
Evaluating Interventional Reasoning Capabilities of Large Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Numerous decision-making tasks require estimating causal effects under inter-
2 ventions on different parts of a system. As practitioners consider using large
3 language models (LLMs) to automate decisions, studying their causal reasoning
4 capabilities becomes crucial. A recent line of work evaluates LLMs ability to
5 retrieve commonsense causal facts, but these evaluations do not sufficiently assess
6 how LLMs reason about interventions. Motivated by the role that interventions
7 play in causal inference, in this paper, we conduct empirical analyses to evaluate
8 whether LLMs can accurately update their knowledge of a data-generating process
9 in response to an intervention. We create benchmarks that span diverse causal
10 graphs (e.g., confounding, mediation) and variable types, and enable a study of
11 intervention-based reasoning. These benchmarks allow us to isolate the ability of
12 LLMs to accurately predict changes resulting from their ability to memorize facts
13 or find other shortcuts. We evaluate six LLMs on the benchmarks, finding that GPT
14 models show promising accuracy at predicting the intervention effects.

15 1 Introduction

16 Large language models (LLMs) have achieved impressive performance on a variety of human-relevant
17 tasks, from summarizing web-based information and answering complex questions, to carrying out
18 tasks as web-based agents [Chen et al., 2021, Brown et al., 2020, Li et al., 2022, Katz et al., 2024,
19 Drouin et al., 2024]. As LLMs become increasingly used to make decisions, which fundamentally
20 require understanding the causal impact various actions can have, there has been a recent push to
21 evaluate whether LLMs demonstrate causal reasoning ability [Cai et al., 2023, Jin et al., 2023, 2024,
22 Kiciman et al., 2024, Liu et al., 2024]. The challenge for this emerging field is to operationalize
23 notions of causal reasoning that can be verbalized as text-based questions for LLMs. Kiciman et al.
24 [2024] address this question by evaluating the ability of LLMs to retrieve known cause-and-effect
25 relations. In contrast, Jin et al. [2023, 2024], Liu et al. [2024], Cai et al. [2023] focus on defining
26 various queries that could only be solved with knowledge of causality and graphical models, testing
27 abstract causal reasoning that could generalize to unseen contexts. In this paper, we contribute to the
28 growing body of work that evaluates the causal reasoning capabilities of LLMs by focusing on an
29 aspect of causality that is both intuitive for humans [Waldmann and Hagmayer, 2005] and crucial for
30 decision-making [Binz and Schulz, 2023]: the ability to adapt our model of the world in response to
31 *interventions*.

32 The notion of intervening on a variable is at the core of causality. In this paper, we focus on
33 perfect interventions, where a variable in a system is manipulated and set to a particular new
34 value. Interventions modify the causal graphical model [Pearl, 2009] of a system by deleting
35 all incoming edges to the intervened variable. Experiments in the field of cognitive psychology
36 [Waldmann and Hagmayer, 2005] suggest that humans instinctively recognize that an action (i.e., an

37 intervention) changes existing causal relationships, correctly making different inferences before and
38 after interventions are performed on various variables.

39 Besides being natural for humans, reasoning about interventions is crucial for causal inference.
40 Consider the example of an LLM agent that acts as a “science assistant.” It knows about some causal
41 relationships and must infer new ones given observational evidence to suggest candidate experiments.
42 For ease, consider three variables, A , B and C , where A is known to cause B and C , but where the
43 effect that B has on C is unknown. Initially, the agent cannot conclude anything about the effect
44 of B on C due to the confounding effect of A , but if it encounters new information that B was
45 intervened on, and B was observed to be correlated with C , the agent should correctly conclude that
46 the intervention severs the confounding effect of A , allowing it to infer that B must actually have a
47 causal effect on C . The ability to understand how interventions affect known causal relationships is
48 central to drawing new causal conclusions when presented with new evidence.

49 The key contribution of this paper is to evaluate intervention reasoning in LLMs by introducing
50 *intervention effects*, binary classification tasks that evaluate which causal relations in a graph are
51 modified by an observed intervention. We provide a framework for verbalizing a wide range of
52 intervention effects to LLMs, varying the choice of causal graphs and names of variables. With
53 this framework, we develop a benchmark of intervention effect tasks that help us disentangle the
54 intervention reasoning capabilities of LLMs from other factors that contribute to performance such as
55 memorization of facts and the ability to extract graphs from text.

56 **Related work.** This paper relates most closely to recent papers that develop benchmarks to evaluate
57 LLMs on various causal reasoning tasks. Kiciman et al. [2024] introduced multiple causal reasoning
58 benchmarks for LLMs, including evaluating the ability of LLMs to recover the bivariate causal
59 DAGs introduced in the Tübingen pairs dataset [Mooij et al., 2016]. Kiciman et al. [2024] found that
60 GPT models recovered known causal relationships with up to 96% accuracy when experimenting
61 with various prompting strategies such as including system prompts. However, evaluating LLMs on
62 their ability to retrieve causal knowledge about known variables constitutes *commonsense causal*
63 *reasoning*. In contrast, this paper contributes to work that evaluates abstract causal reasoning [Binz
64 and Schulz, 2023, Jin et al., 2024, 2023, Cai et al., 2023, Liu et al., 2024], assessing the ability of
65 LLMs to use axioms of causality to solve tasks involving general or even new variables.

66 In the vein of causal reasoning, Jin et al. [2023] studied whether LLMs could correctly infer some
67 causal relationships based on conditional independence statements, comparing LLM predictions
68 to those made by an oracle causal discovery algorithm for observational data, where all causal
69 relations cannot be resolved. In contrast to observational causal discovery, this paper focuses on
70 reasoning about interventions. Also focusing on interventions, Jin et al. [2024] introduced CLadder,
71 a comprehensive benchmark that includes the estimation of causal effects from quantitative data.
72 Causal effect estimation is a complex task that requires solving multiple sub-tasks such as: (i)
73 parsing the prompt to extract a causal DAG, (ii) inferring a function that estimates the effect given
74 the DAG, and (iii) applying that function to the given quantitative data. Concurrently, Cai et al.
75 [2023] introduced a task that asks LLMs to output only causal relationships given a tabular dataset
76 that includes variable names. They focus on disentangling the impacts that prior knowledge (e.g.,
77 variables names) and quantitative data have on LLM performance. The empirical study we conduct to
78 assess whether LLMs are sensitive to the presence of plausibly memorized causal relations is similar
79 to experiments conducted by Cai et al. [2023]. In contrast to these benchmarks, intervention effects
80 target a narrower question than the general estimation of causal effects, since intervention effects
81 involve binary classification only (i.e., the absence/presence of causal relations in DAGs). We argue
82 that the evaluations we design better isolate causal reasoning from sub-tasks like drawing statistical
83 inferences from quantitative data provided in-context, which both CLadder and the work of Cai et al.
84 [2023] require.

85 In focusing on intervention effects, we build on the work of Binz and Schulz [2023], who were
86 motivated by prior work in psychology [Waldmann and Hagmayer, 2005] that shows humans weight
87 collected observational evidence and experimental evidence differently when drawing causal con-
88 clusions. Binz and Schulz [2023] adapted this psychology study for LLMs, creating prompts that
89 describe observational and post-interventional findings to LLMs to see if they update their beliefs
90 about a system after interventions. They found that GPT-3, unlike human subjects, fared poorly
91 at understanding the implications of interventions. Motivated by their focus on intervention-based
92 reasoning, we significantly expand on the evaluations designed by Binz and Schulz [2023], systemati-

93 cally generating intervention effects with varying degrees of difficulty to further explore the effects of
 94 plausible memorization and shortcuts like relation retrieval.

95 2 Understanding changes in causal relationships via Intervention Effects

96 We begin by summarizing causal directed acyclic graphical models (DAGs) and perfect interventions,
 97 the two concepts central to this work. After reviewing these concepts, we introduce intervention
 98 effects, the binary prediction tasks that serve as key contributions of this paper.

99 2.1 Background

100 A causal DAG $G = (V = \{V_1, \dots, V_n\}, E)$ is defined by a vertex set V that consists of random
 101 variables $\{V_1, \dots, V_n\}$ and directed edges in the set E that represent causal relationships between
 102 the variables. A DAG defines a joint distribution over the random variables that factorizes according
 103 to the graph,

$$P(V_1, \dots, V_n) = \prod_{i=1}^n P(V_i | \text{Pa}_i^G), \quad (1)$$

104 where Pa_i^G denotes the parents of V_i in the DAG G . In general, we say that a causal relation exists
 105 between two variables u and v when there is a path from u to v in G .

106 Figure 2 in the appendix illustrates the three causal DAGs that we focus on in this paper: bivariate,
 107 confounding, and mediation. In the confounding graph, the variable A causes both B and C , thus
 108 confounding our ability to infer the causal effect that B has on C from observational data alone. In
 109 the mediation graph, the variable A only has an indirect effect on C via B .

110 Causal DAGs also entail distributions after interventions to these variables, distinguishing them from
 111 standard DAG models. In this paper, we focus on a class of interventions called *perfect interventions*,
 112 represented using the operation $\text{do}(V_i = v)$, which means that the variable V_i is set to the value v .
 113 We define the distribution over variables post-intervention by modifying G to delete all incoming
 114 edges to V_i , severing the links between V_i and its parents.

115 2.2 Intervention effects

116 To formalize intervention effects, the key binary classification task that we consider in this work, we
 117 begin by defining causal relations in DAGs precisely.

118 **Definition 1 (Causal relation)** *Given a causal DAG $G = (V, E)$, we say that there exists a causal*
 119 *relation between V_u and V_v in DAG G if there exists a directed path $e_{ua} \rightarrow \dots \rightarrow e_{bv}$ in G from V_u*
 120 *to V_v (where each edge e_{ij} along the path is in G). We define an indicator variable $C_{uv}(G) \in \{0, 1\}$,*
 121 *such that $C_{uv}(G) = 1$ if and only if there is a causal relation between V_v and V_u .*

122 In words, a causal relation captures whether or not a variable exerts an indirect or direct causal
 123 influence on another variable in a particular DAG G that contains these variables.

124 **Definition 2 (Intervention effect)** *Given a causal DAG $G = (V, E)$, a variable V_i on which we*
 125 *perform a perfect intervention captured by $\text{do}(V_i = *)$ (where we use “*” to indicate that we do not*
 126 *care about the value that V_i is set to), and a query causal relation C_{uv} , an **intervention effect (IE)**
 127 *defines a binary classification task as follows,**

$$\text{IE}_i^G(C_{uv}) = C_{uv}(G) - C_{uv}(G^i), \quad (2)$$

128 where the DAG G^i is the modification of DAG G under an intervention to the variable V_i .

129 An $\text{IE}_i^G(C_{uv})$ is defined with respect to a DAG G and intervention target V_i , and assigns each causal
 130 relation C_{uv} to a binary label of 1 or 0. When a causal relation C_{uv} in G is modified by a perfect
 131 intervention on V_i , i.e., the relation C_{uv} is different in the modified DAG G^i , $\text{IE}_i^G(C_{uv})$ is 1 (and 0
 132 otherwise). Intuitively, since interventions sever edges between a variable and its parents, $\text{IE}_i^G(C_{uv})$
 133 is 1 when the intervention target is a variable along the path from u to v (including v itself).

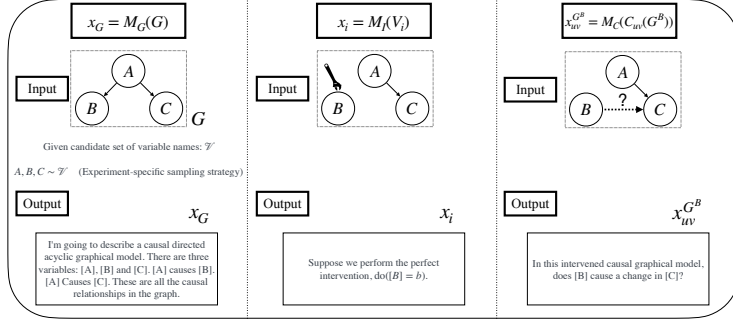


Figure 1: An illustration of the mapping functions that verbalize the information for a single intervention effect estimation task into a prompt.

134 Given a causal DAG G , each possible intervention that we can perform defines a different classification
 135 task over causal relations, defined by a particular $\text{IE}_i^G(\cdot)$, which serves as a labeling function for
 136 that task. We focus on the three causal DAGs in Figure 2 since they sufficiently capture many
 137 real scenarios without introducing unnecessary complexity that could render the empirical findings
 138 ambiguous. The key goal of this paper, addressed in the next section, is to evaluate LLMs on these IE
 139 classification tasks in a zero-shot way (i.e., without training examples) by verbalizing the tasks as
 140 prompts.

141 3 Evaluating LLMs on Intervention Effects

142 To evaluate an LLM on a binary classification problem defined by $\text{IE}_i^G(C_{uv})$, we need to verbalize
 143 the causal DAG G , the concept of an intervention as well as the intervention target V_i , and the causal
 144 relation C_{uv} of interest. Figure 1 illustrates how we verbalize these three steps to generate a complete
 145 prompt.

146 **Step 1: Generating variable names.** To verbalize an input causal DAG $G = (V, E)$, we select
 147 names for the variables V and describe each edge $E_{ij} \in E$ using the format “[i] causes [j].” We
 148 further specify that these are all known causal relationships in the graph to avoid ambiguities, e.g.,
 149 around the presence of unobserved confounding variables. In the empirical studies, we assess multiple
 150 ways of choosing variable names, designing studies to tease apart how well LLMs generalize to novel
 151 contexts instead of relying on variable names and facts that were potentially encountered during
 152 training.

153 **Step 2: Describing an intervention.** In the second step, we specify that a perfect intervention is
 154 performed, during the notation $\text{do}(V_i = v)$ since the do-operator could have been encountered during
 155 training. To further study the robustness of LLMs to memorizing facts from training, the empirical
 156 studies include experiments that use random strings to describe do-operator.

157 **Step 3: Verbalizing the binary classification problem.** In the final step, we verbalize the target
 158 $\text{IE}_i^G(C_{uv})$, using the phrase “does [u] cause a change in [v]?” to query the presence or absence of a
 159 causal relation between u and v .

160 **Evaluation metric.** After verbalizing an intervention effect as a prompt to an LLM, we receive a
 161 yes/no response that we refer to as $\hat{\text{IE}}_i^G(C_{uv})$. Although prediction accuracy is a natural evaluation
 162 metric in this setting, it can be misleading in this setting. To see this, we begin by noting that the LLM
 163 forms some belief about the presence or absence of a causal relation in a DAG G , denoted $\hat{C}_{uv}(G)$.
 164 Consider the scenario where,

$$C_{uv}(G) = 1, C_{uv}(G^i) = 1, \\ \hat{C}_{uv}(G) = 0, \hat{C}_{uv}(G^i) = 0.$$

165 In this example, a target causal relation C_{uv} is true in both the base causal DAG G and its modified,
 166 post-intervention counterpart G^i , but the LLM incorrectly predicts that these causal relations are
 167 false under both graphical scenarios. The accuracy metric misleads us when the LLM correctly
 168 predicts that the causal relation does not vary, but does not parse the causal relation correctly in
 169 either graphical scenario. Thus, to ensure that the accuracy is 0 in such cases, we slightly modify the

Table 1: **IE prediction accuracy on the Random Char benchmark.** GPT-4 variants are the best performing models, while LLaMA-2 appears to struggle with interventional reasoning. Performances that are significantly ($\alpha = 0.05$) worse than 80% accuracy are shown in red and top performances are indicated by bolded figures.

Graph Type Intervened Variable	Bivariate		Confounding			Mediation		
	A	B	A	B	C	A	B	C
GPT-3.5	0.83 ± 0.08	0.87 ± 0.06	0.80 ± 0.09	0.69 ± 0.12	0.36 ± 0.09	0.58 ± 0.11	0.36 ± 0.12	0.67 ± 0.12
GPT-4	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	0.78 ± 0.09	0.82 ± 0.08	0.96 ± 0.03
GPT-4-turbo	1.0 ± 0.0	0.97 ± 0.03	0.96 ± 0.03	0.93 ± 0.05	1.0 ± 0.0	0.98 ± 0.02	1.0 ± 0.0	1.0 ± 0.0
GPT-4o	1.0 ± 0.0	1.0 ± 0.0	0.98 ± 0.02	0.91 ± 0.04	1.0 ± 0.0	0.93 ± 0.05	1.0 ± 0.0	0.87 ± 0.06
LLaMA-2	0.50 ± 0.12	0.40 ± 0.12	0.56 ± 0.12	0.53 ± 0.11	0.16 ± 0.06	0.69 ± 0.09	0.56 ± 0.12	0.64 ± 0.12
LLaMA-3	0.80 ± 0.09	0.83 ± 0.06	0.51 ± 0.1	0.47 ± 0.08	0.96 ± 0.03	0.38 ± 0.1	0.53 ± 0.08	0.4 ± 0.1

170 accuracy metric to be,

$$1[\text{IE}_i^G(C_{uv}) = \hat{\text{IE}}_i^G(C_{uv})] \cdot 1[C_{uv}(G) = \hat{C}_{uv}(G)] \quad (3)$$

171 Now, an LLM’s prediction is only considered correct when it understands the DAG G correctly and
 172 accurately solves an IE problem. To infer the belief $\hat{C}_{uv}(G)$ that an LLM has about a causal relation,
 173 we make a minor modification to the prompt in Figure 1, omitting the second step that verbalizes an
 174 intervention, and asking the LLM about a causal relation but in the “observed causal graphical model”
 175 instead of the intervened one. These additional prompts allow us to evaluate LLMs on IEs.

176 In the next section, we investigate several research questions about LLMs and intervention reasoning
 177 ability using the proposed framework to define a suite of IE problems.

178 4 Empirical Analysis

179 To study the ability of LLMs to understand interventions, we use the introduced framework to develop
 180 three benchmarks that differ in how variable names are selected to verbalize IEs:

- 181 1. **Random Char**: Variable names are randomly chosen English characters.
- 182 2. **Tübingen** : Variables names are chosen from entities that appear in the Tübingen pairs (TP)
 183 dataset of causal relations [Mooij et al., 2016] such that some causal relations exist in the
 184 TP data.
- 185 3. **Random Tübingen** : Variable names are chosen from entities that appear in the TP dataset
 186 so that no causal relations exist in the TP data.

187 See Appendix A.2 for an illustration of the full details about the datasets. For each benchmark and
 188 IE, we sample variable names fifteen times to report significant differences.

189 We investigate four research questions (RQs) on the proposed benchmarks, studying four LLMs: GPT-
 190 3.5, GPT-4, GPT-4-turbo [OpenAI, 2023], and LLaMA-2 [Touvron et al., 2023]. Unless otherwise
 191 specified, for each IE, we aggregate results over the enumerated causal relations. We present the
 192 findings for the first two questions in the main paper and discuss the remaining in Appendix B

193 **RQ1: How accurate are at LLMs at IEs generated with random characters as variables?** To
 194 study this first question, we evaluate the performance of the four LLMs on the Random benchmark.
 195 Table 1 summarize these results. We see that GPT-based models perform notably better than the
 196 LLaMA models, with GPT-4-turbo demonstrating near-perfect accuracy across all effects. LLaMA-2
 197 and LLaMA-3’s performance suggests that they do not reliably model interventions. In what follows
 198 in the main paper, we focus on results for GPT models, deferring LLaMA results to Appendix B.2.

199 **RQ2: To what extent is LLM performance affected by possibly memorized causal rela-**
 200 **tions?** While we might be tempted to conclude that GPT-4 reliably predicts changes to models
 201 after interventions are performed, we consider spurious factors that can affect model performance.
 202 In particular, Kiciman et al. [2024] found that GPT models reliably retrieved information about TP
 203 causal relations, suggesting that these relations could have been included in the training data for
 204 LLMs. (We reproduce their findings in Appendix Table 4.) This leads to a worrying possibility:

Table 2: For causal relations that appear in TP, and for interventions that modify these relations, worse performance on the Tübingen benchmark compared to performance on Random Char or Random Tübingen provides evidence that LLMs might be relying on memorized facts to achieve high accuracy on IEs. Bolded figures indicate performances that are significantly ($\alpha = 0.05$) worse than the corresponding performance on Tübingen .

Model	Graph	Intervened Variable	Random Char	Tübingen	Random Tübingen
GPT-3.5	Bivariate	B	0.8 ± 0.1	0.8 ± 0.1	0.53 ± 0.13
	Confounding	C	0.47 ± 0.13	0.4 ± 0.13	0.47 ± 0.13
		Mediation	B	0.4 ± 0.13	0.4 ± 0.13
		C	0.67 ± 0.12	0.67 ± 0.12	0.6 ± 0.13
GPT-4	Bivariate	B	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
	Confounding	C	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
		Mediation	B	0.8 ± 0.1	0.73 ± 0.11
		C	0.93 ± 0.06	1.0 ± 0.0	0.93 ± 0.06
GPT-4-turbo	Bivariate	B	0.93 ± 0.06	0.8 ± 0.1	0.73 ± 0.11
	Confounding	C	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
		Mediation	B	1.0 ± 0.0	1.0 ± 0.0
		C	1.0 ± 0.0	1.0 ± 0.0	0.93 ± 0.06
GPT-4o	Bivariate	B	1.0 ± 0.0	0.93 ± 0.06	0.73 ± 0.11
	Confounding	C	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
		Mediation	B	1.0 ± 0.0	1.0 ± 0.0
		C	0.93 ± 0.06	0.87 ± 0.09	0.8 ± 0.1

205 after interventions, could LLMs fail to update their beliefs about causal relations that they have
206 potentially memorized? To study this, we consider only the subset of causal relations that appear in
207 TP and the interventions that *sever these relations*. If an LLM achieves good performance in general
208 by having memorized some known causal relations, then it would achieve worse performance on
209 Tübingen for the selected IEs, which sever known causal relations and go against the LLM’s learned
210 biases, compared its performance on Random or Random Tübingen , where post-interventional causal
211 relations do not directly contradict known facts. In Appendix B.2, we specify which causal relations
212 appear in TP for each causal DAG we study, and which interventions modify these relations. Table 2
213 summarizes the results. Bolded figures indicate that the performance on Tübingen drops significantly
214 (with $\alpha = 0.05$) compared to performance on either of the other benchmarks. Interestingly, GPT do
215 not show evidence of relying purely on memorized causal relations, since performance is overall
216 comparable across the benchmarks. Further, the models seem to struggle with the Random Tübingen
217 benchmark more than they do with the other benchmarks, leading to a question about whether they
218 are in fact sensitive to variations on the TP dataset.

219 5 Discussion and Limitations

220 The goal of this paper is to introduce a causal reasoning benchmark that stress-tests the ability of
221 LLMs to accurately predict how knowledge should be updated after interventions are performed,
222 without conflating other aspects of reasoning such as statistical inference on quantitative data. Our
223 findings are optimistic, but we believe that nevertheless, they underscore the continued need for
224 benchmarks and studies that evaluate varied aspects of abstract causal reasoning in LLMs, especially
225 if practitioners wish to use LLMs to generate candidate decisions.

226 While the intervention effect prediction task we define in this paper has the benefit of being easy to
227 evaluate, since it requires binary responses, the findings that this task can suggest are also limited. For
228 example, IE prediction cannot help us assess how accurately LLMs perform *causal identification*, the
229 process of deciding which causal inferences can be made given a causal DAG. Moreover, we focus
230 on evaluation in this paper and do not propose methods for improving causal reasoning in LLMs via
231 few-shot learning or fine-tuning. Both of these limitations point to future research directions that we
232 think are worth exploring.

233 **References**

- 234 [1] Marcel Binz and Eric Schulz. Using cognitive psychology to understand GPT-3. *Proceedings of the*
235 *National Academy of Sciences*, 120(6):e2218523120, 2023. doi: 10.1073/pnas.2218523120. URL
236 <https://www.pnas.org/doi/abs/10.1073/pnas.2218523120>.
- 237 [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
238 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.
239 *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 240 [3] Hengrui Cai, Shengjie Liu, and Rui Song. Is knowledge all large language models needed for causal
241 reasoning? *arXiv:2401.00139*, 2023.
- 242 [4] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan,
243 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models
244 trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 245 [5] Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H Laradji, Manuel Del Verme, Tom Marty,
246 David Vazquez, Nicolas Chapados, and Alexandre Lacoste. Workarena: How capable are web agents at
247 solving common knowledge work tasks? In *Forty-first International Conference on Machine Learning*,
248 2024.
- 249 [6] Zhijing Jin, Jiarui Liu, Zhiheng Lyu, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona Diab, and
250 Bernhard Schölkopf. Can large language models infer causation from correlation? *arXiv:2306.05836*,
251 2023.
- 252 [7] Zhijing Jin, Yuen Chen, Felix Leeb, Luigi Gresele, Ojasv Kamal, Zhiheng Lyu, Kevin Blin, Fernando
253 Gonzalez Adauto, Max Kleiman-Weiner, Mrinmaya Sachan, et al. Cladder: A benchmark to assess causal
254 reasoning capabilities of language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 255 [8] Daniel Martin Katz, Michael James Bommarito, Shang Gao, and Pablo Arredondo. Gpt-4 passes the bar
256 exam. *Philosophical Transactions of the Royal Society A*, 382(2270):20230254, 2024.
- 257 [9] Emre Kiciman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language models:
258 Opening a new frontier for causality. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
259 URL <https://openreview.net/forum?id=mqoxLkX210>. Featured Certification.
- 260 [10] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles,
261 James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode.
262 *Science*, 378(6624):1092–1097, 2022.
- 263 [11] Xiao Liu, Zirui Wu, Xueqing Wu, Pan Lu, Kai-Wei Chang, and Yansong Feng. Are LLMs capable of
264 data-based statistical and causal reasoning? benchmarking advanced quantitative reasoning with data.
265 *arXiv:2402.17644*, 2024.
- 266 [12] Joris M Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. Distinguishing
267 cause from effect using observational data: methods and benchmarks. *Journal of Machine Learning*
268 *Research*, 17(32):1–102, 2016.
- 269 [13] R OpenAI. Gpt-4 technical report. arxiv 2303.08774. *View in Article*, 2(5), 2023.
- 270 [14] Arkil Patel, Satwik Bhattamishra, Siva Reddy, and Dzmitry Bahdanau. Magnifico: Evaluating the in-
271 context learning ability of large language models to generalize to novel interpretations. *arXiv preprint*
272 *arXiv:2310.11634*, 2023.
- 273 [15] Judea Pearl. *Causality*. Cambridge university press, 2009.
- 274 [16] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
275 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. Llama 2: Open foundation and
276 fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 277 [17] Michael R Waldmann and York Hagmayer. Seeing versus doing: two modes of accessing causal knowledge.
278 *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 31(2):216, 2005.

279 **A Experiment Setup Details**

280 **A.1 Causal DAGs**

281 We consider 3 basic causal DAGs: Bivariate, Confounding and Mediation as shown in Figure 2

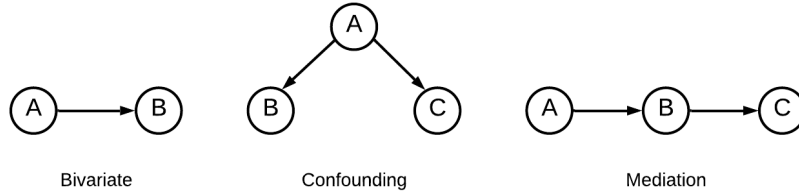


Figure 2: **Causal DAGs.** In the empirical studies, we define intervention effects based on three causal DAGs: bivariate, confounding and mediation graphs.

282 **A.2 Dataset Generation**

283 **Random Char (R).** In the causal graphs used in this benchmark, the nodes are labeled with arbitrary
284 English letters (E.g. 'z', 'b') that are independently generated for each graph.

285 **Tübingen (T).** For the Tübingen benchmark, we select the graph nodes from the Tübingen pairs
286 (TP) dataset in the following manner:

- 287 • **Bivariate.** We sample cause-effect pairs from the TP dataset and each cause-effect pair is
288 used to generate a graph.
289 Example - Altitude (A) - Temperature (B) is a cause-effect pair in the TP dataset, and the
290 input graph states: *Altitude causes temperature.*
- 291 • **Confounding.** For variable *B* and *C*, we first sample cause-effect pairs and then for sampled
292 pair we randomly select a variable (for *A*) from the TP dataset that it is different from the
293 variables in the corresponding cause-effect pair.
294 Example - Altitude (*B*) - Temperature (*C*) is the sampled cause-effect pair and Age (*A*) is
295 the randomly selected variable from the TP dataset. The relationships in the input graph are
296 as follows: *Age causes altitude ; Age causes temperature.*
- 297 • **Mediation.** For variable *A* and *C*, we first randomly sample cause effect pairs and then for
298 sampled pair we randomly select a variable (for *B*) from the TP dataset that it is different
299 from the variables in the corresponding cause-effect pair.
300 Example - Altitude (*A*) - Temperature (*C*) is the sampled cause-effect pair and Age (*B*) is
301 the randomly selected variable from the TP dataset. The relationships in the input graph are
302 as follows: *Altitude causes age ; Age causes temperature.*

303 The rationale behind defining the Tübingen dataset in this manner is elaborated in Appendix B.2.

304 **Random Tübingen (RT).** Similar to the Tübingen case, we select the graph nodes from the set
305 of variables present in the TP dataset. However, none of the causal relations defined in the graphs
306 are present in the TP dataset. Instead the causal relations defined are between randomly selected
307 variables, as follows:

- 308 • **Bivariate.** We randomly sample two unrelated variables from the TP dataset and define a
309 cause-effect relationship in the input graph.
310 Example - Age (*A*) - Temperature (*B*) are variables without a cause-effect relation in the TP
311 dataset, and the input graph states: *Age causes Temperature.*
- 312 • **Confounding.** We randomly sample 3 variables from the TP dataset that no two variables
313 selected have a causal relation in the TP dataset and we define the graph.
314 Example - Altitude (*A*) - Horsepower (*B*) - Cement (*C*) are randomly selected variables
315 from the TP dataset. No two variables have a causal relationship in the TP dataset. The

316 relationships in the input graph is as follows: *Altitude causes Horsepower ; Altitude causes*
317 *Cement*.

318 • **Mediation.** We select the 3 variables exactly as we did for the confounding case and define
319 the graph.

320 Example - Altitude (A) - Horsepower (B) - Cement (C) are randomly selected variables
321 from the TP dataset. The relationships in the input graph is as follows: *Altitude causes*
322 *Horsepower ; Horspower causes Cement*.

323 **A.3 Details regarding LLMs**

324 **A.3.1 Querying LLMs.**

325 To query the GPT models, we used the OpenAI API ¹ through the Langchain interface ². Meanwhile,
326 VLLM library ³ was used for fast inference from LLaMA models.

327 Regarding compute resources, since the GPT models are hosted remotely by OpenAI, we were able
328 to make API calls locally with minimal CPU usage. Conversely, for LLaMA models, we first had
329 to load the model onto a cluster equipped with a GPU (A100/80 GB RAM) before submitting our
330 queries.

331 **A.3.2 LLM Output.**

332 To ensure that LLMs respond with a yes or no to queries, we formalized a response format requiring
333 LLMs to encapsulate their yes or no answers within an `<answer></answer>` tag. If the LLM does
334 not adhere to this format, we initiate a retry, instructing it to comply with the required format. After
335 the first retry, we relax the response format requirements, expecting only a yes or no answer. If after
336 10 retries the LLM fails to meet this criterion, we mark the attempt as a failure and attribute zero
337 accuracy on the intervention effect prediction task.

338 **A.3.3 Model Ids.**

339 The specific model-ids of the LLMs are:

- 340 • GPT-3.5: `gpt-3-turbo-16k`
- 341 • GPT-4: `gpt-4`
- 342 • GPT-4-turbo: `gpt-4-turbo`
- 343 • GPT-4o: `gpt-4o`
- 344 • LLaMA-2: `llama-2-7b`
- 345 • LLaMA-3: `llama-3.1-8b`

346 **A.4 Substitution Task**

347 For this task, we first describe the concept of interventions to LLM and refer to them as some arbitrary
348 string of characters (E.g. 'xyz'). We then follow the same prompt template (Figure 1) to test the
349 intervention effect performance, except we never use the word intervention again and replace it with
350 the chosen substitution word. The prompt verbalization for this task is illustrated in Figure 3.

351 **A.5 Ground Truth Intervention Effects**

352 We provide the ground truth value for every intervention effect task in our analysis in Table 3, i.e., we
353 provide the true change in the causal relations after interventions.

354 This should help the reader understand clearly what queries we considered for the different research
355 questions in our analysis.

¹<https://openai.com/blog/openai-api>

²https://python.langchain.com/docs/expression_language/interface

³<https://docs.vllm.ai/en/latest/>

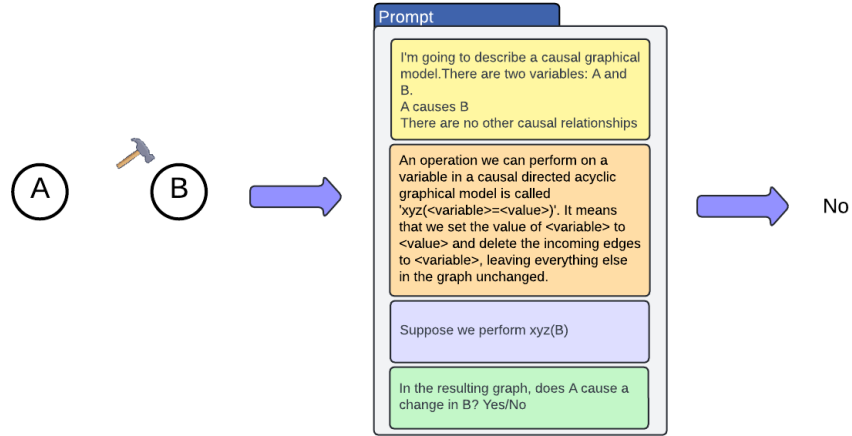


Figure 3: **Prompt design for intervention reasoning in the substitution task.** Instead of using the word intervention directly in the prompt, we illustrate the concept of interventions and define it with a random word, for instance, operation 'xyz' in the prompt template above.

Table 3: Ground truth intervention effects for all the scenarios (causal graph, intervention variable, causal relation) considered in our benchmark.

Graph	Intervention	Questions	Ground Truth IE
Bivariate	A	$A \rightarrow B$	0
		$B \rightarrow A$	0
	B	$A \rightarrow B$	1
		$B \rightarrow A$	0
Confounding	A	$A \rightarrow B$	0
		$A \rightarrow C$	0
		$B \rightarrow C$	0
	B	$A \rightarrow B$	1
		$A \rightarrow C$	0
		$B \rightarrow C$	0
Mediation	A	$A \rightarrow B$	0
		$A \rightarrow C$	0
		$B \rightarrow C$	0
	B	$A \rightarrow B$	1
		$A \rightarrow C$	1
		$B \rightarrow C$	0
C	$A \rightarrow B$	0	
	$A \rightarrow C$	1	
	$B \rightarrow C$	1	

356 **B Additional Results**

357 We now go over the empirical studies in detail and provide the results for the remaining research
 358 questions.

359 **B.1 Reproducing results on the Tübingen dataset**

360 We reproduced the results of Kiciman et al. [2024] on the Tübingen dataset as shown in Table 4.
 361 Since their task did not involve the LLM to reason about the effect of interventions, we term this task
 362 of inferring relations from the input graph in the prompt as relation retrieval.

Table 4: Relation Retrieval accuracy of GPT models on the **Tübingen** dataset.

Model	Accuracy
GPT-3 (<i>text-davinci-003</i>)	0.80
GPT-3.5	0.89
GPT-4	0.96

Table 5: **Relation Retrieval accuracy of models on Random benchmark.** The GPT models show good performance but LLaMA-2 performs poorly. Hence, as per our criteria we drop the LLaMA-2 model for analysis in RQ3.

Model	Bivariate	Confounding	Mediation
GPT-3.5	1.0 ± 0.0	1.0 ± 0.0	0.98 ± 0.02
GPT-4	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
GPT-4-turbo	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
GPT-4o	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
LLaMA-2	0.73 ± 0.11	0.59 ± 0.13	0.84 ± 0.09
LLaMA-3	1.0 ± 0.0	1.0 ± 0.0	0.94 ± 0.06

363 **B.2 Research Questions**

364 **RQ1: How accurate are LLMs at predicting the effects of interventions?**

365 To understand how effective LLMs are at IE prediction, we focused on the **Random (R)** benchmark
 366 in the main paper (Table 1), as we wanted to remove any distracting effect of semantically meaningful
 367 entities as graph nodes. We now present the same results on the **Tübingen (T)** and **Random**
 368 **Tübingen (RT)** benchmarks as well in Table 6.

369 We find that the results in both the cases are very similar to the case with the **Random** benchmark;
 370 GPT-4-turbo performs the best and GPT models outperform LLaMA models by a big margin.

371 **RQ2: To what extent is LLM performance affected by memorized causal relations?** In **RQ1**,
 372 we computed the IE performance by considering all the causal relationship queries for each of the
 373 interventions. In the case of **RQ2**, we only consider specific causal relations that could be potentially
 374 memorized by the LLMs and interventions that sever these relations. For:

- 375 1. Bivariate DAG, we consider $A \rightarrow B$ with the intervention on B .
- 376 2. Confounding DAG, we consider $B \rightarrow C$ with the intervention on C .
- 377 3. Mediation DAG, we consider $A \rightarrow C$ with an intervention on B and C separately.

378 The **Tübingen** benchmark (Appendix A.2) is defined in such a way that, all of the causal relations
 379 under consideration above are cause-effect pairs present in the TP dataset that could be potentially
 380 memorized by the LLM.

381 We provide results for role of memorization in LLaMA models in Table 7. Given their overall poor
 382 IE performance in the random benchmark, it is difficult to draw conclusions about the impact of
 383 memorization but they have similar performance across the benchmarks.

Table 6: **IE prediction accuracy on all three benchmarks.** GPT-4 variants are the best performing models, while LLaMA models struggle with interventional reasoning.

Graph Type	Benchmark	Bivariate		Confounding			Mediation		
		A	B	A	B	C	A	B	C
GPT-3.5	R	0.83 ± 0.08	0.87 ± 0.06	0.8 ± 0.09	0.69 ± 0.12	0.36 ± 0.09	0.58 ± 0.11	0.36 ± 0.12	0.67 ± 0.12
	T	0.87 ± 0.06	0.83 ± 0.06	0.82 ± 0.05	0.67 ± 0.11	0.31 ± 0.07	0.64 ± 0.08	0.42 ± 0.1	0.67 ± 0.12
	RT	0.83 ± 0.1	0.7 ± 0.12	0.87 ± 0.09	0.67 ± 0.12	0.47 ± 0.13	0.64 ± 0.12	0.33 ± 0.12	0.63 ± 0.12
GPT-4	R	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	0.78 ± 0.09	0.82 ± 0.08	0.96 ± 0.03
	T	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	0.84 ± 0.07	1.0 ± 0.0
	RT	1.0 ± 0.0	1.0 ± 0.0	0.93 ± 0.06	0.87 ± 0.09	1.0 ± 0.0	0.98 ± 0.04	0.97 ± 0.05	0.97 ± 0.05
GPT-4-turbo	R	1.0 ± 0.0	0.97 ± 0.03	0.96 ± 0.03	0.93 ± 0.05	1.0 ± 0.0	0.98 ± 0.02	1.0 ± 0.0	1.0 ± 0.0
	T	0.93 ± 0.04	0.9 ± 0.05	0.62 ± 0.09	0.78 ± 0.08	0.96 ± 0.03	0.69 ± 0.09	0.96 ± 0.04	0.96 ± 0.04
	RT	0.8 ± 0.1	0.8 ± 0.1	0.73 ± 0.11	0.6 ± 0.13	1.0 ± 0.0	0.8 ± 0.1	0.97 ± 0.05	0.97 ± 0.05
GPT-4o	R	1.0 ± 0.0	1.0 ± 0.0	0.98 ± 0.02	0.91 ± 0.04	1.0 ± 0.0	0.93 ± 0.05	1.0 ± 0.0	0.87 ± 0.06
	T	0.87 ± 0.06	0.97 ± 0.03	0.89 ± 0.04	0.93 ± 0.03	1.0 ± 0.0	0.84 ± 0.05	0.96 ± 0.04	0.87 ± 0.06
	RT	0.87 ± 0.09	0.83 ± 0.1	0.9 ± 0.08	0.8 ± 0.1	1.0 ± 0.0	0.93 ± 0.06	0.93 ± 0.06	0.87 ± 0.09
LLaMA-2	R	0.5 ± 0.12	0.4 ± 0.12	0.56 ± 0.11	0.53 ± 0.11	0.16 ± 0.06	0.69 ± 0.09	0.56 ± 0.12	0.64 ± 0.12
	T	0.4 ± 0.12	0.2 ± 0.08	0.42 ± 0.11	0.4 ± 0.11	0.09 ± 0.04	0.47 ± 0.12	0.31 ± 0.12	0.44 ± 0.1
	RT	0.43 ± 0.13	0.23 ± 0.11	0.47 ± 0.13	0.47 ± 0.13	0.2 ± 0.1	0.62 ± 0.13	0.4 ± 0.13	0.7 ± 0.12
LLaMA-3	R	0.8 ± 0.1	0.83 ± 0.1	0.47 ± 0.13	0.6 ± 0.13	0.87 ± 0.09	0.38 ± 0.13	0.5 ± 0.13	0.47 ± 0.13
	T	0.73 ± 0.11	0.87 ± 0.09	0.53 ± 0.13	0.53 ± 0.13	1.0 ± 0.0	0.38 ± 0.13	0.6 ± 0.13	0.27 ± 0.11
	RT	0.67 ± 0.12	0.87 ± 0.09	0.47 ± 0.13	0.4 ± 0.13	1.0 ± 0.0	0.4 ± 0.13	0.53 ± 0.13	0.43 ± 0.13

Table 7: **IE prediction performance of LLaMA models on specific scenarios for all the benchmarks to understand the role memorization.** LLaMA models demonstrate relatively similar performance across the three benchmarks.

Model	Graph	Intervened Variable	Random	Tübingen	Random Tübingen
LlaMA-2	Bivariate	B	0.4 ± 0.13	0.4 ± 0.13	0.27 ± 0.11
	Confounding	C	0.0 ± 0.0	0.07 ± 0.06	0.33 ± 0.12
	Mediation	B	0.33 ± 0.12	0.33 ± 0.12	0.4 ± 0.13
LlaMA-3	Bivariate	B	0.67 ± 0.12	0.73 ± 0.11	0.8 ± 0.1
	Confounding	C	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
	Mediation	B	0.67 ± 0.12	0.53 ± 0.13	0.67 ± 0.12
		C	0.53 ± 0.13	0.33 ± 0.12	0.4 ± 0.13

384 **RQ3: Could LLMs be learning a shortcut to predict intervention effects?** Consider the confound-
385 ing DAG (Figure 2) and the causal relation $C_{BC}(G)$ which does not change under intervention on
386 the variable A . In these examples, LLMs that accurately parse causal relations from text descriptions
387 of DAGs would also obtain accurate IE estimates. Hence, predicting causal relations from the input
388 graph in text – call this task relation retrieval – offers a *shortcut*: an LLM can attend to tokens in the
389 context to solve relation retrieval and still perform good at IE estimation, thereby confounding the
390 conclusions that can be drawn with this benchmark.

391 However, the intervention effects we defined to study RQ2 offer an insight into how we can disentangle
392 relation retrieval and accurate IE prediction. Notice that the IEs we defined to study RQ2 characterize
393 scenarios where causal relations differ between the base DAG G and the post-intervention DAG.
394 Thus, LLMs that rely only on relation retrieval cannot accurately estimate IEs. Building on this
395 insight, we divide all the IEs $\kappa_G^i(C_{uv})$ into two groups:

- 396 1. **IE = 0**: Graph doesn’t change as a result of intervention; $C_{uv}(G) = C_{uv}(G^i)$.
- 397 2. **IE = 1**: Graph changes as a result of intervention; $C_{uv}(G) \neq C_{uv}(G^i)$

398 Table 3 clearly categorizes the causal relations under intervention into these two groups.

399 Note that drop in performance on group ($IE = 1$) as compared to the group ($IE = 0$) can indicate
400 reliance on shortcuts based on relation retrieval. We consider the average performance of the LLMs
401 on the effects in each group, selecting only those LLMs which achieve an accuracy ≥ 0.95 on relation
402 retrieval (which we report in Appendix Table 5). We focus on the **Random** benchmark to exclude any
403 impacts due to memorized variable names. Table 8 summarizes this study. We find that the general
404 trend does not show strong reliance on shortcuts across LLMs; only GPT-3.5 has a significant drop

Table 8: **IE prediction performance across sub-cases to isolate the effect of relation retrieval on Random benchmark.** LLMs do not significantly rely on shortcuts related to relation retrieval from the input prompt. Since intervening on variable A never changes the causal graph in any scenario, we don’t consider them for this analysis.

Graph Type		Bivariate		Confounding		Mediation	
Intervened Variable		B	B	C	B	C	
GPT-3.5	IE = 0	0.93 ± 0.06	0.6 ± 0.13	0.47 ± 0.13	0.33 ± 0.12	0.67 ± 0.12	
	IE = 1	0.8 ± 0.1	0.67 ± 0.12	0.2 ± 0.1	0.37 ± 0.12	0.67 ± 0.12	
GPT-4	IE = 0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	0.87 ± 0.09	1.0 ± 0.0	
	IE = 1	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	0.8 ± 0.1	0.93 ± 0.06	
GPT-4-turbo	IE = 0	1.0 ± 0.0	0.96 ± 0.05	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	
	IE = 1	0.93 ± 0.06	0.93 ± 0.06	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	
GPT-4o	IE = 0	1.0 ± 0.0	0.87 ± 0.09	1.0 ± 0.0	1.0 ± 0.0	0.87 ± 0.09	
	IE = 1	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	0.87 ± 0.09	
LlaMA-2	IE = 0	0.4 ± 0.13	0.5 ± 0.13	0.0 ± 0.0	0.33 ± 0.12	0.67 ± 0.12	
	IE = 1	0.4 ± 0.13	0.53 ± 0.13	0.13 ± 0.09	0.33 ± 0.12	0.6 ± 0.13	
LlaMA-3	IE = 0	1.0 ± 0.0	0.4 ± 0.13	1.0 ± 0.0	0.6 ± 0.13	0.27 ± 0.11	
	IE = 1	0.67 ± 0.12	0.6 ± 0.13	0.87 ± 0.09	0.5 ± 0.13	0.47 ± 0.13	

Table 9: **IE prediction accuracy on the Random benchmark for the substitution task.** The performance of GPT-3.5 & GPT-4 is worse under substitution compared to the non-substitution case (Table 1), while GPT-4-turbo does not show significant change.

Graph Type	Bivariate		Confounding			Mediation		
Intervened Variable	A	B	A	B	C	A	B	C
GPT-3.5	0.67 ± 0.12	0.83 ± 0.06	0.56 ± 0.11	0.58 ± 0.11	1.0 ± 0.0	0.62 ± 0.10	0.44 ± 0.10	0.4 ± 0.12
GPT-4	0.97 ± 0.03	1.0 ± 0.0	0.89 ± 0.05	0.82 ± 0.07	1.0 ± 0.0	0.96 ± 0.03	0.76 ± 0.10	1.0 ± 0.0
GPT-4-turbo	0.9 ± 0.07	0.97 ± 0.03	0.96 ± 0.03	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
GPT-4o	0.67 ± 0.12	0.83 ± 0.1	0.67 ± 0.12	0.73 ± 0.11	1.0 ± 0.0	0.76 ± 0.11	0.83 ± 0.1	0.73 ± 0.11
LLaMA-2	0.57 ± 0.13	0.37 ± 0.12	0.6 ± 0.13	0.6 ± 0.13	0.73 ± 0.11	0.62 ± 0.13	0.6 ± 0.13	0.43 ± 0.13
LLaMA-3	0.67 ± 0.12	0.83 ± 0.1	0.6 ± 0.13	0.4 ± 0.13	1.0 ± 0.0	0.4 ± 0.13	0.67 ± 0.12	0.33 ± 0.12

405 in its relative performance on group ($IE = 1$) vs group ($IE = 0$) for the confounding DAG case.
 406 Since, the LLaMA-2 model does not satisfy our constraint of high relation retrieval accuracy, we
 407 don’t consider it for further analysis.

408 **RQ4: Are LLMs robust to descriptions of interventions in-context?** Consider the verbalization
 409 of an intervention in Figure 1. It mentions the $do(\cdot)$ operator and refers to a “perfect intervention”
 410 to prompt an LLM to rely on facts about causal reasoning that could have appeared in the training
 411 dataset. We ask whether LLMs can achieve the same performance on the intervention effects in the
 412 random benchmark if we varied the verbalization to instead “teach” an LLM [Patel et al., 2023] about
 413 a new graphical operation that behaves identically to an intervention. We randomly generate strings
 414 to instantiate this operation. Figure 3 in the appendix illustrates how we generate prompts for this
 415 task which we refer to as the substitution task. Table 9 summarizes IE prediction accuracy for the
 416 substitution task on the **Random** benchmark. Contrasting these results against Table 1, we see that
 417 for GPT-3.5, GPT-4 and GPT-4o, the performance generally suffers. However, GPT-4-turbo appears
 418 to be robust to changes in the way interventions are described.