COAT	3	6	2	0.97	0.86

broader applications of various causal methods. In particular, we argue that LLMs and CDs can be mutually beneficial to mitigate the reliance on high-quality human-annotated causal variables for causal discovery. We presented a new framework called COAT to verify our vision in four challenging real-world case studies.

Looking forward, we envision future studies can be focused on two parts: (i) From the theory sides, how can we extend the previous identifiability theories to accommodate the inputs of LLMs? (ii) From the methodology side, can we design more useful feedback from CDs, and build a more efficient COAT to identify the whole causal graphs? With more sophisticated COAT implementations in the future, the discovered causal graphs will be a foundation for all causal or causality-inspired methods, that empower the exploration of the world and better living of humankind with general artificial intelligence.

440 Impact Statements

441 This work focuses on fully leveraging the rich knowledge
 442 learned by LLMs during pre-training to facilitate causal
 443 discovery from unstructured data, with the hope of em-
 444 powering broader applications and social benefits. Besides,
 445 this paper does not raise any ethical concerns. This study
 446 does not involve any human subjects, practices to data set
 447 releases, potentially harmful insights, methodologies and
 448 applications, potential conflicts of interest and sponsorship,
 449 discrimination/bias/fairness concerns, privacy and security
 450 issues, legal compliance, and research integrity issues.
 451

452 References

453 Aliferis, C. F., Statnikov, A., Tsamardinos, I., Mani, S.,
 454 and Koutsoukos, X. D. Local causal and markov blanket
 455 induction for causal discovery and feature selection for
 456 classification part i: Algorithms and empirical evalua-
 457 tion. *Journal of Machine Learning Research*, 11(7):171–234,
 458 2010. (Cited on page 4)

459 Anonymous. Causal modelling agents: Causal graph dis-
 460 covery through synergising metadata- and data-driven
 461 reasoning. In *The Twelfth International Conference on*
 462 *Learning Representations*, 2024. (Cited on pages 1, 3
 463 and 12)

464 Ban, T., Chen, L., Wang, X., and Chen, H. From query tools
 465 to causal architects: Harnessing large language models
 466 for advanced causal discovery from data. *arXiv preprint*,
 467 arXiv:2306.16902, 2023. (Cited on pages 3 and 12)

468 Berglund, L., Tong, M., Kaufmann, M., Balesni, M., Stick-
 469 land, A. C., Korbak, T., and Evans, O. The reversal curse:
 470 Llms trained on "a is b" fail to learn "b is a". *arXiv*
 471 *preprint*, arXiv:2309.12288, 2023. (Cited on pages 2, 3
 472 and 12)

473 Betker, J., Goh, G., Jing, L., TimBrooks, Wang, J., Li, L.,
 474 LongOuyang, JuntangZhuang, JoyceLee, YufeiGuo, We-
 475 samManassra, PrafullaDhariwal, CaseyChu, YunxinJiao,
 476 and Ramesh, A. Improving image generation with better
 477 captions. 2023. URL <https://cdn.openai.com/papers/dall-e-3.pdf>. (Cited on page 3)

478 Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan,
 479 J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G.,
 480 Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G.,
 481 Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu,
 482 J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M.,
 483 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish,
 484 S., Radford, A., Sutskever, I., and Amodei, D. Language
 485 models are few-shot learners. In *Advances in Neural*
 486 *Information Processing Systems*, 2020. (Cited on pages
 487 1, 3, 5 and 12)

488 Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J.,
 489 Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y.,
 490 Lundberg, S. M., Nori, H., Palangi, H., Ribeiro, M. T.,
 491 and Zhang, Y. Sparks of artificial general intelligence:
 492 Early experiments with GPT-4. *arXiv preprint*,
 493 arXiv:2303.12712, 2023. (Cited on pages 1, 2, 3, 5 and
 494 12)

495 Choi, K., Cundy, C., Srivastava, S., and Ermon, S. Lmpriors:
 496 Pre-trained language models as task-specific priors. *arXiv*
 497 *preprint*, arXiv:2210.12530, 2022. (Cited on pages 3
 498 and 12)

499 Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H.,
 500 Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano,
 501 R., Hesse, C., and Schulman, J. Training verifiers to solve
 502 math word problems. *arXiv preprint*, arXiv:2110.14168,
 503 2021. (Cited on pages 3 and 12)

504 Cui, C., Zhou, Y., Yang, X., Wu, S., Zhang, L., Zou,
 505 J., and Yao, H. Holistic analysis of hallucination in
 506 gpt-4v(ision): Bias and interference challenges. *arXiv*
 507 *preprint*, arXiv:2311.03287, 2023. (Cited on pages 2, 3
 508 and 12)

509 Dong, X., Huang, B., Ng, I., Song, X., Zheng, Y., Jin,
 510 S., Legaspi, R., Spirites, P., and Zhang, K. A versatile
 511 causal discovery framework to allow causally-related hid-
 512 den variables. *arXiv preprint*, arXiv:2312.11001, 2023.
 513 (Cited on pages 3 and 12)

514 Dziri, N., Lu, X., Sclar, M., Li, X. L., Jiang, L., Lin,
 515 B. Y., Welleck, S., West, P., Bhagavatula, C., Bras, R. L.,
 516 Hwang, J. D., Sanyal, S., Ren, X., Ettinger, A., Harchaoui,
 517 Z., and Choi, Y. Faith and fate: Limits of transformers
 518 on compositionality. In *Thirty-seventh Conference on*
 519 *Neural Information Processing Systems*, 2023. (Cited on
 520 pages 3 and 12)

521 Glymour, C., Zhang, K., and Spirtes, P. Review of causal
 522 discovery methods based on graphical models. *Frontiers*
 523 in *Genetics*, 10, 2019. (Cited on pages 3, 5 and 12)

524 Gupta, S., Childers, D., and Lipton, Z. C. Local causal
 525 discovery for estimating causal effects. In *Conference on*
 526 *Causal Learning and Reasoning*, volume 213, pp. 408–
 527 447, 2023. (Cited on page 4)

528 Huang, J., Gu, S., Hou, L., Wu, Y., Wang, X., Yu, H., and
 529 Han, J. Large language models can self-improve. In
 530 *Conference on Empirical Methods in Natural Language*
 531 *Processing*, pp. 1051–1068, 2023. (Cited on page 5)

532 Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary,
 533 B., Bamford, C., Chaplot, D. S., de Las Casas, D., Hanna,
 534 E. B., Bressand, F., Lengyel, G., Bour, G., Lample, G.,
 535 Lavaud, L. R., Saulnier, L., Lachaux, M., Stock, P., Sub-
 536 ramanian, S., Yang, S., Antoniak, S., Scao, T. L., Gervet,

- 495 T., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E.
 496 Mixtral of experts. *arXiv preprint*, arXiv:2401.04088,
 497 2024. (Cited on page 6)
- 498 Jin, Z., Chen, Y., Leeb, F., Gresele, L., Kamal, O., LYU,
 499 Z., Blin, K., Adauto, F. G., Kleiman-Weiner, M., Sachan,
 500 M., and Schölkopf, B. CLadder: A benchmark to
 501 assess causal reasoning capabilities of language models. In
 502 *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a. (Cited on pages 1, 3 and 12)
- 503 Jin, Z., Liu, J., Lyu, Z., Poff, S., Sachan, M., Mihalcea,
 504 R., Diab, M. T., and Schölkopf, B. Can large language
 505 models infer causation from correlation? *arXiv preprint*,
 506 arXiv:2306.05836, 2023b. (Cited on pages 1, 3 and 12)
- 507 Kiciman, E., Ness, R., Sharma, A., and Tan, C. Causal
 508 reasoning and large language models: Opening a new
 509 frontier for causality. *arXiv preprint*, arXiv:2305.00050,
 510 2023. (Cited on pages 1, 3, 7 and 12)
- 511 Lampinen, A. K., Chan, S. C., Dasgupta, I., Nam, A. J. H.,
 512 and Wang, J. X. Passive learning of active causal strategies
 513 in agents and language models. In *Advances in Neural Information Processing Systems*, 2023. (Cited on
 514 pages 3 and 12)
- 515 Lee, P., Bubeck, S., and Petro, J. Benefits, limits, and risks
 516 of gpt-4 as an ai chatbot for medicine. *New England
 517 Journal of Medicine*, 388(13):1233–1239, 2023. (Cited
 518 on page 2)
- 519 Liu, B., Ash, J. T., Goel, S., Krishnamurthy, A., and Zhang,
 520 C. Transformers learn shortcuts to automata. In *The
 521 Eleventh International Conference on Learning Representations*, 2023a. (Cited on pages 3 and 12)
- 522 Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua,
 523 M., Petroni, F., and Liang, P. Lost in the middle:
 524 How language models use long contexts. *arXiv preprint*,
 525 arXiv:2307.03172, 2023b. (Cited on page 4)
- 526 Long, S., Piché, A., Zantedeschi, V., Schuster, T., and
 527 Drouin, A. Causal discovery with language models as im-
 528 perfect experts. *arXiv preprint*, arXiv:2307.02390, 2023a.
 529 (Cited on pages 3 and 12)
- 530 Long, S., Schuster, T., and Piché, A. Can large lan-
 531 guage models build causal graphs? *arXiv preprint*,
 532 arXiv:2303.05279, 2023b. (Cited on pages 3 and 12)
- 533 LYU, Z., Jin, Z., Mihalcea, R., Sachan, M., and Schölkopf,
 534 B. Can large language models distinguish cause from
 535 effect? In *UAI 2022 Workshop on Causal Representation
 536 Learning*, 2022. (Cited on page 12)
- 537 OpenAI. Chatgpt. [https://chat.openai.com/
 538 chat/](https://chat.openai.com/chat/), 2022. (Cited on pages 1, 3, 6 and 12)
- 539 OpenAI. Gpt-4 technical report, 2023. (Cited on pages 1,
 540 3, 6 and 12)
- 541 Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright,
 542 C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama,
 543 K., Ray, A., et al. Training language models to fol-
 544 low instructions with human feedback. *arXiv preprint
 545 arXiv:2203.02155*, 2022. (Cited on pages 3 and 12)
- 546 Pearl, J. Probabilistic reasoning in intelligent systems -
 547 networks of plausible inference. In *Morgan Kaufmann
 548 series in representation and reasoning*, 1991. (Cited on
 549 pages 4, 5 and 6)
- 550 Pearl, J. and Mackenzie, D. *The Book of Why: The New
 551 Science of Cause and Effect*. Basic Books, Inc., USA, 1st
 552 edition, 2018. ISBN 046509760X. (Cited on pages 3, 4
 553 and 12)
- 554 Pearl, J. and Robins, J. M. Causal diagrams for epidemi-
 555 ologic research. *Epidemiology*, 10 1:37–48, 1999. (Cited
 556 on pages 3 and 12)
- 557 Qiao, S., Ou, Y., Zhang, N., Chen, X., Yao, Y., Deng, S., Tan,
 558 C., Huang, F., and Chen, H. Reasoning with language
 559 model prompting: A survey. In *Annual Meeting of the
 560 Association for Computational Linguistics (Volume 1:
 561 Long Papers)*, July 2023. (Cited on page 5)
- 562 Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G.,
 563 Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark,
 564 J., Krueger, G., and Sutskever, I. Learning transferable
 565 visual models from natural language supervision. In
 566 *Proceedings of the 38th International Conference on Ma-
 567 chine Learning*, Proceedings of Machine Learning Re-
 568 search, pp. 8748–8763, 2021. (Cited on page 3)
- 569 Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli,
 570 M., Zettlemoyer, L., Cancedda, N., and Scialom, T. Tool-
 571 former: Language models can teach themselves to use
 572 tools. *arXiv preprint*, arXiv:2302.04761, 2023. (Cited
 573 on page 5)
- 574 Schölkopf, B., Locatello, F., Bauer, S., Ke, N. R., Kalchbren-
 575 ner, N., Goyal, A., and Bengio, Y. Towards causal rep-
 576 resentation learning. *arXiv preprint*, arXiv:2102.11107,
 577 2021. (Cited on pages 1, 3 and 12)
- 578 Spirtes, P., Glymour, C., and Scheines, R. *Causation, Pre-
 579 diction, and Search, Second Edition*. Adaptive computa-
 580 tion and machine learning. MIT Press, 2000. ISBN
 581 978-0-262-19440-2. (Cited on pages 1, 2, 3, 4 and 12)
- 582 Spirtes, P., Glymour, C., Scheines, R., and Till-
 583 man, R. Automated Search for Causal Rela-
 584 tions: Theory and Practice. 2010. URL
 585 [https://kilthub.cmu.edu/articles/
 586 journal_contribution/Automated_](https://kilthub.cmu.edu/articles/journal_contribution/Automated_)

- 550 Search_for_Causal_Relations_Theory_
 551 and_Practice/6490961. (Cited on pages 1, 3, 5
 552 and 12)
- 553 Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux,
 554 M., Lacroix, T., Rozière, B., Goyal, N., Hambro, E.,
 555 Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lam-
 556 pple, G. Llama: Open and efficient foundation language
 557 models. *arXiv preprint*, arXiv:2302.13971, 2023a. (Cited
 558 on page 1)
- 559 Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A.,
 560 Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhos-
 561 ale, S., Bikell, D., Blecher, L., Canton-Ferrer, C., Chen,
 562 M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W.,
 563 Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn,
 564 A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez,
 565 V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S.,
 566 Lachaux, M., Lavril, T., Lee, J., Liskovich, D., Lu, Y.,
 567 Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog,
 568 I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi,
 569 K., Schelten, A., Silva, R., Smith, E. M., Subramanian,
 570 R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan,
 571 J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kam-
 572 badur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov,
 573 S., and Scialom, T. Llama 2: Open foundation and fine-
 574 tuned chat models. *arXiv preprint*, arXiv:2307.09288,
 575 2023b. (Cited on page 6)
- 576 Tu, R., Zhang, K., Bertilson, B. C., Kjellström, H., and
 577 Zhang, C. Neuropathic pain diagnosis simulator for
 578 causal discovery algorithm evaluation. In *Advances
 in Neural Information Processing Systems*, pp. 12773–
 579 12784, 2019. (Cited on pages 5, 7 and 8)
- 580 Tu, R., Ma, C., and Zhang, C. Causal-discovery perfor-
 581 mance of chatgpt in the context of neuropathic pain diag-
 582 nosis. *arXiv preprint*, arXiv:2301.13819, 2023. (Cited
 583 on pages 2, 3, 7 and 12)
- 584 Vowels, M. J., Camgoz, N. C., and Bowden, R. D'ya like
 585 dags? a survey on structure learning and causal discovery.
 586 *ACM Computing Survey*, 55(4), 2022. ISSN 0360-0300.
 587 (Cited on pages 1, 3, 5 and 12)
- 588 Wei, J., Wang, X., Schuurmans, D., Bosma, M., brian ichter,
 589 Xia, F., Chi, E. H., Le, Q. V., and Zhou, D. Chain of
 590 thought prompting elicits reasoning in large language
 591 models. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho,
 592 K. (eds.), *Advances in Neural Information Processing
 Systems*, 2022. (Cited on pages 2, 3, 5 and 12)
- 593 Willig, M., Zečević, M., Dhami, D. S., and Kersting, K. Can
 594 foundation models talk causality? In *UAI 2022 Workshop
 on Causal Representation Learning*, 2022. (Cited on
 595 pages 3 and 12)
- 596 Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B.,
 597 Zhang, M., Wang, J., Jin, S., Zhou, E., Zheng, R., Fan, X.,
 598 Wang, X., Xiong, L., Zhou, Y., Wang, W., Jiang, C., Zou,
 599 Y., Liu, X., Yin, Z., Dou, S., Weng, R., Cheng, W., Zhang,
 600 Q., Qin, W., Zheng, Y., Qiu, X., Huan, X., and Gui, T. The rise and potential of large language model based
 601 agents: A survey. *arXiv preprint*, arXiv:2309.07864,
 602 2023. (Cited on page 5)
- 603 Zečević, M., Willig, M., Dhami, D. S., and Kersting, K. Causal parrots: Large language models may talk causality
 604 but are not causal. *Transactions on Machine Learning
 Research*, 2023. ISSN 2835-8856. (Cited on pages 1, 3
 605 and 12)
- 606 Zhang, C., Bauer, S., Bennett, P., Gao, J., Gong, W., Hilmkil,
 607 A., Jennings, J., Ma, C., Minka, T., Pawlowski, N., and
 608 Vaughan, J. Understanding causality with large language
 609 models: Feasibility and opportunities. *arXiv preprint*,
 610 arXiv:2304.05524, 2023a. (Cited on pages 1, 3 and 12)
- 611 Zhang, Y., Fitzgibbon, B., Garofolo, D., Kota, A., Papen-
 612 hausen, E., and Mueller, K. An explainable AI approach
 613 to large language model assisted causal model auditing
 614 and development. *arXiv preprint*, arXiv:2312.16211,
 615 2023b. (Cited on page 12)
- 616 Zhang, Y., Li, Y., Cui, L., Cai, D., Liu, L., Fu, T., Huang, X.,
 617 Zhao, E., Zhang, Y., Chen, Y., Wang, L., Luu, A. T., Bi,
 618 W., Shi, F., and Shi, S. Siren's song in the AI ocean: A
 619 survey on hallucination in large language models. *arXiv
 preprint*, arXiv:2309.01219, 2023c. (Cited on pages 1, 2,
 620 3 and 12)

605 A. More Related Work

606 **Causal discovery** aims to discover the *unknown* causal relations from the observational data (Spirtes et al., 2010; 2000),
 607 which is critical to both real-world applications and scientific discoveries (Pearl & Robins, 1999; Pearl & Mackenzie, 2018).
 608 Despite recent theoretical and empirical improvements (Glymour et al., 2019; Vowels et al., 2022), most existing causal
 609 discovery approaches rely on well-structured data with human-crafted factors as inputs, and can easily suffer from the low
 610 quality of annotated data (e.g., latent confounders) (Dong et al., 2023). In this work, we show that incorporating LLMs that
 611 learn world knowledge from massive training data could effectively relieve the need for artificial causal factor annotation.
 612 Meanwhile, COAT also opens up a new line to learn causal representations with rich pre-trained knowledge as well as
 613 feedback from causal discovery algorithms (Schölkopf et al., 2021),
 614

615 **Reasoning with LLMs** has achieved remarkable performance across a variety of tasks with few demonstrations of the
 616 samples (Brown et al., 2020; OpenAI, 2022; Ouyang et al., 2022; OpenAI, 2023). The strong capabilities of LLMs show
 617 that it is evident that LLMs could acquire and understand commonsense knowledge about the world (Bubeck et al., 2023).
 618 The power of LLMs can be further unlocked with suitable context as inputs (Wei et al., 2022). Nevertheless, LLMs have
 619 also been found to make mistakes in basic algorithmic reasoning (Cobbe et al., 2021; Dziri et al., 2023), easy to hallucinate
 620 nonfactual results (Zhang et al., 2023c), tend to learn shortcuts or dataset biases (Liu et al., 2023a; Berglund et al., 2023;
 621 Cui et al., 2023), and fall short in complex planning and reasoning tasks (Bubeck et al., 2023). The drawbacks of LLMs
 622 render it *risky* to rely on the reasoning results to derive any rigorous results. Therefore, we do not directly derive the results
 623 from LLMs. Rather, we merely leverage the learned world knowledge in LLMs to find useful causal factors by constructing
 624 proper instructions based on causal feedback.
 625

626 **Causal learning with LLMs** has received lots of attention results from the community (Kiciman et al., 2023; Zhang et al.,
 627 2023a). Kiciman et al. (2023) find that LLMs can recover the pairwise causal relations very well. Lampinen et al. (2023)
 628 show that transformer-based agents can learn causal strategies passively if allowed intervention during tests. Choi et al.
 629 (2022); Long et al. (2023a); Ban et al. (2023); Anonymous (2024) propose to incorporate the causal discovery results by
 630 LLMs as a prior or constraint to improve the performance of data-driven causal discovery algorithms. However, Willig
 631 et al. (2022); Zečević et al. (2023) find that LLMs can not understand causality while simply retelling the causal knowledge
 632 contained in the training data. Zhang et al. (2023a); Jin et al. (2023b;a) find that LLMs can hardly provide satisfactory
 633 answers for discovering new knowledge or decision-making tasks. Although Long et al. (2023b) find that LLMs can build
 634 causal graphs with 3-4 nodes, Tu et al. (2023) find that the performance of LLMs in more complex causal discovery remains
 635 limited as LLMs can hardly understand new concepts and knowledge. The aforementioned debate implies the limitations in
 636 directly adopting the causal discovery results by LLMs, which motivates us to incorporate the existing causal discovery
 637 algorithms with rigorous guarantees instead of LLMs to learn the causal relations.
 638

639 LYU et al. (2022) find that it is crucial to incorporate prompts aligned with the underlying causal story for LLMs to do
 640 pairwise causal relation inference. Zhang et al. (2023b) propose to leverage LLMs to audit causal discovery results in an
 641 explainable way.
 642

643 The closest works to ours are Choi et al. (2022); Long et al. (2023a); Ban et al. (2023); Anonymous (2024) which also
 644 incorporates LLMs into the pipeline of causal discovery. Nevertheless, all of the existing combinations of LLMs and causal
 645 discovery still focus on artificially curated structured data and rely on the capability of LLMs to infer causal relations,
 646 therefore, limited in both the reliability and the utility of LLMs in causal learning.
 647

648 B. More Details about Experiments

649 B.1. More Details on Constructing AppleGastronome

650 In the AppleGastronome benchmark, we consider the target variable as a rating score of the apple by several gastronomes.
 651 Each apple has its own attributes, including size, smell, and taste (or sweetness). Each gastronome has a unique preference
 652 for some attributes of the apple. They will give and rating as well as write a review according to the matchness of the apple
 653 with respect to their preference. We generate the review using GPT-4 by fetching GPT-4 the preferences and the apple
 654 attributes.
 655

656 The prompts for generating the unstructured inputs are given in Fig. 8. The additional results on Relation Extraction are
 657 given in Table. 3.
 658

Discovery of the Hidden World with Large Language Models

LLM	METHOD	RELATION EXTRACTION		
		RECALL	PRECISION	F1
GPT-4	PAIR-WISE	0.75	0.27	0.40
	COAT	1.00	1.00	1.00
GPT-3.5	PAIR-WISE	0.80	0.40	0.53
	COAT	1.00	0.75	0.86
LLAMA2	PAIR-WISE	0.33	0.13	0.18
	COAT	0.67	0.80	0.73
MISTRAL	META	0.75	0.50	0.60
	COAT	1.00	1.00	1.00

Table 4. Independence tests of the annotation noises with annotated features and other noises AppleGastronome.

LLM	TEST OBJECT	T	P-VALUE
GPT-4	FEATURE	0.4803	0.0325
	FEATURE	0.2828	0.9997
GPT-3.5	NOISE	0.3030	0.1594
	NOISE	0.06181	0.9745

Examples of AppleGastronome are given in Fig. 9.

B.2. More Details on Prompts for AppleGastronome

The prompts for factor proposal are given in Fig. 10.

The prompts for factor annotation are given in Fig. 11, Fig. 12, Fig. 13.

The prompts for constructing feedback are given in Fig. 14.

B.3. More Details of Results on AppleGastronome

The detailed causal graph results are given from Fig. 15 to Fig. 19.

B.4. Implementation of the FCI algorithm

We use a third-party open-sourced Python library to perform the FCI algorithm: <https://causal-learn.readthedocs.io/en/latest/>

We set $\alpha = 0.05$, and `independence_test_method="fisherz"` throughout all experiments. Other parameters are kept as the default.

B.5. More Details on Constructing Neuropathic

In the Neuropathic benchmark, we convert the dataset into a clinical diagnosis task. In the original dataset, there are three levels of causal variables, including the symptom-level, radiculopathy-level and the pathophysiology-level. In experiments, we mainly consider the target variable of right shoulder impingement. When generating the clinical diagnosis notes as x using GPT-4, we will avoid any mentioning of variables other than symptoms.

As we intend to leverage the Neuropathic benchmark to simulate the real-life diagnosis, after the factor proposal stage, we directly incorporate external experts that measure the values of the candidate factors. The prompts to generate the diagnosis records are given in Fig. 20.

Examples of Neuropathic are given in Fig. 21.

B.6. More Details of Results on Neuropathic

The detailed causal graph results are given from Fig. 22 to Fig. 26.

715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737

You are a picky gastronome on apples. You are ready to **evaluate apples and write reviews**. Your writing should be **clear, solid and convincing to suppliers and customers**.

Task

Please write a short review about the evaluation results for a given apple.

Evaluation Results:

- vibrant and uniformly colored
- musty or rotten
- more sour than sweet
- dry and lacking moisture
- might not bring the expected returns and could even lead to losses

Additional Requirement:

- Combine those evaluation results into more detailed reasoning.
- Single paragraph; No quotation marks.
- Modern English
- Begin this part by a mark [Review Begin].
- no more than 60 words.

Figure 8. Illustration of prompts for generating AppleGastronome.

738 **B.7. Additional Real-World Case Studies**

739 **Brain Tumor** ([kaggle/brain-tumor-classification-mri](#)). As shown in Fig. 27: We use gpt-4-vision-preview to handle image samples. Each sample is a scanning MRI of a human brain. The interesting variable is the tumor type. Two visual factors appeared in the final causal graph. The proposed two factors, the contrast enhancement (brightness w.r.t. background) and the mass effect (pushing or displacing the surrounding brain tissue), are results of tumors' activities; and therefore can be used to determine the tumor type. These factors can be verified by the medical literature.

740
741
742
743
744
745 **News & Stock** ([kaggle/stock-price-and-news-realted-to-it](#)). As shown in Fig. 28: We analyze the potential factors behind the news and their relation to future return rates (using data before 2007-09). Each sample is the news about Microsoft at a certain date. The interesting variable is the stock's future return rate. Factor values are post-processed by moving averaging and rolling standardization with fixed time windows. The “innovation and technology focus” is identified as one potential causal factor of future return rates for the Microsoft stock and it yields the highest sharp ratio in the out-of-sample trading test using data after 2007-09.

751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769

```

770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
    ## Score: -4
    Unfortunately, this batch of apples fails to meet quality standards. Their small size, extensive
    surface blemishes, and dry texture fall short of expectations. The off-putting musty scent
    suggests potential rot, and the sour taste overpowers any sweetness. Suppliers should beware;
    these apples could result in financial disappointment if offered to discerning customers.

    This particular apple, while boasting a striking and consistent hue, falls short of expectations
    due to its unappealing small stature, an off-putting musty odor suggesting decay, and a flavor
    profile that tilts disappointingly toward sourness. The lack of juiciness further detracts from
    its desirability, making it a risky choice for suppliers and a potential loss-inducer in the
    market.

    ## Score: -3
    This apple, regrettably, fails to deliver on multiple fronts, its diminutive stature only
    overshadowed by an off-putting musty undertone that hints at decay. The fruit's unbalanced
    profile skewing toward sourness compounded by a disappointing dryness ensures a less than
    palatable experience. Given these attributes, its market performance is predictably
    underwhelming, posing a significant risk for monetary losses.

```

Figure 9. Illustration of examples in AppleGastronome.

```

797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
    You are an excellently helpful AI
    assistant for analyzing, abstracting,
    and processing data. Now try your
    best.

    Rating=0
    ...
    Rating=1
    ...
    Rating=2
    ...

```

Tasks: Factor Abstraction.

What are the high-level factors behind
the text that contribute to the
allocation of groups?

- Propose as many as possible factors
based on your knowledge.
- Adjust your candidate factors by
your observation on given data.
- The proposed factors should be
diverse and different semantically.
Their overlapping should be
minimized.

Figure 10. Illustration of the prompt for factor proposal.

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
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879

You are an excellently helpful AI
assistant for analyzing, abstracting,
and processing data. Now try your best.

Data

The reviewer comments on an apple that
are randomly picked:

{text}

Tasks

For the given sample, what is the most
appropriate facotr value based on the
criterion?

{factor_state}

**Your final output should follow this
template**:

The value is: ____.

Figure 11. Illustration of the prompt for factor annotation (part 1).

```

880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934

```

```

## Output

### Part 1. analysis

In this part, feel free to write down
the process of your considerations.

Hint
- You need to abstract, identify, and
choose suitable values. You may write
down your hypothesis.
- each value has one specific
criterion, you should choose the most
suitable one.
- **Your final output should follow this
template**:

***The value is: ____.****

direct copy helpful information from
sample that related to factor criterion,
if any.

...

```

Figure 12. Illustration of the prompt for factor annotation (part 2).

```

915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934

```

```

### Part 4. Final output

In this part, you are required to report
the value

**Your final output should follow this
template**:

***The value is: ____.****

```

Figure 13. Illustration of the prompt for factor annotation (part 3).

935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967

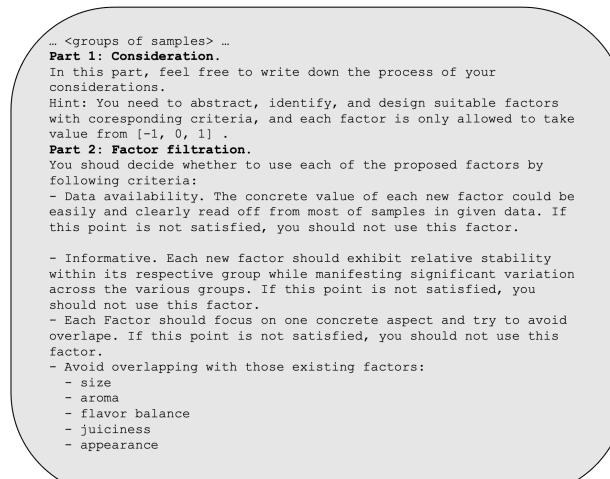


Figure 14. Illustration of the prompt for feedback.

958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989

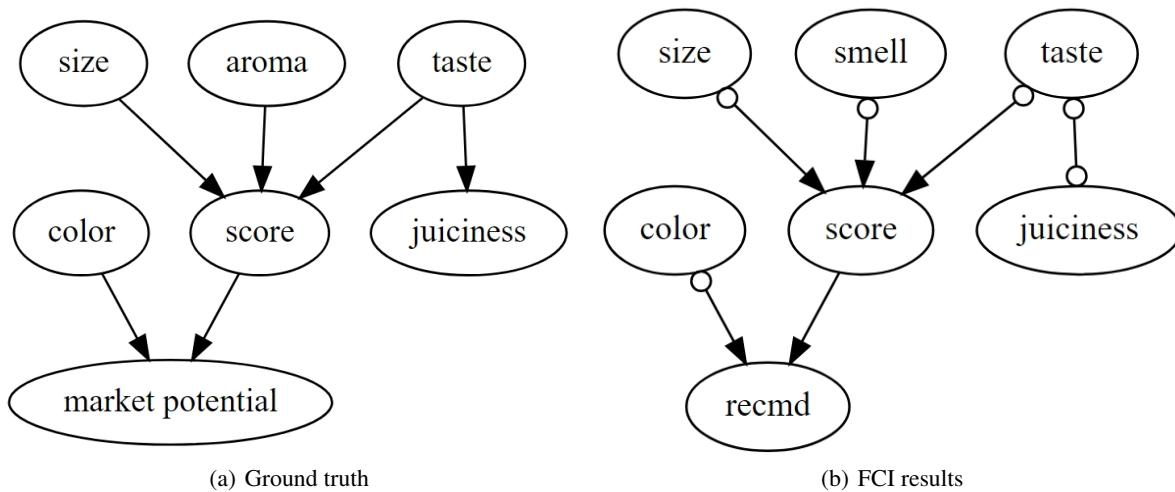
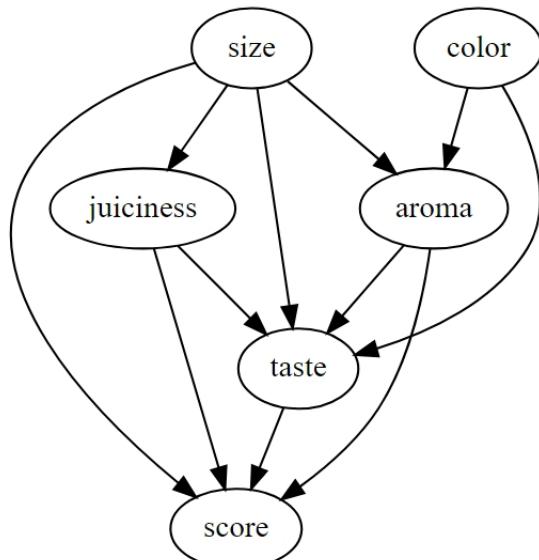
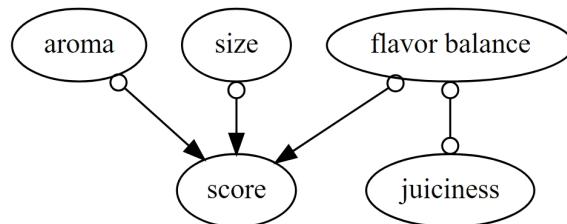


Figure 15. Ground truth and faithful (via FCI algorithm) causal graphs in AppleGastronome.

990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011



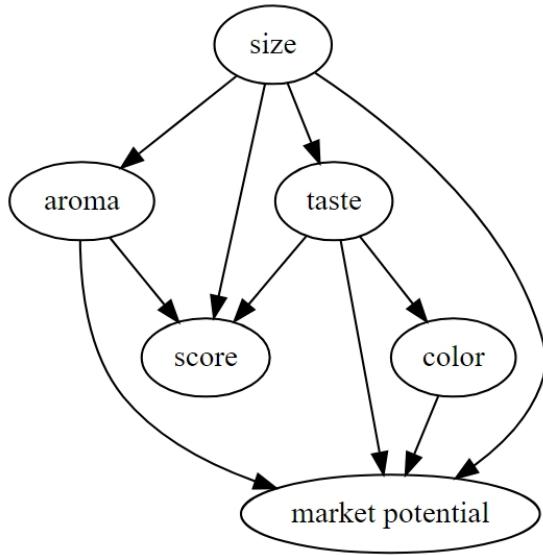
(a) GPT-4 reasoning



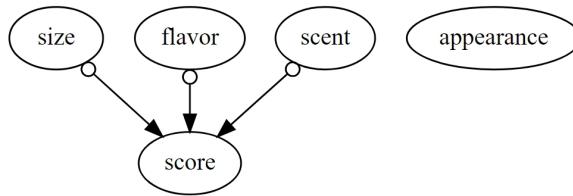
(b) GPT-4 COAT

Figure 16. Causal graphs with GPT-4 in AppleGastronome.

1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044



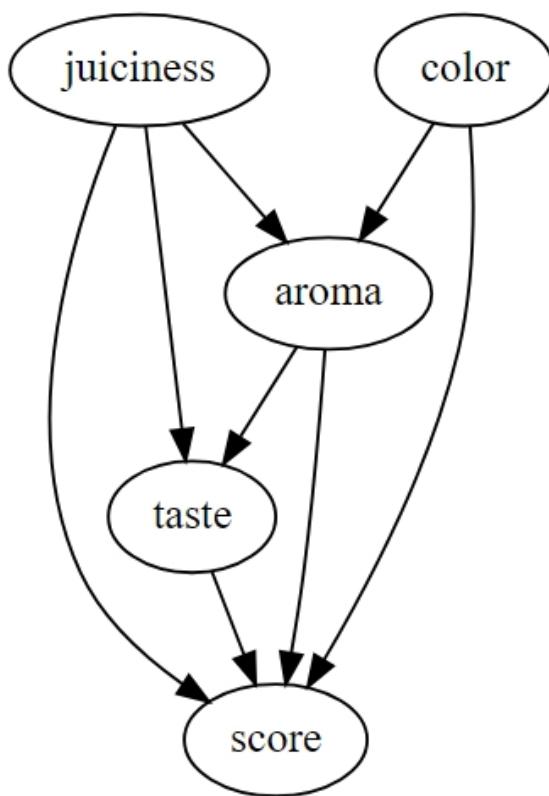
(a) GPT-3.5 reasoning



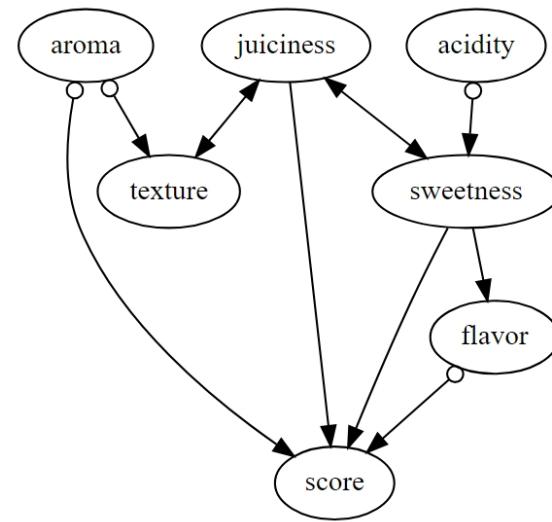
(b) GPT-3.5 COAT

Figure 17. Causal graphs with GPT-3.5 in AppleGastronome.

1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071



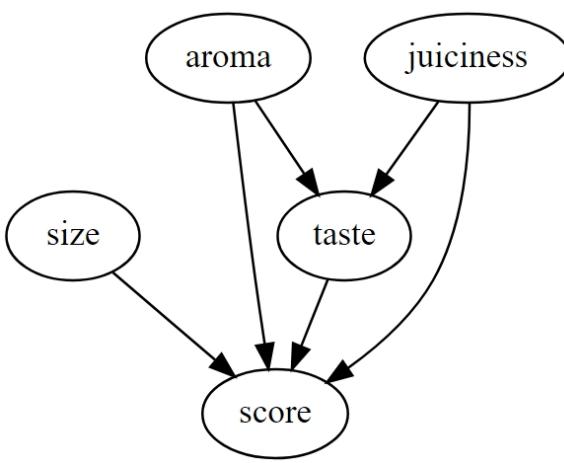
(a) LLaMA2 reasoning



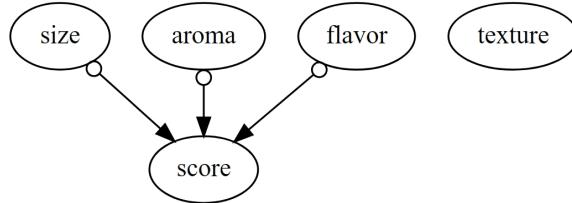
(b) LLaMA2 COAT

Figure 18. Causal graphs with GPT-3.5 in AppleGastronome.

1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099



(a) Mistral reasoning



(b) Mistral COAT

Figure 19. Causal graphs with GPT-3.5 in AppleGastronome.

```

1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154

```

You are a helpful and professional assistant to Neurologists. Try your best!

Write a clinical notes based on following informaiton:

```

Patient id: Neuro_{idx:05}
Age: {10+np.random.choice(50)}
Gender: {get_gender()}

---
```

Write one summary (one paragraph, do not do inference) about **symptom diagnosis**:
{symptom_notes}

```

---
```

L: Left
R: Right
DLI: discoligamentous injury
Radi: radiculopathy

```

---
```

Your clinical notes should surrounded by <note> ... </note>

Figure 20. Illustration of prompts for generating Neuropathic.

```

1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154

```

```

## Group with 'R shoulder impingement'=0
-----
Patient ID: Neuro_00000
Age: 57
Gender: Male

Symptom Diagnosis Summary:
The patient presents with a left shoulder impingement and associated intracapsular problems. The symptoms include limited and painful range of motion with a positive outcome when performing a "thumbs up" test on the left side. Right upper limb assessment also reveals issues; the patient wears an armband suggesting ongoing management for right lower arm disorders. The patient reports lumbago and pain localized to the left thigh. Physical therapy assessment (L pta) is noted, indicating involvement in a therapeutic intervention. Joint assessments reveal a left hip joint condition and right hip arthritis, accompanied by right hamstring discomfort.

```

Figure 21. Illustration of examples in AppleGastronome.

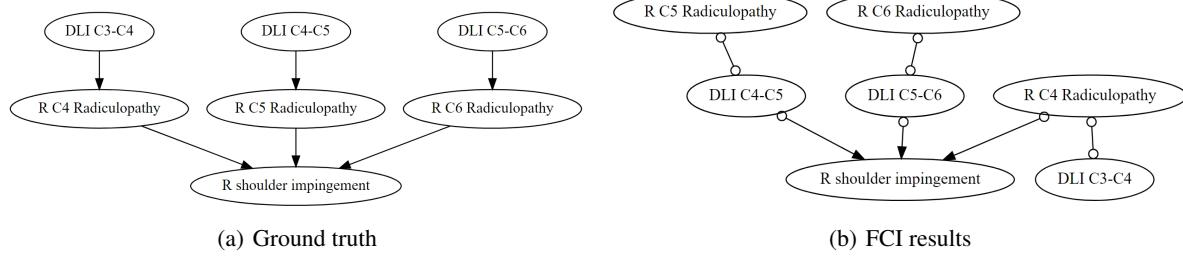


Figure 22. Ground truth and faithful (via FCI algorithm) causal graphs in Neuropathic.

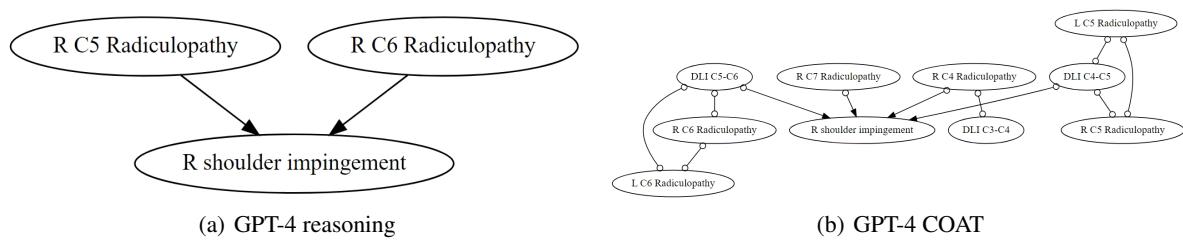


Figure 23. Causal graphs with GPT-4 in Neuropathic.

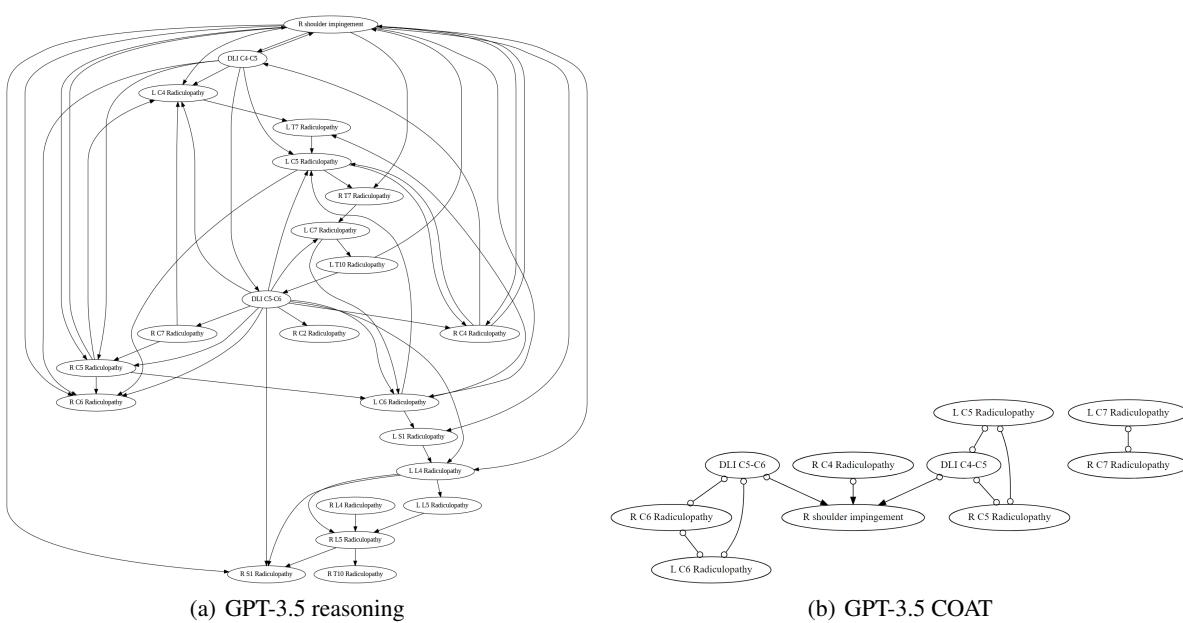
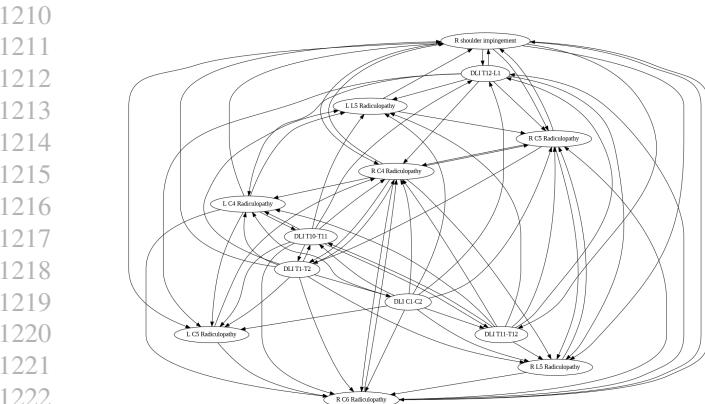
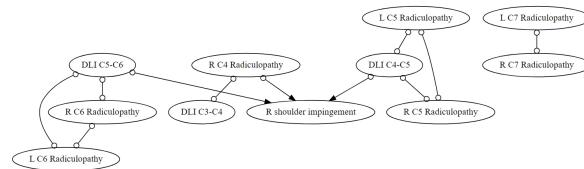


Figure 24. Causal graphs with GPT-3.5 in Neuropathic. (part 1)

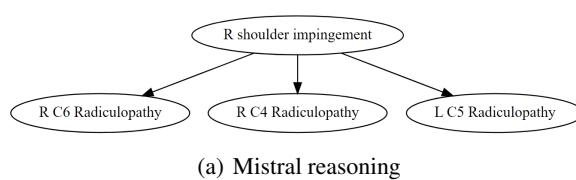


(a) LLaMA2 reasoning

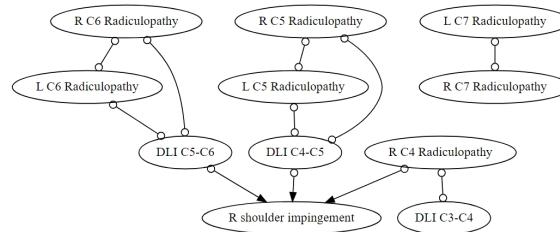


(b) LLaMA2 COAT

Figure 25. Causal graphs with GPT-3.5 in Neuropathic. (part 2)



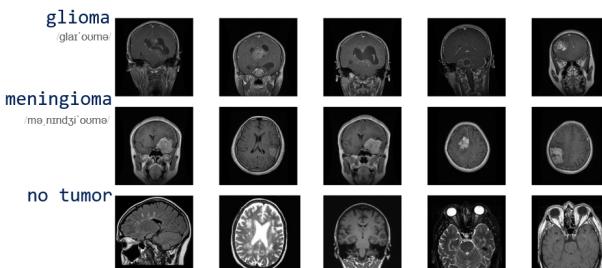
(a) Mistral reasoning



(b) Mistral COAT

Figure 26. Causal graphs with GPT-3.5 in Neuropathic. (part 3)

Brain Tumor

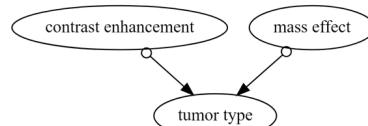


Contrast Enhancement

- 1: Strong, uniform enhancement typical of a **meningioma**.
- -1: Heterogeneous enhancement with or without areas of necrosis typical of a **glioma**.
- 0: No abnormal enhancement or no mass present.

Mass Effect

- 1: Mass causing significant displacement of normal brain structures (can occur with both meningioma and glioma but is more pronounced in **glioma** due to intra-axial location).
- -1: Mass causing minimal or no displacement of brain structures (more common in **meningioma** due to extra-axial location).
- 0: No mass is present, or the specific criteria for 1 or -1 are not met.



"They (meningiomas) almost always demonstrate **uniform strong enhancement** on post-contrast imaging ..."

Lyndon, Daniel, et al. "Dural masses: meningiomas and their mimics." *Insights into imaging* 10.1 (2019): 11.

"High-grade (gliomas) tumors were **heterogeneous contrast enhancement** with other medical imaging aspects of MRI such as bleeding, necrosis, and edema"

Haydar, Nisreen, et al. "Role of Magnetic Resonance Imaging (MRI) in grading gliomas comparable with pathology: A cross-sectional study from Syria." *Annals of Medicine and Surgery* 82 (2022): 104679.

"**Glioblastoma multiforme (GBM)** is the most common primary brain tumor and is associated with a poor prognosis. One of its defining characteristics is the significant deformation of surrounding tissue, the so-called "mass effect"."

Tunc, Birkan, et al. "Modeling of glioma growth with mass effect by longitudinal magnetic resonance imaging." *IEEE Transactions on Biomedical Engineering* 68.12 (2021): 3713-3724.

"Meningiomas typically appear as lobular, **extra axial masses** with well-circumscribed margins. They typically have a broad-based dural attachment and, if sufficiently large, inward displacement of the cortical grey matter."

Watts, J., et al. "Magnetic resonance imaging of meningiomas: a pictorial review." *Insights into imaging* 5 (2014): 113-122.

Figure 27. Visual Illustration for the Brain Tumor Case Study

1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295
 1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319

News & Stock: Results & Evaluation

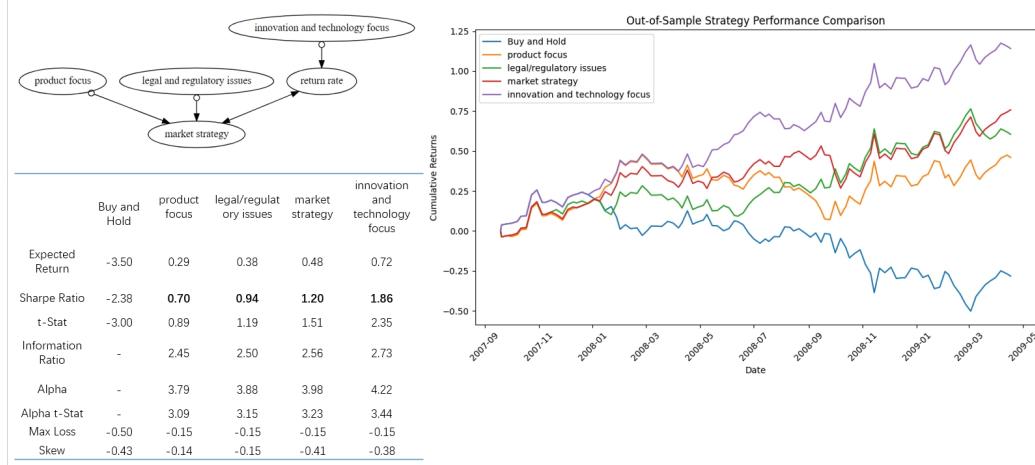


Figure 28. Visual Illustration for the News & Stock Case Study