# FROM CORPORA TO CAUSALITY: UNVEILING CAUSAL COMPREHENSION IN LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

## ABSTRACT

This study investigates the efficacy of Large Language Models (LLMs) in causal discovery. Using newly available open-source LLMs, OLMo and BLOOM, which provide access to their pre-training corpora, we explore three research questions aimed at understanding how LLMs process causal discovery. These questions focus on the impact of memorization versus generalization, the influence of incorrect causal relations in pre-training data, and the role of contexts of causal relations. Our findings indicate that while LLMs are effective in recognizing causal relations that occur frequently in pre-training data, their ability to generalize to new or rare causal relations is limited. Moreover, the presence of incorrect causal relations significantly undermines the confidence of LLMs in corresponding correct causal relations, and the context of a causal relation markedly affects the performance of LLMs to identify causal relations. This study shows that LLMs possess a limited capacity to generalize novel causal relations. It also highlights the importance of managing incorrect causal relations in pre-training data and integrating contextual information to optimize LLM performance in causal discovery tasks.

025 026 027

028

004

010 011

012

013

014

015

016

017

018

019

021

## 1 INTRODUCTION

029 Identification and understanding of causal relations hold fundamental importance in human cogni-030 tion and science, as those relations form the basis of causal models, which are utilized to answer 031 observational, interventional, and counterfactual questions (Zanga et al., 2022; Wan et al., 2024). 032 The task of identifying causal relations among a set of random variables is known as causal dis-033 *covery*, where a random variable may refer to an event in daily life, a medical treatment, or a drug 034 effect, etc. (Pearl, 2009; Peters et al., 2017; Nogueira et al., 2021). For decades, various statistical methods have been developed to identify causal relations from observational or interventional data 035 Heckerman et al. (1995); Chickering (2002); Koivisto & Sood (2004); Mooij et al. (2016a). However, algorithms that can accurately recover true causal structures from observational data remain 037 elusive. Neal (2020).

With the rise of Large Language Models (LLMs), recent studies exploit the potential of LLMs for causal discovery by evaluating them on benchmark datasets Willig et al. (2022); Ban et al. (2023). 040 Closed-source LLMs, such as GPT-3 and GPT-4, surpass the state-of-the-art (SOTA) statistical 041 methods on several publicly available datasets (Kıcıman et al., 2023). However, Romanou et al. 042 (2023) notice both GPT-3 and GPT-4 have a performance drop on the causal relations involving 043 real-world events occurring post-Jan 2022, compared to the ones before Jan 2022. Kıcıman et al. 044 (2023) find out that given part of a data table in the Tübingen cause-effect pairs dataset (Mooij et al., 045 2016b), GPT-4 can recover 61% of the remaining part. Zečević et al. (2023) conjecture that LLMs 046 may just recall causal knowledge in their large pre-training corpora by acting as "causal parrots". 047 However, there are no solid experiments to verify to what extent memorization and generalization 048 affect model performance in causal discovery tasks because the pre-training corpora of those LLMs are not accessible and the high-performing LLMs are closed-source.

The recently released open-source LLMs OLMo and BLOOM make their respective pre-training corpora Dolma and ROOTS publicly available Groeneveld et al. (2024); Workshop et al. (2023).

<sup>&</sup>lt;sup>1</sup>The code and data are available at https://anonymous.4open.science/r/causality\_ llm-5FD3

This provides the opportunity for us to investigate the correlations between model outputs and the frequency of relations mentioned in their pre-training corpora. In this work, we focus on three research questions and try to conduct experiments to answer them. *RQ 1*) What is the difference in *performance on causal discovery tasks for LLMs when recognizing relations through memorization compared to inferring them through generalization? RQ 2 How does the occurrence of incorrect causal relations affect LLMs performance in causal discovery?* and *RQ 3 How does the context of a causal relation influence LLM performance in causal discovery tasks?* 

- Our experiments reveal the following findings.
  - Although LLMs are proficient at recognizing causal relations through memorization, their ability to generalize novel causal relations is highly limited. This limitation poses significant challenges for deploying LLM-based causal discovery methods in scenarios where causal relations are rarely or not included in their pre-training data.
  - The presence of incorrect causal relations, such as the reversal of correct causal relations, adversely impacts LLMs' confidence in identifying correct causal relations. This find-ing highlights the necessity of minimizing conflicting causal information in pre-training datasets to enhance the performance of LLMs.
  - The validity and strength of causal relations can vary significantly across different contexts. This variability suggests that LLM-based causal discovery methods should incorporate the context of causal relations as input to ensure accuracy, particularly to avoid misleading contexts that could substantially degrade performance.
- 074 075

063

064

065

066 067

068

069

070 071

073

# 076 2 BACKGROUND

078 Causal discovery aims to identify causal relations among a given set of random variables. For each 079 pair of variables X and Y, it identifies whether  $X \leftarrow Y, Y \leftarrow X$ , or there is no causal influence between them, where  $\leftarrow$  denotes the direction of causality. The traditional algorithms for this task 081 are statistical methods that perform causal discovery on tabular data, which are capable of unveiling previously unknown or uncertain causal relations that are not *explicitly* mentioned anywhere 083 in text (e.g., "sea level pressure causally influences zonal wind at 10 m" Huang et al. (2021)). In contrast, prior NLP methods focus on either extracting mentions of known causal relations from 084 documents Yang et al. (2022) or answering questions related to causality Oh et al. (2013). The gold 085 standard for causal discovery is experimental approaches such as randomized controlled trials and A/B testing Fisher (1935). However, such experiments are often not feasible due to ethical or fi-087 nancial constraints, which necessitates the use of alternative methods that rely solely on statistics 088 collected from observational data. 089

The statistical causal discovery methods are conventionally categorized into constraint-based meth-090 ods, such as Peter and Clark (PC) Spirtes et al. (2000) and inductive causation (IC) Pearl (2009), 091 and score-based methods Heckerman et al. (1995); Chickering (2002); Koivisto & Sood (2004); 092 Mooij et al. (2016a). Those methods rely on statistics calculated from tabular data to infer causal graphs, in which random variables are depicted as nodes and their causal relations are represented as 094 edges. However, a significant drawback of these approaches is their dependency on extensive data 095 collection to construct reliable tabular data, a process that can be both time-consuming and costly. 096 Furthermore, a theoretical limitation of these statistical methods is their inability to precisely predict 097 ground-truth causal graphs, unless strong assumptions are made. Instead, they typically yield an 098 equivalence class of true causal graphs Spirtes et al. (2000); Pearl (2009).

099 Recent advances of LLMs provide new opportunities to tackle the task without accessing tabular data 100 by formulating it as a pairwise causal relation prediction task Kıcıman et al. (2023); Zečević et al. 101 (2023); Long et al. (2022). Given a pair of variable names, an LLM is instructed to identify which 102 is the cause and which is the effect using prompts K1c1man et al. (2023); Zečević et al. (2023), 103 by distilling such knowledge directly from the LLM. However, the reliability of such methods is 104 under scrutiny. Zečević et al. (2023) argue that LLMs are "causal parrots", which may depend 105 on *memorization* to recall the causal relations present in their training data. In other words, LLMs may not generalize well to detect causal relations that seldom or never occur in pre-training data. 106 If this argument holds, LLMs may primarily excel at reproducing causal relations known in their 107 training data rather than uncovering novel ones. However, there is no solid empirical justification

of this argument because prior works employ either commercial LLMs or open-source LLMs that have no access to their training data. The current techniques for understanding and investigating memorization in LLMs are still in their infancy Speicher et al. (2024).

111

#### 112 113 114

# 3 Methodology

115 116

We aim to investigate the key limitations of LLMs for causal discovery by answering three research questions. The first research question aims to collect stronger empirical evidence to verify the "causal parrots" hypothesis. The second research question investigates to what extend the presence of incorrect causal relations in the training data, which are oriented in the opposite direction to their true counterparts, influence the performance of LLMs. Instead of only feeding two variable names to LLMs, the last research question is concerned with the *first* quantitative study on how the context of a causal relation impacts the predictive performance of LLMs.

124 Unlike prior works, we collect evidence of memorization from the pre-training datasets of LLMs 125 and investigate their statistical properties in relation to LLMs' predictive performance. As it is al-126 most infeasible to collect all mentions of a causal relation from a dataset, we curate a synthetic 127 causal relation dataset to further investigate to what extent LLMs can generate to unseen causal 128 relations. Herein, we select OLMo-7b-Instruct and BLOOM-7b1, which are the LLMs that have 129 their pre-training data publicly available Groeneveld et al. (2024); Workshop et al. (2023). To determine whether LLMs primarily rely on memorization or generalization, we classify causal relations 130 into various occurrence intervals, ensuring that each interval contains a comparable number of rela-131 tions. We then assess the LLMs' performance in recognizing causal relations across these defined 132 intervals. Evaluations are performed by transforming causal relations into yes-no questions, such 133 as "does smoking cause lung cancer?". We employ accuracy and F1 score metrics to assess perfor-134 mance. If LLMs mainly utilize memorization to identify causal relations, we anticipate observing 135 high accuracy and F1 scores for relations that frequently occur in the pre-training data, with a notable 136 decline in performance for less frequently occurring relations. This experimental method aligns with 137 the approaches stated in Razeghi et al. (2022). 138

In addition to examining the frequency of causal relations, we also investigate how the presence 139 of incorrect causal relations impacts LLMs' confidence in corresponding correct causal relations. 140 For example, we want to explore how the occurrence of "lung cancer causes smoking" might affect 141 an LLM's confidence in the correct relation "smoking causes lung cancer." To this end, we have 142 devised a novel experimental setup. We assess the confidence of LLMs in correct causal relations 143 under varying frequencies of corresponding incorrect causal relations. We hypothesize that a higher 144 presence of incorrect causal relations diminishes the LLMs' confidence in the correct causal rela-145 tions. The confidence level of the LLMs is measured by the proportion of responses that affirm the 146 correct causal relation out of multiple generated responses for one query.

Due to the impracticality of exhaustively retrieving all semantically equivalent mentions of target causal relations in pre-training data, we create new pre-training corpora including synthetic causal and incorrect causal relations. These relations, such as "blaonge causes goloneke," utilize terms that do not exist in the original pre-training corpora. We then integrate these synthetic relations into the pre-training data at various frequencies. This approach allows us to re-evaluate our experimental results on real-world causal relations, thereby validating the reliability of our findings under controlled conditions.

154 One distinction between causal discovery and numerical reasoning tasks Razeghi et al. (2022) is 155 context dependency. Numerical reasoning, such as 3 + 4 = 7, exhibits consistency across various 156 contexts. However, causal relations might not have this consistency. For example, the causal relation 157 "rain causes flooding" may be true during a heavy downpour in a city with poor drainage but may not 158 be true during light rain in areas with good drainage systems. Therefore, we assess the performance 159 of LLMs on causal discovery across varying contexts. For each selected causal relation from humanannotated datasets, we employ GPT-40 to generate five positive contexts that affirm the relation and 160 five negative contexts that challenge it. The LLMs' ability to recognize these causal relations is then 161 evaluated in these contexts.



Figure 1: The average F1 score and accuracy of OLMo-7b-Instruct by occurrence interval on full causal discovery tasks, where F1 and accuracy are computed from 0 to 4 ICL examples. The occurrence data of (a) and (b) are derived from the exact matching query, while the occurrence data of (c) and (d) are derived from the "event A"  $\Rightarrow$  "causes"  $\Rightarrow$  "event B" query. An asterisk (\*) denotes that the p-value of the correlation coefficients is less than 0.05.



Figure 2: The average F1 score and accuracy of BLOOM-7b1 by occurrence interval on full causal discovery tasks, averaged across 0 to 4 ICL examples. The occurrence data are derived from the exact matching query.

# 4 EXPERIMENTAL SETUP

In this section, we outline the details of our experimental setup.

# 4.1 DATASETS

**Tasks.** Following (Kıcıman et al., 2023), we consider the following two causal discovery tasks. *Causal Direction Identification*. Given two causally related variables (X, Y), the causal direction identification task involves deciding whether  $X \to Y$  or  $X \leftarrow Y$  is true. *Full Causal Discovery*. Given a set of random variables **X**, for each possible pair of variables  $(X_i, X_j)$ , an LLM is instructed to identify whether:  $X_i \to X_j$ ,  $X_i \leftarrow X_j$ , or no causal relation between  $X_i$  and  $X_j$ . The causal direction identification and full causal discovery tasks can be treated as classification tasks. Therefore, we evaluate the results using F1 and accuracy.

204 205

177

179

181

183

185

187

188

189 190 191

192 193

194 195

196

4.1.1 REAL-WORLD DATA

 Causal Direction Identification. For this task, we consider two datasets derived from Concept-Net Speer et al. (2017) and CauseNet Heindorf et al. (2020). From ConceptNet, we select the top 1,900 causal relations based on confidence and generate an equal number of reverse-causal relations by swapping the cause and effect, resulting in 3,800 causal and reverse-causal relations. From CauseNet, we select 814 high-confidence causal relations and create an equal number of reversecausal relations, totaling 1,628 relations. These procedures are detailed in Appendix A.2.

212

Full Causal Discovery. We consider six datasets for this task. We utilize four small causal graphs within the medical literature as our ground-truth causal graphs, which include Alcohol, Cancer, Diabetes, and Obesity (see Fig. 10) Hernán et al. (2004); Long et al. (2022). We also use a causal graph from atmospheric science, named Arctic Sea Ice Huang et al. (2021). This causal graph



Figure 3: The average F1 score and accuracy of OLMo-7b-Instruct by occurrence interval on causal direction identification task, averaged across 0 to 4 ICL examples. The occurrence data are derived from the exact matching query in the Dolma pre-training corpus.



Figure 4: The average F1 score and accuracy of OLMo-7b-Instruct by occurrence interval on causal direction identification task, averaged across 0 to 4 ICL examples. The occurrence data are derived from the "event A"  $\Rightarrow$  "causes"  $\Rightarrow$  "event B" query in the Dolma pre-training corpus.

explores the factors influencing arctic sea ice coverage. The Arctic Sea Ice is based on expert knowledge and consists of a causal graph with 12 variables and 46 edges, each edge derived from textbooks and peer-reviewed publications (see Fig. 11). Then, we employ a larger causal graph used for evaluating car **Insurance** risks Binder et al. (1997), which comprises 27 variables and 52 edges (see Fig. 12).

249 250 251

226

227

228

238

239

240

#### 4.1.2 SYNTHETIC DATA

252 **Causal Direction Identification.** We create a pre-training dataset including synthetic correct and 253 incorrect causal relations that are absent in the original corpora. This dataset includes 100,000 documents randomly sampled from Dolma, with incorrect causal relations that either swap the positions 254 of cause and effect or use negation templates such as "X does not cause Y." We generate 100 ar-255 tificial causal relations using fictitious terms like 'blaonge' and 'goloneke'. Utilizing predefined 256 templates listed in Table 5 in Appendix A.5, we craft mentions for both correct and incorrect causal 257 relations. Then we create positive documents containing correct causal relations and negative doc-258 uments containing incorrect causal relations by inserting these mentions between sentences within 259 the documents. We adopt three approaches for the insertion of mentions. Correct Relation Scaling: 260 we vary the insertion of each correct causal relation from 0 to 1,000 occurrences. **Reverse Relation** 261 Scaling: we first insert 1000 occurrences of each correct causal relation followed by inserting the 262 corresponding reverse causal relations from 0 to 1,000 occurrences. Negated Relation Scaling: 263 After inserting 1,000 occurrences of each correct causal relation, we insert negations of these causal 264 relations, from 0 to 1,000 occurrences. We then fine-tune OLMo-7b-Instruct utilizing LoRA Hu et al. (2022) on synthetic datasets, with details provided in Appendix A.6. 265

266 267

- 4.2 MODELS
- Large Language Models. We conduct experiments using the following language models: OLMo-7b-Instruct Groeneveld et al. (2024), BLOOM-7b1 Workshop et al. (2023), Llama2-7b-chat Meta

Spearman



Pearson (r=0.61); Spear-

man (r=0.90\*)

man (r=1.0\*)

278 279

(r=0.61);

(r=0.90\*)

280 281

282

283 284

287

289

290

291

292 293

294

295 296 297

298

299

300

301

302

303

Figure 5: The average F1 score and accuracy of BLOOM-7b1 by occurrence interval on causal direction identification task, averaged across 0 to 4 ICL examples. The occurrence data are derived from the exact matching query in the ROOTS pre-training corpus.

(r=0.07);

(r=0.81\*)

Spearman

son (r=0.01); Spearman

(r=0.69\*)



Pearson (r=0.94); Spearrence; Pearson (r=0.97\*); Spearman(r=1.0\*)

Figure 6: The average F1 score and accuracy of fine-tuned OLMo-7b-Instruct by various occurrences on synthetic causal relations, averaged across 0 to 4 ICL examples.

(2023), Llama3-8b-Instruct Meta (2024), GPT-3.5-turbo OpenAI (2022) and GPT-40 OpenAI (2024). OLMo-7b-Instruct and BLOOM-7b1 provide access to both their pre-training corpora and model weights. Llama2-7b-chat and Llama3-8b-Instruct have only released their model weights. GPT-3.5-turbo and GPT-4o are closed-source models. OLMo-7b-Instruct was pre-trained using the Dolma dataset Soldaini et al. (2024), while BLOOM-7b1 utilized the ROOTS corpus Laurençon et al. (2022). The release of corresponding search tools, WIMBD Elazar et al. (2024) for Dolma and ROOTS Search Piktus et al. (2023) for ROOTS, enables the searching for causal relations.

304 **In-Context Learning and Prompt.** For both the causal direction identification and the full causal 305 discovery tasks, we utilize similar in-context learning demonstrations and prompts, detailed further 306 in Appendix A.3. When evaluating a pair of variables (X, Y), we pose two questions to the LLMs: 307 "Does X cause Y?" and "Does Y cause X?" The LLMs are expected to generate step-by-step expla-308 nations and provide a final response of either 'yes' or 'no'. 309

310 4.3 RETRIEVAL QUERY 311

The pre-training corpus for OLMo-7b-Instruct is Dolma Soldaini et al. (2024), which has a search 312 tool named WIMBD Elazar et al. (2024). In our usage of WIMBD, we implement two search 313 queries: an exact match for "event A causes event B"; an ordered phrase search for "event A"  $\Rightarrow$ 314 "causes"  $\Rightarrow$  "event B". Here,  $X \Rightarrow Y$  indicates that X occurs before Y within a predefined window 315 of text. The search tool for BLOOM-7b1 pre-training corpus ROOTS Laurençon et al. (2022) is 316 ROOTS Search Piktus et al. (2023). Due to its limited search capability, we only utilize exact match 317 in ROOTS Search. In Table 3, 4 in Appendix A.4, we detail the methods used to create queries for 318 retrieving causal relations. 319

320 321

- 5 **EXPERIMENTAL RESULTS**
- **Research Question 1.** What is the difference in performance on causal discovery tasks when LLMs 323 recognize relations through memorization compared to inferring them through generalization?

331

332

333

369



Figure 7: The average confidence of correct causal relations on OLMo-7b-Instruct (a) and BLOOM-7b1 (b) by reverse casual relation occurrence ratio intervals on full causal discovery tasks.

Relations frequently occurring in pre-training data are likely memorized by LLMs. However, relations that are seldom or never present in pre-training data require LLMs to generalize these relations.

To address RQ 1, we evaluate LLMs on causal relations across different occurrence intervals, which contain the similar number of causal relations. Causal relations with high occurrences are likely to be memorized by LLMs, whereas those with low occurrences reveal LLMs' generalization ability Carlini et al. (2023). We then analyze the correlation between the occurrence of causal relations and the performance of LLMs on these causal relations.

Real-World Data We compute the average F1 and accuracy at each occurrence interval over various 345 numbers of ICL examples (i.e., from 0-shot to 4-shot). The results are plotted with the x-axis 346 representing occurrence intervals and the y-axis representing F1 or accuracy. Fig. 1, 2, 3, 4 and 347 5 show that both F1 and accuracy exhibit a strong positive correlation with occurrence in the pre-348 training corpora. For instance, in the full causal discovery task, the Spearman correlation between 349 F1 scores and occurrence rates is 0.9 using OLMo-7b-Instruct and its pre-training data. Compared 350 to highly frequent causal relations, LLMs exhibit significantly poorer performance when identifying 351 low-frequency causal relations. For instance, in a full causal discovery task, OLMo-7b-Instruct achieves an F1 score of 0.88 in the highest occurrence interval, but only 0.2 in the lowest occurrence 352 interval. In the causal direction identification task, OLMo-7b-Instruct reaches a 0.93 F1 score at the 353 highest occurrence interval, compared to just 0.35 at the lowest. These results indicate that LLMs 354 have limited generalization ability in causal discovery tasks. 355

*Synthetic Data* We fine-tune OLMo-7b-Instruct with Correct Relation Scaling. Fig. 6 demonstrates
 that both F1 and accuracy have a strong positive correlation with occurrence within the pre-training
 corpora, which aligns with real-world data.

359 Discussion These results demonstrate that while LLMs excel at recognizing causal relations through 360 memorization, their capacity to generalize from less frequent or entirely novel data remains highly 361 constrained. This limitation highlights the challenges in deploying LLMs in scenarios where causal 362 relations are novel and absent from their pre-training data. Furthermore, this suggests the necessity 363 of traditional statistical methods for causal discovery that rely solely on statistics to determine causal relations, irrespective of the novelty of causal relations. This insight suggests that future research 364 might explore integrating traditional statistical methods with LLMs to enhance their generalization 365 ability. 366

Research Question 2. How does the occurrence of incorrect causal relations affect LLMs in causal discovery tasks?

incorrect causal relations include reversals of correct causal relations (e.g., lung cancer causes smoking) and negations of correct causal relations (e.g., smoking does not cause lung cancer).

We hypothesize that when both correct and incorrect causal relations are frequent, LLMs may struggle to discern the correct relations, thereby reducing their confidence in correct causal relations. To investigate this, we examine the correlation between the occurrence ratio of incorrect causal relations and LLMs' confidence in correct causal relations. The occurrence ratio is defined as the number of incorrect causal relations divided by the number of corresponding correct causal relations. Confidence in correct causal relations (*i.e.*, affirmative confidence) is measured by the proportion of affirmative responses among multiple generated responses, where a response is considered affirmative



(a) Confidence vs Occurrence at ConceptNet; Pearson (r=-0.98\*), Spearman (r=-0.4)

385

386

387 388

389

390

391 392 393

396 397

398

399

400 401

402

403

404

429

(b) Confidence vs Occurrence at CauseNet; Pearson (r=-0.71\*), Spearman (r=-0.51\*) (c) Confidence vs Occurrence at ConceptNet; Pearson (r=-0.86), Spearman (r=-0.89\*)

(d) Confidence vs Occurrence at CauseNet; Pearson (r=-0.58), Spearman (r=-0.6)

Figure 8: The average confidence of correct causal relations on OLMo-7b-Instruct (a,b) and BLOOM-7b1 (c,d) by reverse casual relation occurrence ratio intervals on causal direction identification task, averaged across 0 to 4 ICL examples.



Figure 9: The average confidence of correct causal relations on fine-tuned OLMo-7b-Instruct by reverse casual relation occurrence ratio (a) and negation casual relation occurrence ratio (b) on synthetic causal relations, averaged across 0 to 4 ICL examples.

if it contains "yes" and negative if it contains "no". If neither "yes" nor "no" appears in an answer,
we classify it as a 'fail'. The average proportion of 'fail' across all datasets is 0.03, indicating that
most responses are either 'yes' or 'no'. For example, if the phrase "smoking causes lung cancer"
appears 13,652 times and its reverse "lung cancer causes smoking" appears 99 times, the resulting
occurrence ratio is approximately 0.007. If the query "Does smoking cause lung cancer?" results
in affirmative responses in 8 out of 10 generation samples, the affirmative confidence for "smoking
causes lung cancer" is 0.8. In this experiment, we sample 10 responses for each query.

<u>*Real-World Data*</u> We calculate and plot the correlation between different intervals of occurrence ratios of incorrect causal relations and affirmative confidence. The experiment results, shown in Fig. 7 and 8, indicate a negative correlation, showing that LLMs' confidence in correct causal relations decreases as the occurrence ratio of incorrect causal relations increases.

416 Synthetic Data We fine-tune OLMo-7b-Instruct employing both Reverse Relation Scaling and
 417 Negated Relation Scaling. Fig. 9 shows a similar negative correlation with real-world data: as
 418 the occurrence of incorrect causal relations increases, there is a decline in the LLMs' confidence in
 419 the corresponding correct causal relations.

<u>Discussion</u> This negative correlation suggests that while LLMs excel at memorizing frequently oc curring information, they struggle to discern the correct relation when confronted with high fre quencies of conflicting data. This inability leads to a loss of confidence in correct causal relations.
 This finding underscores the necessity of not only enhancing the presence of correct information
 but also of eliminating misinformation in pre-training data. Furthermore, these results pave the way
 for future research aimed at developing models that can manage conflicting information within their
 pre-training data.

Research Question 3. How does the context of a causal relation influence LLM performance in causal discovery tasks?

We hypothesize the strength and validity of causal relations can vary across different contexts. Thus,
 when a causal discovery question is presented with different contexts, LLMs might provide different and sometimes opposite answers to the causal relation's validity.

432		Full Ca	usal Disc	covery
433		w/o Ctx	P.Ctx	N.Čtx
434	OLMo-7b-Instruct (3 ICL)	0.65	0.875	0.421
435	BLOOM-7b1 (3 ICL)	0.629	0.76	0.597
436	Llama2-7b-chat (3 ICL)	0.682	0.852	0.255
437	Llama3-8b-Instruct (3 ICL)	0.67	0.738	0.207
438	GPT-3.5-turbo (3 ICL)	0.652	0.86	0.242
439	GPT-4o (3 ICL)	0.69	0.92	0.272
440		Co	onceptNe	t
441		w/o Ctx	P.Ctx	N.Ctx
442	OLMo-7b-Instruct (3 ICL)	0.9	0.95	0.624
443	BLOOM-7b1 (3 ICL)	0.79	0.81	0.704
440	Llama2-7b-chat (3 ICL)	0.79	0.952	0.318
444	Llama3-8b-Instruct (3 ICL)	0.66	0.85	0.104
445	GPT-3.5-turbo (3 ICL)	0.77	0.906	0.338
446	GPT-40 (3 ICL)	0.87	0.962	0.346
447		C	auseNet	
448		w/o Ctx	P.Ctx	N.Ctx
449	OLMo-7b-Instruct (3 ICL)	0.89	0.992	0.616
450	BLOOM-7b1 (3 ICL)	0.72	0.784	0.632
451	Llama2-7b-chat (3 ICL)	0.92	0.998	0.472
452	Llama3-8b-Instruct (3 ICL)	0.88	0.946	0.144
453	GPT-3.5-turbo (3 ICL)	0.93	0.982	0.674
454	GPT-4o (3 ICL)	0.98	0.998	0.602

Table 1: Affirmative ratio of LLMs on causal relations across different contexts.

457 From ConceptNet and CauseNet, we select 100 high-confidence correct causal relations from each. 458 Since both ConceptNet and CauseNet lack context information, for each causal relation, we use 459 GPT-40 to generate five positive contexts that enhance it and five negative contexts that weaken 460 it. Then we hire thirteen annotators to evaluate these causal relations under different contexts in three rounds. The prompt and evaluation details are presented in Appendix A.7. The agreement 461 between annotators and GPT-40 is 0.76 using Krippendorff's Alpha Castro (2017). We then assess 462 the performance of LLMs on these causal relations within positive and negative contexts. The query 463 format is similar to Table 2, except we provide context information using the phrase "Given the 464 scenario: {description}". We assess LLM performance on correct causal relations within various 465 contexts using the affirmative ratio. This ratio is calculated by dividing the number of correct causal 466 relations identified by the LLM by the total number of correct causal relations presented. 467

<u>Observation</u> From the results in Table 1, we observe that all LLMs are more likely to identify causal relations in positive contexts compared to no context. In contrast, adding negative contexts significantly decreases LLMs' ability to identify causal relations compared to no context. These results indicate that the validity and strength of causal relations can vary in different contexts.

<u>Discussion</u> The significant variation in causal relation identification across positive and negative
 contexts indicates the context sensitivity of LLM-based causal discovery methods. This observa tion suggests that LLM-based algorithms should explicitly provide contextual information to enable
 LLMs to better understand the scenario and thereby make more accurate predictions. It is par ticularly crucial for these algorithms to avoid misleading contexts, as our results demonstrate that
 negative contexts can substantially impair LLM performance. Furthermore, investigating the under lying mechanisms of how different contexts influence the strength and validity of causal relations
 could be a promising direction for future research.

479 480

# 6 RELATED WORK

481 482

Causality with LLMs Kıcıman et al. (2023); Zečević et al. (2023); Long et al. (2022); Feng et al. (2023) explore the inference of causal relations by submitting pairwise queries about variable pairs to LLMs. These queries are either structured as option selection questions Kıcıman et al. (2023) or yes-no questions Long et al. (2022); Zečević et al. (2023). Results from these experiments demon-

486 strate that the LLM-based approach surpasses traditional statistical algorithms in performance. Re-487 markably, the LLM-based method requires only the names of the variables, without needing their 488 statistical data. However, the approach of pairwise queries may lead to inefficiencies in time and 489 computation, as identifying all possible relations among a set n of variables necessitates  $O(n^2)$ 490 queries. To address this, Jiralerspong et al. (2024) have proposed a breadth-first search strategy that significantly reduces the number of queries to a linear scale. Additionally, beyond exploring 491 relationships among observable variables, Liu et al. (2024) has developed a framework capable of 492 uncovering high-level hidden variables from unstructured data using LLMs, and subsequently infer-493 ring causal relationships. 494

495

496 Influence of Pre-training Data on Language Models. Research conducted by Kassner et al. 497 (2020) and Wei et al. (2021) involving controlled variations in pretraining data sheds light on its 498 impact on language models' (LM) capabilities to memorize factual information and understand syn-499 tactic rules. Their findings confirm that the frequency of data plays a crucial role in determining a model's ability to remember specific facts or grammatical structures about verb forms. Furthermore, 500 Sinha et al. (2021); Min et al. (2022) show that altering the word order during pretraining barely af-501 fects the LMs' performance in subsequent tasks, and mixing up labels in in-context learning scenar-502 ios does not significantly affect the models' few-shot learning accuracy. These studies collectively 503 indicate that the efficacy of LMs predominantly hinges on their capacity to process complex word 504 co-occurrence patterns. Additionally, Carlini et al. (2023; 2019); Song & Shmatikov (2019) have 505 identified that LMs can retain sensitive information from their training datasets, even when such in-506 stances are infrequent. The experiments of Razeghi et al. (2022) demonstrate that models are more 507 accurate on numerical reasoning questions whose terms are more prevalent in pre-training data.

508 509

# 7 CONCLUSION

514

515

516

517

In this study, we investigate the factors that impact the performance of LLMs in causal discovery tasks. Our results show that the frequency of causal relations within a model's pre-training data has a positive correlation with LLM performance, while the presence of incorrect causal relations can negatively affect the models' confidence in correct causal relations. Furthermore, our experiments reveal that the context of causal relations significantly affects the validity of causal relations. To facilitate a deeper understanding of LLMs, we strongly advocate for the release of both model weights and pre-training data by more LLM providers.

518 519 520

# 521 522

523

524

525

526

527

528

# Reproducibility Statement

We release our code and scripts at https://anonymous.4open.science/r/ causality\_llm-5FD3. Appendix A.1 presents the ground-truth causal graphs used in our full causal discovery task. Appendix A.3 provides examples of in-context learning and the prompts used during our experiments. Appendix A.4 outlines the queries utilized for searching within the pre-training data. Appendix A.7 details the human evaluation process for assessing causal relations under various contexts. These resources ensure transparency and facilitate the replication of our research findings.

529 530 531

# 532 ETHICS STATEMENT

533

The ability of LLMs to identify and generalize causal relations could significantly impact various fields, from healthcare to social sciences, where understanding causality is crucial. However, we acknowledge that relying on LLMs for causal discovery may perpetuate existing biases present in training data, potentially leading to misleading or harmful conclusions if deployed without proper safeguards. The limitations we've identified, particularly regarding incorrect causal relations and context dependency, underscore the need for careful human oversight when applying these models in real-world scenarios.

# 540 REFERENCES

- Taiyu Ban, Lyvzhou Chen, Xiangyu Wang, and Huanhuan Chen. From query tools to causal architects: Harnessing large language models for advanced causal discovery from data. *arXiv preprint arXiv:2306.16902*, 2023.
- John Binder, Daphne Koller, Stuart Russell, and Keiji Kanazawa. Adaptive probabilistic networks with hidden variables. *Machine Learning*, 29:213–244, 1997.
- 548 Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer:
  549 evaluating and testing unintended memorization in neural networks. In *Proceedings of the 28th* 550 USENIX Conference on Security Symposium, SEC'19, pp. 267–284, USA, 2019. USENIX Asso-551 ciation. ISBN 9781939133069.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=TatRHT\_1cK.
- Santiago Castro. Fast Krippendorff: Fast computation of Krippendorff's alpha agreement measure.
   https://github.com/pln-fing-udelar/fast-krippendorff, 2017.
- Jiuhai Chen and Jonas Mueller. Quantifying uncertainty in answers from any language model and enhancing their trustworthiness, 2023. URL https://arxiv.org/abs/2308.16175.
- David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002.
- Yanai Elazar, Akshita Bhagia, Ian Helgi Magnusson, Abhilasha Ravichander, Dustin Schwenk,
   Alane Suhr, Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hannaneh Ha jishirzi, Noah A. Smith, and Jesse Dodge. What's in my big data? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?
   id=RvfPnOkPV4.
- Tao Feng, Lizhen Qu, and Gholamreza Haffari. Less is more: Mitigate spurious correlations for open-domain dialogue response generation models by causal discovery. *Transactions of the Association for Computational Linguistics*, 11:511–530, 2023. doi: 10.1162/tacl\_a\_00561. URL https://aclanthology.org/2023.tacl-1.30.
  - Tao Feng, Lizhen Qu, Zhuang Li, Haolan Zhan, Yuncheng Hua, and Gholamreza Haffari. Imo: Greedy layer-wise sparse representation learning for out-of-distribution text classification with pre-trained models, 2024. URL https://arxiv.org/abs/2404.13504.
  - R. A. Fisher. *The Design of Experiments*. Oliver and Boyd, 1935.
  - Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koeppl, Preslav Nakov, and Iryna Gurevych. A survey of confidence estimation and calibration in large language models, 2024. URL https://arxiv.org/abs/2311.08298.
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. Olmo: Accelerating the science of language models, 2024. URL https://arxiv.org/abs/2402.00838.
- 592

559

575

576

577 578

579 580

581

582

583

<sup>3</sup> David Heckerman, Dan Geiger, and David M Chickering. Learning bayesian networks: The combination of knowledge and statistical data. *Machine learning*, 20:197–243, 1995.

- Stefan Heindorf, Yan Scholten, Henning Wachsmuth, Axel-Cyrille Ngonga Ngomo, and Martin Potthast. Causenet: Towards a causality graph extracted from the web. In *Proceedings of the* 29th ACM International Conference on Information & Knowledge Management, CIKM '20, pp. 3023–3030, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450368599. doi: 10.1145/3340531.3412763. URL https://doi.org/10.1145/ 3340531.3412763.
- Miguel A Hernán, Sonia Hernández-Díaz, and James M Robins. A structural approach to selection bias. *Epidemiology*, 15(5):615–625, 2004.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?
   id=nZeVKeeFYf9.
- Yiyi Huang, Matthäus Kleindessner, Alexey Munishkin, Debvrat Varshney, Pei Guo, and Jianwu Wang. Benchmarking of data-driven causality discovery approaches in the interactions of arctic sea ice and atmosphere. *Frontiers in big Data*, 4:642182, 2021.
- Thomas Jiralerspong, Xiaoyin Chen, Yash More, Vedant Shah, and Yoshua Bengio. Efficient
   causal graph discovery using large language models, 2024. URL https://arxiv.org/abs/
   2402.01207.
- Nora Kassner, Benno Krojer, and Hinrich Schütze. Are pretrained language models symbolic reasoners over knowledge? In Raquel Fernández and Tal Linzen (eds.), *Proceedings of the 24th Conference on Computational Natural Language Learning*, pp. 552–564, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.conll-1.45. URL https://aclanthology.org/2020.conll-1.45.
- Mikko Koivisto and Kismat Sood. Exact bayesian structure discovery in bayesian networks. *The Journal of Machine Learning Research*, 5:549–573, 2004.
- Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language models: Opening a new frontier for causality, 2023. URL https://arxiv.org/abs/2305.00050.
- 626 Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, 627 Jörg Frohberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gérard Dupont, Stella 628 Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, Giada Pistilli, Olivier Nguyen, Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa, Paulo Villegas, Tris-630 tan Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel Romero Muñoz, Jian Zhu, 631 Daniel Van Strien, Zaid Alyafeai, Khalid Almubarak, Vu Minh Chien, Itziar Gonzalez-Dios, Aitor 632 Soroa, Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan, Shamik Bose, David Ifeoluwa 633 Adelani, Long Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim, Violette Lepercq, Suzana Ilic, 634 Margaret Mitchell, Sasha Luccioni, and Yacine Jernite. The bigscience ROOTS corpus: A 1.6TB 635 composite multilingual dataset. In Thirty-sixth Conference on Neural Information Processing 636 Systems Datasets and Benchmarks Track, 2022. URL https://openreview.net/forum? 637 id=UoEw6KigkUn.
- Chenxi Liu, Yongqiang Chen, Tongliang Liu, Mingming Gong, James Cheng, Bo Han, and Kun
   Zhang. Discovery of the hidden world with large language models, 2024. URL https://
   arxiv.org/abs/2402.03941.
- 642
   643
   644
   645
   Stephanie Long, Tibor Schuster, and Alexandre Piché. Can large language models build causal graphs? In *NeurIPS 2022 Workshop on Causality for Real-world Impact*, 2022. URL https: //openreview.net/forum?id=LQQoJGw8JD1.
- 646 Meta. Llama 2: Open foundation and fine-tuned chat models, 647 2023. URL https://ai.meta.com/research/publications/ llama-2-open-foundation-and-fine-tuned-chat-models/.

Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024. URL https://ai.meta.com/blog/meta-llama-3/.

- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
  Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In
  Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 11048–11064, Abu Dhabi, United
  Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/
  2022.emnlp-main.759. URL https://aclanthology.org/2022.emnlp-main.759.
- Joris M. Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. Distinguishing cause from effect using observational data: methods and benchmarks. *J. Mach. Learn. Res.*, 17(1):1103–1204, jan 2016a. ISSN 1532-4435.
- Joris M Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. Distinguishing cause from effect using observational data: Methods and benchmarks. *Journal of Machine Learning Research*, 17(32):1–102, 2016b.
- Brady Neal. Introduction to Causal Inference from a Machine Learning Perspective. 2020. URL https://www.bradyneal.com/Introduction\_to\_Causal\_ Inference-Dec17\_2020-Neal.pdf.
- Ana Rita Nogueira, João Gama, and Carlos Abreu Ferreira. Causal discovery in machine learning: Theories and applications. *Journal of Dynamics and Games*, 8(3):203–231, 2021. ISSN 2164-6066. doi: 10.3934/jdg.2021008. URL /article/id/ 1efb412f-212e-4043-a129-7729066412c5.
- Jong-Hoon Oh, Kentaro Torisawa, Chikara Hashimoto, Motoki Sano, Stijn De Saeger, and Kiyonori
   Ohtake. Why-question answering using intra-and inter-sentential causal relations. In *Proceedings* of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long
   Papers), pp. 1733–1743, 2013.
- 675 676 OpenAI. Introducing chatgpt, 2022. URL https://openai.com/index/chatgpt/.
- 677 OpenAI. Hello gpt-40, 2024. URL https://openai.com/index/hello-gpt-40/.
- <sup>678</sup> Judea Pearl. *Causality*. Cambridge university press, 2009.

- Jonas Peters, Dominik Janzing, and Bernhard Schlkopf. *Elements of Causal Inference: Foundations and Learning Algorithms*. The MIT Press, 2017. ISBN 0262037319.
- Aleksandra Piktus, Christopher Akiki, Paulo Villegas, Hugo Laurençon, Gérard Dupont, Sasha Luccioni, Yacine Jernite, and Anna Rogers. The ROOTS search tool: Data transparency for LLMs. In Danushka Bollegala, Ruihong Huang, and Alan Ritter (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pp. 304–314, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-demo.29. URL https://aclanthology.org/2023.acl-demo.29.
- Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. Impact of pretraining term frequencies on few-shot numerical reasoning. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 840–854, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.59. URL https://aclanthology.org/2022.findings-emnlp.59.
- Angelika Romanou, Syrielle Montariol, Debjit Paul, Léo Laugier, Karl Aberer, and Antoine Bosselut. Crab: Assessing the strength of causal relationships between real-world events. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 15198–15216, 2023.
- 700 Marco Scutari. Learning bayesian networks with the bnlearn r package. Journal of Statistical 701 Software, 35(3):1-22, 2010. doi: 10.18637/jss.v035.i03. URL https://www.jstatsoft. org/index.php/jss/article/view/v035i03.

Koustuv Sinha, Robin Jia, Dieuwke Hupkes, Joelle Pineau, Adina Williams, and Douwe Kiela. Masked language modeling and the distributional hypothesis: Order word matters pre-training for little. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pp. 2888–2913, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.230. URL https: //aclanthology.org/2021.emnlp-main.230.

709 Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, 710 Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Harsh 711 Jha, Sachin Kumar, Li Lucy, Xinxi Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas 712 Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle 713 Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Pete Walsh, Luke 714 Zettlemoyer, Noah A. Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and 715 Kyle Lo. Dolma: an open corpus of three trillion tokens for language model pretraining research, 2024. URL https://arxiv.org/abs/2402.00159. 716

- Congzheng Song and Vitaly Shmatikov. Auditing data provenance in text-generation models. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '19, pp. 196–206, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450362016. doi: 10.1145/3292500.3330885. URL https://doi. org/10.1145/3292500.3330885.
- Robyn Speer, Joshua Chin, and Catherine Havasi. Conceptnet 5.5: an open multilingual graph of
   general knowledge. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*,
   AAAI'17, pp. 4444–4451. AAAI Press, 2017.
- Till Speicher, Mohammad Aflah Khan, Qinyuan Wu, Vedant Nanda, Soumi Das, Bishwamittra Ghosh, Krishna P Gummadi, and Evimaria Terzi. Understanding memorisation in Ilms: Dynamics, influencing factors, and implications. *arXiv preprint arXiv:2407.19262*, 2024.
- Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, Prediction, and Search.* MIT press, 2000.
- Guangya Wan, Yuqi Wu, Mengxuan Hu, Zhixuan Chu, and Sheng Li. Bridging causal discovery and large language models: A comprehensive survey of integrative approaches and future directions. *arXiv preprint arXiv:2402.11068*, 2024.
- Jason Wei, Dan Garrette, Tal Linzen, and Ellie Pavlick. Frequency effects on syntactic rule learning in transformers. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 932–948, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.72. URL https://aclanthology.org/2021.emnlp-main.72.

742

743

- Moritz Willig, Matej Zečević, Devendra Singh Dhami, and Kristian Kersting. Probing for correlations of causal facts: Large language models and causality. 2022.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, 745 Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick 746 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, 747 Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural 748 language processing. In Qun Liu and David Schlangen (eds.), Proceedings of the 2020 Confer-749 ence on Empirical Methods in Natural Language Processing: System Demonstrations, pp. 38– 750 45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020. 751 emnlp-demos.6. URL https://aclanthology.org/2020.emnlp-demos.6. 752
- BigScience Workshop, :, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić,
  Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé,
  Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji

756 Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, 758 Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, 760 Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar Natan, 761 Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza 762 Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Frohberg, Joseph Tobing, 764 Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon 765 Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, 766 Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, 767 Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subra-768 mani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo 769 Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harliman, 770 Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian 771 Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, So-772 maieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, 774 Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, An-775 drea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chh-776 ablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, 777 Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Tee-778 han, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Fevry, Trishala Neeraj, Urmish Thakker, Vikas Raunak, Xiangru Tang, Zheng-Xin Yong, 780 Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, 781 Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared 782 Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myr-783 iam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre 784 Cornette, Pierre François Lavallée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwa, Amanpreet Singh, Anas-785 tasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aurélie Névéol, Charles Lovering, 786 Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli 787 Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, 788 Jessica Zosa Forde, Jordan Clive, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, 789 Miruna Clinciu, Najoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, 790 Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdeněk Kasner, Alice 793 Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony 794 Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh Haji-Hosseini, Bahareh Behroozi, Benjamin Ajibade, Bharat Saxena, Carlos Muñoz Ferrandis, Daniel McDuff, Danish Contractor, David Lansky, Davis David, Douwe Kiela, Duong A. Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatima Mirza, Frankline Ononiwu, Habib Rezanejad, Hessie Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jesse 798 Passmore, Josh Seltzer, Julio Bonis Sanz, Livia Dutra, Mairon Samagaio, Maraim Elbadri, Mar-799 got Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, Muhammed 800 Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nour Fahmy, Olanre-801 waju Samuel, Ran An, Rasmus Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas 802 Wang, Sourav Roy, Sylvain Viguier, Thanh Le, Tobi Oyebade, Trieu Le, Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, An-804 tonio Miranda-Escalada, Ayush Singh, Benjamin Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, 805 Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel León Periñán, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrimann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Imane Bello, Ishani Dash, Jihyun Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthik Rangasai Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, Maria A Castillo, Marianna Nezhurina, Mario Sänger, Matthias Samwald, Michael Cullan, Michael Weinberg, Michiel De



Figure 10: Four causal graphs illustrating well-known exposure-outcome effects in the medical literature. This figure is from Long et al. (2022).

Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patrick Haller, Ramya Chandrasekhar, Renata Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aroonsiri, Srishti Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash Venkatraman, Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. Bloom: A 176b-parameter open-access multilingual language model, 2023. URL https://arxiv.org/abs/2211.05100.

Jie Yang, Soyeon Caren Han, and Josiah Poon. A survey on extraction of causal relations from natural language text. *Knowledge and Information Systems*, 64(5):1161–1186, 2022.

Alessio Zanga, Elif Ozkirimli, and Fabio Stella. A survey on causal discovery: Theory and practice. *International Journal of Approximate Reasoning*, 151:101–129, 2022.

- Matej Zečević, Moritz Willig, Devendra Singh Dhami, and Kristian Kersting. Causal parrots: Large language models may talk causality but are not causal. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id= tv46tCzs83.

A APPENDIX

A.1 GROUND-TRUTH CAUSAL GRAPHS

Figure 10, 11, 12 demonstrate ground-truth causal graphs for the causal discovery task.

A.2 CAUSAL DIRECTION IDENTIFICATION TASK

ConceptNet is a knowledge graph that connects natural language concepts via labeled edges. It
includes the "[A, /r/Causes, B]" relation, indicating that event A causes event B. Each relation in
ConceptNet also has a weight attribute, reflecting the confidence level of the relation; a higher
weight suggests broader agreement across sources. From ConceptNet, we selected the top 1,900
causal relations by weight and generated an equal number of reverse-causal relations by swapping
the cause and effect. This process yielded a total of 3,800 causal and reverse-causal relations.

CauseNet is a large-scale knowledge base containing claimed causal relations between concepts.
 We extract 814 high-confidence causal relations from CauseNet, each supported by at least 100 web sources and 10 extraction patterns. By swapping the cause and effect, we generate an equivalent number of reverse-causal relations. We then create a dataset containing 1,628 causal and reverse-causal relations.



Figure 11: The causal graph between key atmospheric variables and sea ice over the Arctic based on literature review. This figure is from Huang et al. (2021).



Figure 12: The causal graph for evaluating car insurance risks. This figure is sourced from Scutari (2010).
917

#### 918 A.3 IN-CONTEXT LEARNING AND PROMPT 919

920		
921	Demonstra	tions for in-context learning
922	User:	This task is to determine the cause-and-effect relationship between two events based on commonsense knowledge. We are interested in the causal relationship between 'it is raining' and 'carrying an umbrella'. Does 'it is raining' cause 'carrying an umbrella'?
923		Let's provide a step-by-step explanation, then give your final answer using yes or no. Step-by-Step Explanation:
924	Assistant:	<ol> <li>Observation: The presence of rain often results in individuals using umbrellas.</li> <li>Rationale: People use umbrellas to protect themselves from getting wet during rain.</li> </ol>
925		3. Conclusion: Rain acts as a cause for the use of an umbrella. The final answer is Yes. Final Answer: Yes
926	Llear	This task is to determine the cause-and-effect relationship between two events based on commonsense knowledge. We are interested in the causal relationship between 'it is raining' and 'carrying an umbrella'.
927	0501.	Does 'carrying an umbrella' causes 'it is raining'? Let's provide a step-by-step explanation, then give your final answer using yes or no.
928		Step-by-Step Explanation: 1. Observation: The act of opening an umbrella does not influence weather patterns or conditions.
929	Assistant:	<ol> <li>Rationale: Rain is a meteorological event determined by atmospheric conditions. The use of an umbrella does not have the capability to influence weather conditions.</li> <li>Conclusion: There is no causal relationship where the umbrella causes rain. The answer is No.</li> </ol>
930		Final Answer: No This task is to determine the cause-and-effect relationship between two events based on commonsense knowledge.
931	User:	We are interested in the causal relationship between 'Eating food' and 'feeling Hunger'. Does 'Eating food' causes 'feeling Hunger'?
932		Let's provide a step-by-step explanation, then give your final answer using yes or no. Step-by-Step Explanation:
933	Assistant:	<ol> <li>Observation: Eating involves consuming food, which typically reduces the feeling of hunger.</li> <li>Rationale: The act of eating is intended to satiate hunger, not induce it.</li> </ol>
934		3. Conclusion: Eating does not cause hunger. Instead, it alleviates it. The answer is No. Final Answer: No
935	User	This task is to determine the cause-and-effect relationship between two events based on commonsense knowledge. We are interested in the causal relationship between 'Eating food' and 'feeling Hunger'.
936	0.301.	Does 'Feeling hunger' causes 'Eating food'? Let's provide a step-by-step explanation, then give your final answer using yes or no.
937		Step-by-Step Explanation: 1. Observation: The sensation of hunger typically motivates individuals to eat.
938	Assistant:	<ol> <li>Rationale: Hunger is a biological signal that the body needs nutrients. This sensation prompts the behavior of eating to replenish energy and nutrients.</li> <li>Conclusion: It is reasonable to conclude that hunger causes eating. The answer is Yes.</li> </ol>
939	Prompt	Final Answer: Yes
940		This task is to determine the cause-and-effect relationship between two events based on commonsense knowledge. We are interested in the causal relationship between {cause}' and '{effect}'.
941	User:	Does '{cause}' cause '{effect}'? Let's provide a step-by-step explanation, then give your final answer using yes or no.

#### Table 2: Demonstrations for in-context learning and the prompt for new input.

943 944 945

946

947

948

942

For the causal direction identification task and the causal discovery task, we employ similar incontext learning demonstrations and prompts, detailed in Table 2. When presented with a pair of nodes (A, B), we generate two questions: "Does A cause B?" and "Does B cause A?".

In the causal direction identification task, the ground-truth instances are formatted as  $(A \rightarrow B, true)$ 949 and  $(A \leftarrow B, false)$ . These yes-no questions are directly transformed into such instances, aligning 950 perfectly with the binary nature of the task. In the causal discovery task, the ground-truth instances 951 are structured as (A, B, l), where the label l can take one of four possible values:  $\leftarrow, \rightarrow, \times, \leftrightarrow$ . 952 Here,  $\times$  denotes no causal relation, and  $\leftrightarrow$  indicates a bi-directional causal relation. We include 953 bi-directional causal relation because it exists in some ground-truth causal graphs such as Arctic Sea 954 Ice. The conversion of yes-no responses to these four-way labels is handled as follows. If only one 955 of the questions receives a 'yes' answer, it translates directly to the corresponding causal direction 956  $(i.e., \leftarrow \text{ or } \rightarrow)$ . If both questions are answered with 'no', this indicates no causal relation  $(i.e., \times)$ . If both questions receive a 'yes' response, this suggests a bi-directional relation (*i.e.*,  $\leftrightarrow$ ). 957

To determine the most confident answer, each LLM should generate ten distinct responses Chen & Mueller (2023); Geng et al. (2024). We then extract 'yes' or 'no' from each output. If the count of 'yes' responses is greater than or equal to the count of 'no' responses, the final answer is 'yes'. If 'no' responses predominate, the final answer is 'no'. This methodology ensures a robust approach to determining causal relationships in both tasks.

The decoding hyperparameters are configured as follows: the top-p sampling parameter is set to
0.9, the repetition penalty is 1.25, the temperature is 0.8, and the maximum number of new tokens
generated does not exceed the maximum input length. We employ the Hugging Face library to load
LLMs and generate responses Wolf et al. (2020). All experiments were conducted on NVIDIA A100
GPUs.

968 969

970 A.4 QUERY FOR SEARCH ENGINE

971

The queries for searching can be found in Table 3, 4.

natch for "event A causes event B"
es = [f"{cause} causes {effect}", f"{effect} is caused by {cause}", f"{cause} leads to {effect}",
e) results in {effect}", f'{cause} triggers {effect}", f'{effect} is triggered by {cause}",
a) affacts (affact)" f"(affact) is affacted by (cause)" f"(cause) impacts (affact)"
e is impacted by leffect?" f'(cause) is responsible for leffect?"
e} is the reason for {effect}". f"The effect of {cause} is {effect}".
esult of {cause} is {effect}". f"The consequence of {cause} is {effect}".
t} is a consequence of {cause}", f'{effect} is a result of {cause}", f'{effect} is an effect of {cause}"]
· · · · · · · · · · · · · · · · · · ·
match_phrase query for each template
list = []
ise in templates:
$phrase = \{$
_phrase": {
: phrase,
int(ien(prrase.split())*0.25),
list annend(match_nhrase)
.{
r <sup>*</sup> : should_list.
um_should_match": 1
YNTHETIC CAUSAL RELATIONS
demonstrates templates for creating mentions of synthetic causal relations and anti-causal s.
RAINING DETAILS
tuning OLMo-7b-Instruct using LoRA on synthetic datasets, utilizing the official code from $Ao$ repository <sup>2</sup> . The model was trained on two NVIDIA A100 GPUs with a batch size of 2 J, and a total batch size of 128. We set the LoRA rank and alpha to 256, with a dropout rate The learning rate was configured to 1e-4, employing a linear scheduler for rate adjustments. ning was conducted over one epoch.
IUMAN EVALUATION FOR CAUSAL RELATION WITH CONTEXTS
mpt of generation contexts of causal relations is shown in Table 6. In this task, we require an-
to evaluate causal relations with different contexts. Below we show detailed task instruction ators.

<sup>1025 &</sup>lt;sup>2</sup>We employed the official OLMo code available at https://github.com/allenai/ open-instruct.

1026	Ordered phrase search for "event A" $\Rightarrow$ "causes" $\Rightarrow$ "event B"
1027	causal_mentions = ["causes", "leads to", "results in", "triggers", "induces", "influences", "affects", "impacts",
1028	"is responsible for", "is the reason for", "cause", "lead to", "result in", "trigger", "induce",
1029	"influence", "affect", "impact", "are responsible for", "are the reason for"]
1030	# create cause clause in span term format
1031	cause_clauses = []
1032	for item in cause.split():
1033	cause_clauses.append({''span_term'': {''text'': item}})
1034	# create effect clause in span term format
1035	effect_clauses = []
1036	for item in effect.split():
1037	effect_clauses.append({"span_term": {"text": item}})
1038	# create causal relation clause in span term format
1039	all_relation_clauses = []
1040	for rel in causal_mentions:
1041	relation_clauses = []
1042	Ior term in rel.split():
1043	all_relation_clauses.append(relation_clauses)
1044	
1045	# for each causal relation clause, create a query
1046	for relation_clauses in all_relation_clauses:
1047	query = { "snan near": ∫
1048	"clauses":
1049	{
1050	"span_near": {
1051	"clauses": cause_clauses,
1052	stop : 0, "in order": True
1052	}
1053	},
1055	
1055	"span_near": {
1050	"slop": 0.
1057	"in_order": True
1050	}
1059	}, {
1060	{ "cman_near": ∫
1061	"clauses": effect_clauses.
1062	"slop": 0,
1063	"in_order": True
1064	
1065	} 1
1066	"slop": 32, # window size
1067	"in_order": True
1068	}
1069	}
1070	Table 4: "event $\Lambda$ " $\rightarrow$ "causes" $\rightarrow$ "event <b>B</b> " query for WIMBD
1071	Table 4. EVENT $A \rightarrow$ causes $\rightarrow$ EVENT D query for WIWIDD.
1072	
1073	
1074	Annotation Steps. Below is suggested annotation steps to annotators.
1075	
1076	1. 1. Read the Scenario Carefully: Each scenario provides a specific context. Understand the
1077	details and implications of the scenario.
1078	
1079	2. 2. Review the Question: Each question asks you to assess the likelihood of a causal relation
	occurring, given the provided scenario.

Correct causal relations	Reverse causal relations	Negation of causal relations $f_{1}^{(1)} = f_{1}^{(2)} f_{2}^{(2)} + f$
f'effect is caused by cause.",	f'cause is caused by effect.".	f'effect is not caused by cause.".
f'cause leads to effect.",	f'effect leads to cause.",	f"cause does not lead to effect.",
f''cause results in effect.", f''cause triggers effect."	f'effect results in cause.", f'effect triggers cause "	f''cause does not result in effect.", f''cause does not trigger effect."
f'effect is triggered by cause.",	f'cause is triggered by effect.",	f'effect is not triggered by cause.",
f'cause induces effect.", f'cause influences effect."	f'effect induces cause.", f'effect influences cause "	f''cause does not induce effect.",
f'effect is influenced by cause.",	f'cause is influenced by effect.",	f'effect is not influenced by cause.",
f'cause affects effect.", f'affect is affected by cause."	f'effect affects cause.", f'agusa is affected by affect "	f''cause does not affect effect.", f''effect is not affected by cause "
f'cause impacts effect.",	f'effect impacts cause.",	f'cause does not impact effect.",
f'cause is impacted by effect.",	f'effect is impacted by cause.",	f''cause is not impacted by effect.",
f'cause is the reason for effect.",	f'effect is the reason for cause.",	f'cause is not the reason for effect.",
f"The effect of cause is effect.",	f"The effect of effect is cause.",	f"The effect of cause is not effect.",
f"The consequence of cause is effect.",	f"The consequence of effect is cause.",	f"The consequence of cause is not effect.",
f"effect is a consequence of cause.",	f'cause is a consequence of effect.",	f"effect is not a consequence of cause.",
f"effect is an effect of cause.", ]	f'cause is an effect of effect.", ]	f'effect is not an effect of cause.",]
Table 5: Templates for creatin	g mentions of imaginary causal	relations and anti-causal relations.
1		
Prompt for generating contexts of o	causal relations	as where 'cause' does not cause 'effect'
Each scenario should be distinctly	and clearly described, categorized up	nder the respective headings.
Response Format:		lider die respective neueniger
Scenarios where 'cause' causes 'ef	fect':	
Heading:		
Reason:		
Description: Reason: 		
Table 6: Pro	ompt for generating contexts of	causal relations.
3. 3. Select the Appropria probability range that I	ate Answer: Based on your und best represents the likelihood of	erstanding of the scenario, select the the stated causal relation occurring.
for each question, we have belo	ow options	
• 100%: The causal rela	tion definitely occurs.	
• 81-99%: The causal re	lation almost certainly occurs.	
• 51-80%: The causal re	lation is likely to occur.	
• 50%: The causal relati	on has 50	
• 20-49%: The causal re	lation somewhat likely to occur	
• 1-19%: The causal rela	ation rarely occurs.	
• 0%: The causal relatio	n never occurs.	
• The scenario does not	make sense. If the scenario cor	tradicts common sense or could not
occur in the real world	or it is not a scenario at all, ple	ase select this option.
Annotation Examples In Ta		
	ble 7, we show some annotation	examples to help annotators have a
better understanding of this task	ble 7, we show some annotation	n examples to help annotators have a
better understanding of this task	ble 7, we show some annotation c.	n examples to help annotators have a

**Acceptance Policy.** We will only reject a job if there is clear evidence of malicious behavior, such as random clicking, which suggests non-compliance with task guidelines.

1134	Annotation examples
1135	Question: to what extent do you think 'soaking in a hotspring' causes 'relaxation'?
1137	81-99%
1138	Reason: The warm water of a hot spring helps to raise the body's temperature, which can relax muscle tension and soothe aches and pains in the joints and muscles
1139	This physical relaxation naturally leads to mental relaxation.
1140	###relation_91###
1141	The water temperature in the hotspring is excessively hot, making the individual feel uncomfortable.
1142 1143 1144	Answer the following question ONLY based on information described in above scenario and your common sense. Question: under above scenario, to what extent do you think 'soaking in a hotspring' causes 'relaxation'?
1145	Reason: Uncomfortably high temperatures can cause overheating, dizziness, or discomfort, preventing relaxation.
1146	Scenario — Entertaining Friends:
1147	During a casual get-together with friends, you crack jokes and everyone bursts into laughter.
1149 1150	Answer the following question ONLY based on information described in above scenario and your common sense. Question: under above scenario, to what extent do you think 'making people laugh' causes 'you have fun too'?
1151	81-99% Reason: The shared joy and camaraderie among friends create a fun and enjoyable atmosphere.
1152	
1153	Table 7: Examples of causal relation evaluation under different contexts.
1154	
1155	<b>Defense Defense</b> Operations of the first state of the life of the second state of the
1157	not publish your name email address or any other personal information. If you have concerns about
1158	how we handle your personal data, please contact the project manager.
1159	
1160	
1161	B MORE EXPERIMENT RESULTS
1162	
1163	B.1 EVALUATING BOTH OPEN- AND CLOSED-SOURCE LLMS ON CAUSAL DISCOVERY TASKS.
1164	Causal questions indicate both causal direction identification task and causal discovery task.
1165	Kıcıman et al. (2023); Zečević et al. (2023); Feng et al. (2024); Jiralerspong et al. (2024) have re-
1167	ported that closed-source LLMs (e.g., GPT-3.5-turbo, GPT-4) achieve state-of-the-art performance
1168	in causal direction identification task and causal discovery tasks. However, their analyses predomi-
1169	nantly focus on specific closed-source models and offer a limited examination of open-source LLMs.
1170	In this section, we employ closed-source and open-source LLMs to conduct causal relation identifi- cation and causal discovery tasks. We aim to compare and analyze the performance disparities when
1171	utilizing different models. Table 8, 9, 10, 11, 12 and 13 show the results of causal discovery ex-
1172	periments on the Arctic Sea Ice, Insurance, Alcohol, Cancer, Diabetes, and Obesity causal graphs.
1173	Table 14 and 15 show the results of causal direction identification tasks on the ConceptNet and
1174	CauseNet datasets.
1175	We employ the Normalized Hamming Distance (NHD) as one metric for full causal discovery. A
1176	notable issue with NHD is that due to the typically sparse nature of causal graphs, models that predict
1177	no edges can still achieve a low NHD. This setup inadvertently penalizes models that predict a larger
1178	number of edges, even true edges may be predicted. To address this, following the methodologies
1179	outlined by Kiciman et al. (2023) and Jiralerspong et al. (2024), we calculate the ratio between the
1180	being incorrect. The lower the ratio the better the model performs compared to the worst baseline
1101	that outputs the same number of edges. Therefore, we report NHD ratio ( <i>i.e.</i> , NHD / baseline
1102	NHD), along with the number of predicted edges, to provide a more comprehensive evaluation of
1103	model performance in the full causal discovery task.
1185	Due to the transparency of OLMo-7h-Instruct and the robust canabilities of its search tool. OLMo-
1186	7b-Instruct serves as our primary analysis model. Therefore, we explored various numbers of in-

1187 context learning examples to identify the optimal example number. In seven out of eight datasets, OLMo-7b-Instruct with three demonstration examples achieves the highest F1, compared to other

numbers of demonstration examples tested. Therefore, to ensure a fair comparison, other LLMs also utilized three demonstration examples for in-context learning.

Considering all LLMs, GPT-40 outperforms others in six of the eight datasets evaluated, specif-1191 ically Arctic Sea Ice, Insurance, Alcohol, Obesity, ConceptNet, and CauseNet. In the remaining 1192 two datasets, Cancer and Diabetes, GPT-40 ranks as the second-best model, with only a slight per-1193 formance differential from the top model. These experiment results show that GPT-40 is the most 1194 effective model for causal discovery and causal direction identification tasks in both closed- and 1195 open-source models. Among open-source models exclusively, Llama3-8b-Instruct excels, achieving 1196 the highest F1 scores in six datasets: Insurance, Alcohol, Cancer, Diabetes, Obesity, and CauseNet. 1197 Meanwhile, Llama2-7b-chat achieves the highest F1 in two datasets, Arctic Sea Ice and Obesity. 1198 In the ConceptNet dataset, OLMo-7b-Instruct, configured with three in-context learning examples, records the best F1 score. 1199

		<b>Precision</b> <sup>↑</sup>	<b>Recall</b> ↑	F1↑	<b>Accuracy</b> ↑	Predict edges	NHD↓	Baseline NHD	Ratio (NHD/Baseline NHD)↓
)1	OLMo-7b-Instruct (0 ICL)	0.4259	0.5	0.46	0.625	54	0.375	0.6944	0.54
2	OLMo-7b-Instruct (1 ICL)	0.3928	0.4782	0.4314	0.5972	56	0.4027	0.7083	0.5686
) <	OLMo-7b-Instruct (2 ICL)	0.4615	0.1304	0.2034	0.6736	13	0.3263	0.4097	0.7966
3	OLMo-7b-Instruct (3 ICL)	0.5555	0.1087	0.1818	0.6875	9	0.3125	0.3819	0.8181
	OLMo-7b-Instruct (4 ICL)	0.5417	0.2826	0.3714	0.6944	24	0.3055	0.4861	0.6285
4	BLOOM-7b1 (3 ICL)	0.3934	0.5217	0.4485	0.5902	61	0.4097	0.7430	0.5514
_	Llama2-7b-chat (3 ICL)	0.4444	0.5217	0.48	0.6388	54	0.3611	0.6944	0.52
5	Llama3-8b-Instruct (3 ICL)	1.0	0.1956	0.3272	0.7430	9	0.2569	0.3819	0.6727
3	GPT-3.5-turbo (3 ICL)	0.7647	0.2826	0.4126	0.7431	17	0.2569	0.4375	0.5873
0	GPT-40 (3 ICL)	0.5178	0.6304	0.5686	0.6944	56	0.3055	0.7083	0.4313

1208Table 8: Causal discovery results for the Arctic Sea Ice causal graph, with 12 nodes and 46 edges.1209GPT-40 surpasses all competing models, achieving an F1 score of 0.5686 and an NHD ratio of12100.4313. The second-best performing model is an open-source LLM, Llama2-7b-chat. (# ICL) indicates the number of demonstration examples for in-context learning.

	<b>Precision</b> ↑	<b>Recall</b> ↑	F1↑	<b>Accuracy</b> ↑	Predict edges	NHD↓	Baseline NHD	Ratio (NHD/Baseline NHD)↓
OLMo-7b-Instruct (0 ICL)	0.0873	0.7692	0.1568	0.4101	458	0.5898	0.6995	0.8431
OLMo-7b-Instruct (1 ICL)	0.0963	0.9038	0.1740	0.3882	488	0.6117	0.7407	0.8259
OLMo-7b-Instruct (2 ICL)	0.0901	0.5961	0.1565	0.5418	344	0.4581	0.5432	0.8434
OLMo-7b-Instruct (3 ICL)	0.1254	0.6731	0.2114	0.6419	279	0.3580	0.4540	0.7885
OLMo-7b-Instruct (4 ICL)	0.1093	0.7884	0.1920	0.5267	375	0.4732	0.5857	0.8079
BLOOM-7b1 (3 ICL)	0.0710	0.7115	0.1291	0.3155	521	0.6844	0.7860	0.8708
Llama2-7b-chat (3 ICL)	0.1245	0.7115	0.2120	0.6227	297	0.3772	0.4787	0.7879
Llama3-8b-Instruct (3 ICL)	0.2656	0.3269	0.2931	0.8875	64	0.1124	0.1591	0.7069
GPT-3.5-turbo (3 ICL)	0.1575	0.5	0.2396	0.7736	165	0.2263	0.2976	0.7603
GPT-40 (3 ICL)	0.2287	0.6730	0.3414	0.8148	153	0.1851	0.2812	0.6585

Table 9: Causal discovery results for the Insurance causal graph, with 27 nodes and 52 edges. GPT-40 surpasses all competing models, achieving an F1 score of 0.3414 and an NHD ratio of 0.6585. The second-best performing model is an open-source LLM, Llama3-8b-Instruct.

25		<b>Precision</b> ↑	<b>Recall</b> ↑	F1↑	Accuracy↑	Predict edges	NHD↓	Baseline NHD	Ratio (NHD/Baseline NHD)↓
-	OLMo-7b-Instruct (0 ICL)	0.5	1.0	0.6667	0.6667	6	0.3333	1.0	0.3333
26	OLMo-7b-Instruct (1 ICL)	0.6	1.0	0.75	0.7778	5	0.2222	0.8889	0.25
07	OLMo-7b-Instruct (2 ICL)	0.5	1.0	0.6667	0.6667	6	0.3333	1.0	0.3333
27	OLMo-7b-Instruct (3 ICL)	0.6	1.0	0.75	0.7778	5	0.2222	0.8889	0.25
28	OLMo-7b-Instruct (4 ICL)	0.6	1.0	0.75	0.7778	5	0.2222	0.8889	0.25
20	BLOOM-7b1 (3 ICL)	0.5	1.0	0.6667	0.6667	6	0.3333	1.0	0.3333
29	Llama2-7b-chat (3 ICL)	0.75	1.0	0.8571	0.8889	4	0.1111	0.7778	0.1429
	Llama3-8b-Instruct (3 ICL)	1.0	1.0	1.0	1.0	3	0	0.6667	0
<b>60</b>	GPT-3.5-turbo (3 ICL)	1.0	1.0	1.0	1.0	3	0	0.6667	0
<b>1</b> -1	GPT-40 (3 ICL)	1.0	1.0	1.0	1.0	3	0	0.6667	0

Table 10: Causal discovery results for the Alcohol causal graph, with 3 nodes and 3 edges. Llama38b-Instruct, GPT-3.5-turbo, and GPT-4 accurately predict the ground-truth causal graph. The second-best performing model is Llama2-7b-chat.

1235 1236

1200

1212

1224

1237

B.2 DO PRE-TRAINING CORPORA CONTAIN MORE CORRECT CAUSAL RELATIONS?

Given the effective performance of LLMs on causal discovery tasks, a pertinent research question arises: Why can LLMs perform so well? We posit that a significant factor is the nature of the pretraining data, which contains more correct causal relations than incorrect ones, leading LLMs to primarily memorize correct causal relations.

	<b>Precision</b> ↑	<b>Recall</b> ↑	F1↑	<b>Accuracy</b> ↑	Predict edges	NHD↓	Baseline NHD	Ratio (NHD/Baseline NHD)↓
OLMo-7b-Instruct (0 I	CL) 0.4166	1.0	0.5882	0.5625	12	0.4375	1.0	0.4375
OLMo-7b-Instruct (1 I	CL) 0.4	0.8	0.5333	0.5625	10	0.4375	0.9375	0.4667
OLMo-7b-Instruct (2 I	CL) 0.5	0.8	0.6153	0.6875	8	0.3125	0.8125	0.3846
OLMo-7b-Instruct (3 I	CL) 0.5714	0.8	0.6667	0.75	7	0.3125	0.9375	0.3333
OLMo-7b-Instruct (4 I	CL) 0.5	1.0	0.6667	0.6875	10	0.3125	0.9375	0.3333
BLOOM-7b1 (3 ICL)	0.4	0.4	0.4	0.625	5	0.375	0.625	0.6
Llama2-7b-chat (3 ICL	) 0.4166	1.0	0.5882	0.5625	12	0.4375	1.0	0.4375
Llama3-8b-Instruct (3	CL) 1.0	0.8	0.8889	0.9375	4	0.0625	0.5625	0.1111
GPT-3.5-turbo (3 ICL)	1.0	0.8	0.8889	0.9375	4	0.0625	0.5625	0.1111
GPT-4o (3 ICL)	0.8	0.8	0.8	0.875	5	0.125	0.625	0.2

1249 Table 11: Causal discovery results for the Cancer causal graph, with 4 nodes and 5 edges. Llama3-1250 8b-Instruct and GPT-3.5-turbo surpass all other models. The second-best performing model is GPT-1251 40. 1252

-	<b>Precision</b> ↑	<b>Recall</b> ↑	F1↑	Accuracy↑	Predict edges	NHD↓	Baseline NHD	Ratio (NHD/Baseline NHD)↓
OLMo-7b-Instruct (0 ICL)	0.4166	1.0	0.5882	0.5625	12	0.4375	1.0625	0.4117
OLMo-7b-Instruct (1 ICL)	0.4166	1.0	0.5882	0.5625	12	0.4375	1.0625	0.4117
OLMo-7b-Instruct (2 ICL)	0.4166	1.0	0.5882	0.5625	12	0.4375	1.0625	0.4117
OLMo-7b-Instruct (3 ICL)	0.5	1.0	0.6666	0.6875	10	0.3125	0.9375	0.3333
OLMo-7b-Instruct (4 ICL)	0.4545	1.0	0.625	0.625	11	0.375	1.0	0.375
BLOOM-7b1 (3 ICL)	0.4285	0.6	0.5	0.625	7	0.375	0.75	0.5
Llama2-7b-chat (3 ICL)	0.5556	1.0	0.7142	0.75	9	0.25	0.875	0.2857
Llama3-8b-Instruct (3 ICL)	1.0	0.8	0.8889	0.9375	4	0.0625	0.5625	0.1111
GPT-3.5-turbo (3 ICL)	1.0	1.0	1.0	1.0	5	0	0.625	0
GPT-40 (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909

1260 Table 12: Causal discovery results for the Diabetes causal graph, with 4 nodes and 5 edges. GPT-1261 3.5-turbo accurately predict the ground-truth causal graph. The second-best performing model is 1262 GPT-40. 1263

1264

1265 Research Question 4. Do pre-training corpora contain more correct causal relations than incorrect 1266 ones? 1267

1268 Humans fundamentally rely on causal relations to understand and generate text. Therefore, it is 1269 reasonable that pre-training corpora, which are collected from human-generated texts, are likely to inherently contain a higher proportion of correct causal relations. 1270

1271 Observation We count the total occurrence of correct and incorrect causal relations in Dolma and 1272 ROOTS corpora. The results are shown in Table 16. We use exact matching to count correct and 1273 incorrect causal relations. We observe that the occurrence of causal relations is, on average, 12 times 1274 higher than that of incorrect causal relations in Dolma and ROOTS corpora. From our observation, 1275 most incorrect causal relations do not exist in an affirmation context. They are usually in a question or negation context. For example, "Which option is correct? A. smoking causes cancer B. cancer 1276 causes smoking" or "Which means that either smoking causes cancer or cancer causes smoking." 1277

1278 Discussion In conclusion, these experimental results show that correct causal relations are more 1279 frequently represented than incorrect ones in pre-training corpora. This also explain why LLMs can 1280 identify many causal relations in causal discovery tasks.

1281

1285

1282 B.3 INFLUENCE OF MODEL SIZE ON LLMS' PERFORMANCE IN CAUSAL DISCOVERY TASKS 1283

1284 **Research Question 5.** Do larger models perform better on causal discovery tasks?

1286 We assume that within the same architectural framework, increasing the model size (i.e., the number of parameters) leads to improved performance on causal discovery tasks. The rationale is that larger 1287 models can memorize more information from the pre-training data than their smaller models. 1288

1289 Observation We select models from the Llama2 and Llama3 series, each varying in size. These 1290 models are evaluated on causal discovery and causal direction identification tasks, with results doc-1291 umented in Table 17 and 18. The findings indicate that for both the Llama2 and Llama3 models, 1292 there is a positive correlation between the number of parameters and performance. However, dis-1293 crepancies arise when comparing across architectures. For example, a small Llama3 model (e.g., Llama3-8b-Instruct) can outperform a significantly larger Llama3 model (e.g., Llama2-70b-chat). 1294 Notably, across most datasets, Llama3-70b-Instruct either matches or surpasses the performance of 1295 the currently leading closed-source LLM, GPT-40.

1296		<b>Precision</b> ↑	Recall↑	F1↑	Accuracy <sup>↑</sup>	Predict edges (46)	NHD↓	Baseline NHD	Ratio (NHD/Baseline NHD)↓
1007	OLMo-7b-Instruct (0 ICL)	0.5714	0.8	0.6666	0.75	7	0.3125	0.9375	0.3333
1297	OLMo-7b-Instruct (1 ICL)	0.5	1.0	0.6666	0.6875	10	0.3125	0.9375	0.3333
1000	OLMo-7b-Instruct (2 ICL)	0.5555	1.0	0.7142	0.75	9	0.25	0.875	0.2857
1290	OLMo-7b-Instruct (3 ICL)	0.8	0.8	0.8	0.875	5	0.125	0.625	0.2
1299	OLMo-7b-Instruct (4 ICL)	0.5555	1.0	0.7142	0.75	9	0.25	0.875	0.2857
1200	BLOOM-7b1 (3 ICL)	0.4444	0.8	0.5714	0.625	9	0.375	0.875	0.4285
1300	Llama2-7b-chat (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
	Llama3-8b-Instruct (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
1301	GPT-3.5-turbo (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
1000	GPT-40 (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909

Table 13: Causal discovery results for the Obesity causal graph, with 4 nodes and 5 edges. Llama2 7b-chat, Llama3-8b-Instruct, GPT-3.5-turbo and GPT-40 outperform all other models. The second best performing method is OLMo-7b-Instruct (3 ICL).

1307		<b>Precision</b> ↑	<b>Recall</b> ↑	F1↑	Accuracy↑
1308	OLMo-7b-Instruct (0 ICL)	0.5482	0.8831	0.6765	0.5778
1309	OLMo-7b-Instruct (1 ICL)	0.5491	0.8184	0.6573	0.5734
1310	OLMo-7b-Instruct (2 ICL)	0.5771	0.7825	0.6643	0.6047
1311	OLMo-7b-Instruct (3 ICL)	0.6612	0.8427	0.7410	0.7053
1212	OLMo-7b-Instruct (4 ICL)	0.5294	0.8721	0.6589	0.5486
1012	BLOOM-7b1 (3 ICL)	0.5027	0.7248	0.5937	0.5041
1313	Llama2-7b-chat (3 ICL)	0.6197	0.7774	0.6897	0.6503
1314	Llama3-8b-Instruct (3 ICL)	0.7659	0.6575	0.7076	0.7282
1315	GPT-3.5-turbo (3 ICL)	0.6732	0.7308	0.7008	0.6891
1316	GPT-40 (3 ICL)	0.8141	0.8342	0.8240	0.8224
1317					

1318Table 14: Causal direction identification results on the ConceptNet dataset, with 1900 causal rela-1319tions and 1900 reverse causal relations. GPT-40 outperforms all competing methods, achieving an1320F1 score of 0.8240. The second-best performing method is OLMo-7b-Instruct (3 ICL), with an F11321score of 0.7410.

1322
 1323
 1324
 1324
 1325
 *Discussion* The experiment results lead to a critical consideration of the 'bigger is better' paradigm in LLM research. Future research should thus not only focus on scaling up the size but also on refining the architecture and learning algorithms to better leverage increased model capacity.

1350					
1951		<b>Precision</b> ↑	Recall↑	F1↑	Accuracy↑
1551	OLMo-7b-Instruct (0 ICL)	0.5461	0.9657	0.6977	0.5815
1352	OLMo-7b-Instruct (1 ICL)	0.5359	0.9606	0.6881	0.5644
1353	OLMo-7b-Instruct (2 ICL)	0.5610	0.9091	0.6938	0.5988
1354	OLMo-7b-Instruct (3 ICL)	0.6568	0.8771	0.7511	0.7094
1355	OLMo-7b-Instruct (4 ICL)	0.5860	0.9410	0.7223	0.6382
1356	BLOOM-7b1 (3 ICL)	0.5067	0.6928	0.5853	0.5092
1357	Llama2-7b-chat (3 ICL)	0.7030	0.8931	0.7867	0.7582
1358	Llama3-8b-Instruct (3 ICL)	0.8838	0.8296	0.8558	0.8602
1359	GPT-3.5-turbo (3 ICL)	0.8990	0.8857	0.8923	0.8931
1360	GPT-4o (3 ICL)	0.8596	0.9557	0.9051	0.8998

Table 15: Causal direction identification results on the CauseNet dataset, with 814 causal relations and 814 reverse causal relations. GPT-40 outperforms all competing methods, achieving an F1 score of 0.9051. The second-best performing method is GPT-3.5-turbo, with an F1 score of 0.8923.

1365				
1000			Correct Causal Relations	Incorrect Causal Relations
1366	Coursel Discoursery (all detects)	Dolma	28812	1127
1367	Causal Discovery (all datasets)	ROOTS	814	118
1368	Caused Direction Identification (ConcentNat)	Dolma	41407	3410
1360	Causar Direction Identification (Conceptiver)	ROOTS	1176	131
1000	Caused Direction Identification (Cause Nat)	Dolma	949427	107070
1370	Causal Direction Identification (CauseNet)	ROOTS	24591	4236
1371				

Table 16: Occurrences of correct and incorrect causal relations in the Dolma and ROOTS corpora.

Arctic Sea Ice								
	<b>Precision</b> ↑	<b>Recall</b> ↑	F1↑	<b>Accuracy</b> ↑	Predict edges	NHD↓	Baseline NHD	Ratio (NHD/Baseline NHD)↓
Llama2-7b-chat (3 ICL)	0.4444	0.5217	0.48	0.6388	54	0.3611	0.6944	0.52
Llama2-13b-chat (3 ICL)	0.4478	0.6522	0.5309	0.6319	67	0.3681	0.7847	0.4690
Llama2-70b-chat (3 ICL)	0.3606	0.9565	0.5238	0.4444	122	0.5556	1.0	0.5556
Llama3-8b-Instruct (3 ICL)	1.0	0.1956	0.3272	0.7430	9	0.2569	0.3819	0.6727
Llama3-70b-Instruct (3 ICL)	0.5689	0.7174	0.6346	0.7361	58	0.2639	0.7222	0.3653
GPT-3.5-turbo (3 ICL)	0.7647	0.2826	0.4126	0.7431	17	0.2569	0.4375	0.5873
GPT-40 (3 ICL)	0.5178	0.6304	0.5686	0.6944	56	0.3055	0.7083	0.4313
				Insura	ance			
Llama2-7b-chat (3 ICL)	0.1245	0.7115	0.2120	0.6227	297	0.3772	0.4787	0.7879
Llama2-13b-chat (3 ICL)	0.1338	0.7307	0.2262	0.6433	284	0.3566	0.4609	0.7738
Llama2-70b-chat (3 ICL)	0.1619	0.7692	0.2675	0.6995	247	0.3004	0.4102	0.7324
Llama3-8b-Instruct (3 ICL)	0.2656	0.3269	0.2931	0.8875	64	0.1124	0.1591	0.7069
Llama3-70b-Instruct (3 ICL)	0.2183	0.5961	0.3195	0.8189	142	0.1811	0.2661	0.6804
GPT-3.5-turbo (3 ICL)	0.1575	0.5	0.2396	0.7736	165	0.2263	0.2976	0.7603
GPT-40 (3 ICL)	0.2287	0.6730	0.3414	0.8148	153	0.1851	0.2812	0.6585
				Alco	hol			
Llama2-7b-chat (3 ICL)	0.75	1.0	0.8571	0.8889	4	0.1111	0.7778	0.1429
Llama2-13b-chat (3 ICL)	0.75	1.0	0.8571	0.8889	4	0.1111	0.7778	0.1429
Llama2-70b-chat (3 ICL)	0.75	1.0	0.8571	0.8889	4	0.1111	0 7778	0 1429
Lama3-8b-Instruct (3 ICL)	1.0	1.0	10	1.0	3	0	0.6667	0
Lama3-70b-Instruct (3 ICL)	1.0	1.0	1.0	1.0	3	Ő	0.6667	ő
GPT-3 5-turbo (3 ICL)	1.0	1.0	1.0	1.0	3	0	0.6667	ů 0
GPT-40 (3 ICL)	1.0	1.0	1.0	1.0	3	0	0.6667	ů 0
GI I 40 (5 ICE)	1.0	1.0	1.0	Can	cer	0	0.0007	
Llama2-7b-chat (3 ICL)	0.4166	1.0	0.5882	0.5625	12	0 4375	1.0	0.4375
Llama2-13b-chat (3 ICL)	0.5556	1.0	0.71/13	0.75	0	0.4575	0.875	0.2857
Llama2-70b-chat (3 ICL)	0.5556	1.0	0.7143	0.75	ó	0.25	0.875	0.2857
Llama2 8h Instruct (2 ICL)	1.0	1.0	0.7145	0.75	4	0.25	0.5625	0.2037
Llama 2 70b Instruct (2 ICL)	1.0	0.8	0.0007	0.9375	4	0.0625	0.5625	0.1111
CDT 2.5 turba (2 ICL)	1.0	0.0	0.0009	0.9373	4	0.0023	0.5625	0.1111
GF 1-5.5-IUFD0 (5 ICL)	1.0	0.8	0.0009	0.93/3	4	0.0025	0.3023	0.1111
Gr 1-40 (3 ICL)	0.8	0.8	0.8	0.875	3	0.123	0.025	0.2
Llomo2 7h shot (2 ICL)	0 5556	1.0	0.7142	Diab	etes	0.25	0.975	0 2857
Liama2-/D-Chat (SICL)	0.5550	1.0	0.7142	0.75	9	0.25	0.8/5	0.2857
Liama2-13b-chat (3 ICL)	0.625	1.0	0.7692	0.8125	8	0.18/5	0.8125	0.2307
Llama2-70b-chat (3 ICL)	0.625	1.0	0.7692	0.8125	8	0.1875	0.8125	0.2307
Llama3-8b-Instruct (3 ICL)	1.0	0.8	0.8889	0.9375	4	0.0625	0.5625	0.1111
Llama3-70b-Instruct (3 ICL)	1.0	1.0	1.0	1.0	5	0	0.625	0
GPI-3.5-turbo (3 ICL)	1.0	1.0	1.0	1.0	5	0	0.625	0
GPT-40 (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
	0.0000	1.0	0.000 -	Obes	sity	0.0(0-	0.6075	0.0000
Llama2-7b-chat (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
Llama2-13b-chat (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
Llama2-70b-chat (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
Llama3-8b-Instruct (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
Llama3-70b-Instruct (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
GPT-3.5-turbo (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909
GPT-40 (3 ICL)	0.8333	1.0	0.9091	0.9375	6	0.0625	0.6875	0.0909

<sup>1403</sup> Table 17: Performance on causal discovery task using Llama2 and Llama3 models of different sizes.

4.400					
1406					
1407					
1408					
1409					
1410					
1411					
1412					
1413					
1414					
1415					
1416					
1417					
1418					
1419					
1420					
1421					
1422		ConceptNe	t		
1423		<b>Precision</b> ↑	<b>Recall</b> ↑	F1↑	<b>Accuracy</b> ↑
1404	Llama2-7b-chat (3 ICL)	0.6197	0.7774	0.6897	0.6503
	$I_{1}$ I amo 2 12h shot (2 ICI)	0 (010	0 8605	0 7077	0.6647
1424	Liama2-150-chat (5 ICL)	0.6010	0.0005	0.7077	0.00.1
1425	Llama2-70b-chat (3 ICL)	0.6010	0.8742	0.7380	0.6897
1425 1426	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL)	0.6010 0.6384 0.7659	0.8005 0.8742 0.6575	0.7380 0.7076	0.6897 0.7283
1424 1425 1426 1427	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL)	$\begin{array}{c} 0.6010 \\ 0.6384 \\ 0.7659 \\ 0.8555 \end{array}$	$\begin{array}{c} 0.8003 \\ 0.8742 \\ 0.6575 \\ 0.8253 \end{array}$	0.7380 0.7076 <b>0.8401</b>	0.6897 0.7283 <b>0.8430</b>
1424 1425 1426 1427 1428	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL)	$\begin{array}{c} 0.6010 \\ 0.6384 \\ 0.7659 \\ 0.8555 \\ 0.6732 \end{array}$	$\begin{array}{c} 0.8005\\ 0.8742\\ 0.6575\\ 0.8253\\ 0.7308\end{array}$	0.7380 0.7076 <b>0.8401</b> 0.7008	0.6897 0.7283 <b>0.8430</b> 0.6891
1424 1425 1426 1427 1428 1429	Llama2-70b-chat (3 ICL) Llama3-70b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL)	$\begin{array}{c} 0.6010\\ 0.6384\\ 0.7659\\ 0.8555\\ 0.6732\\ 0.8141 \end{array}$	$\begin{array}{c} 0.8005\\ 0.8742\\ 0.6575\\ 0.8253\\ 0.7308\\ 0.8342\end{array}$	0.7076 0.7076 0.8401 0.7008 0.8240	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224
1424 1425 1426 1427 1428 1429 1430	Llama2-70b-chat (3 ICL) Llama3-70b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL)	0.6010 0.6384 0.7659 0.8555 0.6732 0.8141 CauseNet	$\begin{array}{c} 0.8003\\ 0.8742\\ 0.6575\\ 0.8253\\ 0.7308\\ 0.8342\end{array}$	0.7380 0.7076 <b>0.8401</b> 0.7008 0.8240	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224
1424 1425 1426 1427 1428 1429 1430 1431	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL)	0.6010 0.6384 0.7659 0.8555 0.6732 0.8141 CauseNet Precision↑	0.8003 0.8742 0.6575 0.8253 0.7308 0.8342 Recall↑	0.7380 0.7076 <b>0.8401</b> 0.7008 0.8240 <b>F1</b> ↑	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224 Accuracy↑
1424 1425 1426 1427 1428 1429 1430 1431 1432	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL)	0.6010 0.6384 0.7659 0.8555 0.6732 0.8141 CauseNet Precision↑ 0.7030	0.8003 0.8742 0.6575 0.8253 0.7308 0.8342 <b>Recall</b> ↑ 0.8931	0.7380 0.7076 <b>0.8401</b> 0.7008 0.8240 <b>F1</b> ↑ 0.7867	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224 Accuracy↑ 0.7582
1424 1425 1426 1427 1428 1429 1430 1431 1432 1433	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL) Llama2-7b-chat (3 ICL) Llama2-13b-chat (3 ICL)	0.6010 0.6384 0.7659 0.8555 0.6732 0.8141 CauseNet Precision↑ 0.7030 0.6625	0.8003 0.8742 0.6575 0.8253 0.7308 0.8342 <b>Recall</b> ↑ 0.8931 0.9213	0.7380 0.7076 <b>0.8401</b> 0.7008 0.8240 <b>F1</b> ↑ 0.7867 0.7708	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224 Accuracy↑ 0.7582 0.7260
1424 1425 1426 1427 1428 1429 1430 1431 1432 1433 1434	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL) Llama2-7b-chat (3 ICL) Llama2-13b-chat (3 ICL) Llama2-70b-chat (3 ICL)	0.6010 0.6384 0.7659 0.8555 0.6732 0.8141 CauseNet <b>Precision</b> ↑ 0.7030 0.6625 0.7359	0.8003 0.8742 0.6575 0.8253 0.7308 0.8342 <b>Recall</b> ↑ 0.8931 0.9213 0.9521	0.7380 0.7076 <b>0.8401</b> 0.7008 0.8240 <b>F1</b> ↑ 0.7867 0.7708 0.8302	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224 <b>Accuracy</b> ↑ 0.7582 0.7260 0.8053
1424 1425 1426 1427 1428 1429 1430 1431 1432 1433 1434 1435	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL) Llama2-7b-chat (3 ICL) Llama2-70b-chat (3 ICL) Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL)	0.6010 0.6384 0.7659 0.8555 0.6732 0.8141 CauseNet <b>Precision</b> ↑ 0.7030 0.6625 0.7359 0.8838 0.9220	0.8003 0.8742 0.6575 0.8253 0.7308 0.8342 <b>Recall</b> ↑ 0.8931 0.9213 0.9521 0.8296 0.8422	0.7380 0.7076 <b>0.8401</b> 0.7008 0.8240 <b>F1↑</b> 0.7867 0.7708 0.8302 0.8558	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224 Accuracy↑ 0.7582 0.7260 0.8053 0.8602
1424 1425 1426 1427 1428 1429 1430 1431 1432 1433 1434 1435 1436	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL) Llama2-7b-chat (3 ICL) Llama2-70b-chat (3 ICL) Llama2-70b-chat (3 ICL) Llama3-70b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL)	0.6010 0.6384 0.7659 0.8555 0.6732 0.8141 CauseNet Precision↑ 0.7030 0.6625 0.7359 0.8838 0.8939 0.8020	0.8003 0.8742 0.6575 0.8253 0.7308 0.8342	0.7380 0.7076 <b>0.8401</b> 0.7008 0.8240 <b>F1↑</b> 0.7867 0.7708 0.8302 0.8558 <b>0.9175</b>	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224 <b>Accuracy</b> ↑ 0.7582 0.7260 0.8053 0.8602 <b>0.9152</b> 0.921
1424 1425 1426 1427 1428 1429 1430 1431 1432 1433 1434 1435 1436 1437	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL) Llama2-7b-chat (3 ICL) Llama2-70b-chat (3 ICL) Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-3.5-turbo (3 ICL)	0.6010 0.6384 0.7659 0.8555 0.6732 0.8141 CauseNet Precision↑ 0.7030 0.6625 0.7359 0.8838 0.8939 0.8990 0.2506	0.8003 0.8742 0.6575 0.8253 0.7308 0.8342	0.7380 0.7076 <b>0.8401</b> 0.7008 0.8240 <b>F1↑</b> 0.7867 0.7708 0.8302 0.8558 <b>0.9175</b> 0.8923 0.8923	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224 <b>Accuracy</b> ↑ 0.7582 0.7260 0.8053 0.8602 <b>0.9152</b> 0.8931
1424 1425 1426 1427 1428 1429 1430 1431 1432 1433 1434 1435 1436 1437 1438	Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL) Llama2-7b-chat (3 ICL) Llama2-70b-chat (3 ICL) Llama3-8b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) Llama3-70b-Instruct (3 ICL) GPT-3.5-turbo (3 ICL) GPT-4o (3 ICL)	0.6010 0.6384 0.7659 0.8555 0.6732 0.8141 CauseNet Precision↑ 0.7030 0.6625 0.7359 0.8838 0.8939 0.8990 0.8596	0.8003 0.8742 0.6575 0.8253 0.7308 0.8342	0.7380 0.7076 <b>0.8401</b> 0.7008 0.8240 <b>F1↑</b> 0.7867 0.7708 0.8302 0.8558 <b>0.9175</b> 0.8923 0.9051	0.6897 0.7283 <b>0.8430</b> 0.6891 0.8224 <b>Accuracy</b> ↑ 0.7582 0.7260 0.8053 0.8602 <b>0.9152</b> 0.8931 0.8998

Table 18: Performance on causal direction identification task using Llama2 and Llama3 models of different sizes.