

CausalEval: Towards Better Causal Reasoning in Language Models

Longxuan Yu^{*†}, Delin Chen^{*‡}, Siheng Xiong^{*+}, Qingyang Wu[§],
Dawei Li[♡], Zhikai Chen[♣], Xiaoze Liu[♠], Liangming Pan[◇]

[†]UC San Diego [‡]UMass Amherst ⁺Georgia Tech [§]UCLA

[♡]Arizona State University [♣]Michigan State University

[♠]Purdue University [◇]University of Arizona

loy004@ucsd.edu delinchen@umass.edu liangmingpan@arizona.edu

Abstract

Causal reasoning (CR) is a crucial aspect of intelligence, essential for problem-solving, decision-making, and understanding the world. While language models (LMs) can generate rationales for their outputs, their ability to reliably perform causal reasoning remains uncertain, often falling short in tasks requiring a deep understanding of causality. In this paper, we introduce CausalEval, a comprehensive review of research aimed at enhancing LMs for causal reasoning, coupled with an empirical evaluation of current models and methods. We categorize existing methods based on the role of LMs: either as reasoning engines or as helpers providing knowledge or data to traditional CR methods, followed by a detailed discussion of methodologies in each category. We then assess the performance of current LMs and various enhancement methods on a range of causal reasoning tasks, providing key findings and in-depth analysis. Finally, we present insights from current studies and highlight promising directions for future research. We aim for this work to serve as a comprehensive resource, fostering further advancements in causal reasoning with LMs.¹

1 Introduction

Large Language Models (LLMs) (OpenAI et al., 2023; Dubey et al., 2024) have achieved significant success across various domains. Causal reasoning (CR), which includes causal discovery and inference, is a fundamental aspect of human intelligence, facilitating effective problem-solving, decision-making (Cavenaghi et al., 2024), and understanding of the world (Yao et al., 2021), with applications in healthcare (Mahmooda et al., 2014; Prospero et al., 2020), economics (Celli, 2022; Vermuri, 2015), biology (Ramsey and Andrews, 2018),

* Equal contribution.

¹Resources are available at <https://github.com/chend102/Awesome-LLM-causal-reasoning>.

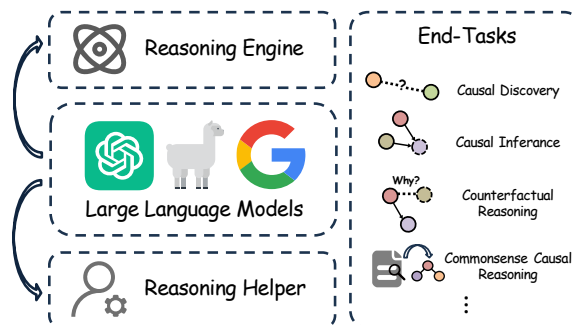


Figure 1: Large language models for causal reasoning: serving as reasoning engines or providing support to traditional methods in various end-tasks.

among other fields. While LLMs can often generate rationales for their outputs (Huang et al., 2022), it remains controversial whether they genuinely understand causal relationships (Kıcıman et al., 2023a). For example, Liu et al. (2024a) found that LLMs struggle to integrate pre-trained causal knowledge in statistical reasoning tasks.

Enhancing the causal reasoning capabilities of LLMs differs from improving their general reasoning skills. CR focuses on understanding causal relationships, while general reasoning primarily relies on identifying correlations (Bao et al., 2024). Unlike general reasoning, which often involves simpler, more static relationships, CR requires deeper analysis and managing uncertainty to determine how events are causally linked (Pearl, 2009). Additionally, CR must account for confounding variables (Greenland et al., 1999), where an observed correlation between two variables is due to the influence of a third variable. Counterfactual reasoning (Pearl, 2009), often necessary for causality verification, is also challenging, as it involves reasoning about hypothetical scenarios that did not actually occur.

Existing research on enhancing the CR capabilities of LLMs can be broadly categorized based on their role: either as the **engine** directly performing

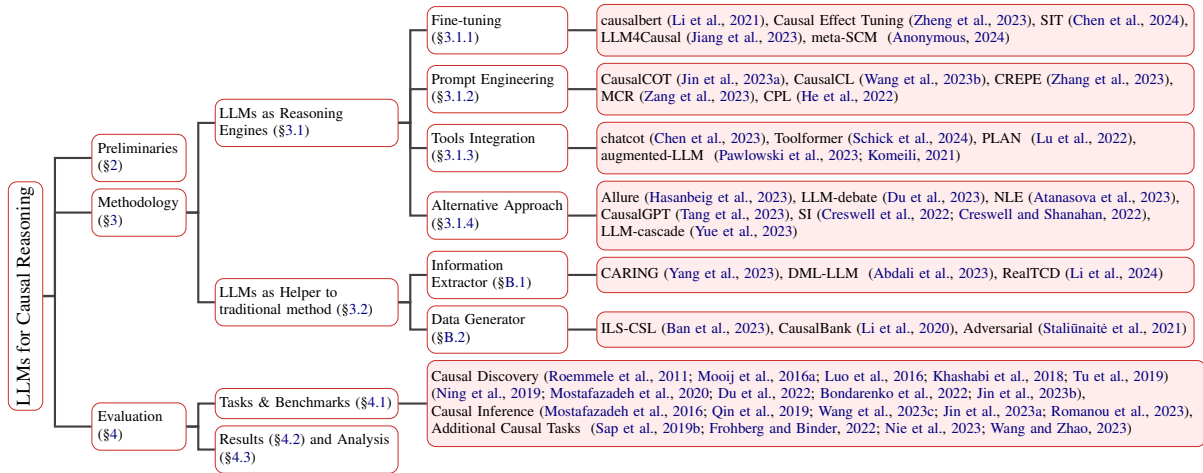


Figure 2: Structure overview. We categorize the role of LLMs in causal reasoning into two main functions: as reasoning engine and as helper. Each function is further divided into specific methodologies. We also outline the evaluation process, including tasks, benchmarks, results and analysis.

causal inferences or as the **helper** providing extra information to traditional CR methods (Figure 1). Despite significant progress, challenges such as the shallow depth of causal reasoning skills and the scarcity of high-quality data still require further exploration (Zečević et al., 2023). Addressing these challenges is crucial, as enhancing causal reasoning not only improves the general performance of LLMs but also provides substantial benefits for tackling complex reasoning problems and supporting decision-making (Liu et al., 2024c). Thus, integrating causal inference into the model architecture (Melnychuk et al., 2022) or during the pre-training phase (Rohekar et al., 2023) presents a promising direction for future research.

In this work, we focus on methodologies to enhance the causal reasoning capabilities of LLMs, exploring practical challenges and identifying potential opportunities. We introduce a fine-grained **taxonomy** (§3) that systematically reviews and summarizes current techniques employed in LLMs to improve CR. To evaluate these capabilities, we conduct extensive **assessments** (§4) of widely-used LLMs on existing benchmarks, providing key findings and detailed analysis. Additionally, we discuss the **limitations** of current approaches (§5) and propose promising directions for **future research** (§6). Previous work (Liu et al., 2024b; Wan et al., 2024) reviews causal reasoning in LLMs, focusing more on the perspective of *traditional causal reasoning*. In contrast, our focus lies in the *methodological advancements* that enhance the causal reasoning abilities of LLMs. More importantly, we conduct *practical evaluations* of LLMs on various

Levels	Equation	Query
Association	$P(y x)$	What is?
Intervention	$P(y \text{do}(x), z)$	What if?
Counterfactual	$P(y_x x', y')$	Why?

Table 1: The Ladder of Causality

CR tasks, offering valuable experimental results and insights that are often lacking in related literature. By bridging theoretical advancements and practical evaluations, our study fills a critical gap in the literature and offers insights to guide the development of LLMs with enhanced causal reasoning.

2 Preliminaries

CR in the context of machine learning (ML) involves understanding and modeling cause-and-effect relationships within data, which goes beyond the traditional ML focus on correlations. Traditional ML methods often identify patterns or correlations, limiting their applicability in understanding the mechanisms behind complex systems (Peters et al., 2017). Causal ML, on the other hand, addresses this gap by focusing on cause-and-effect relationships, which is crucial for making informed decisions, predicting the effects of actions, and understanding complex systems (Prosperi et al., 2020; Fernández-Loría and Provost, 2022; Lagemann et al., 2023; Berrevoets et al., 2023).

The CR framework can be categorized into three levels: association, intervention, and counterfactual (Pearl et al., 2000; Pearl and Mackenzie, 2018), which form “the ladder of causality” (Table 1). As-

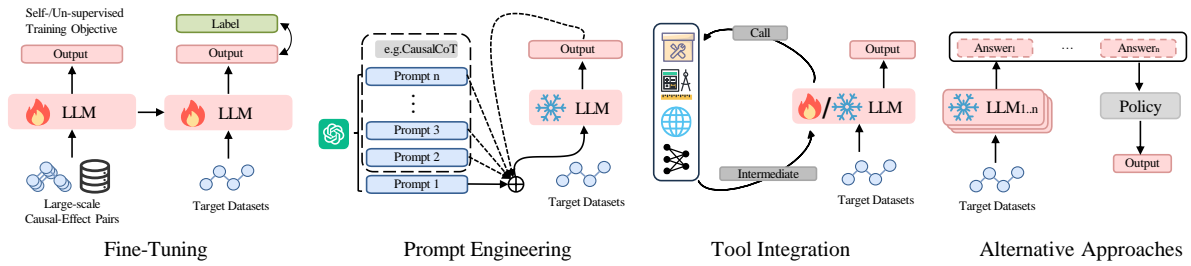


Figure 3: Overview of methods for LLMs as causal reasoning engines. (a) Fine-Tuning: Adapting LLMs using large-scale causal-effect pairs and target datasets. (b) Prompt Engineering: Crafting targeted prompts to elicit the internal CR capabilities. (c) Tool Integration: Leveraging external tools to support LLMs in performing intermediate steps. (d) Alternative Approaches: Implementing additional methods, such as iterative improvement protocols, multi-agent systems, and rationale-based evaluation.

sociation is formalized in statistical terms through conditional probability between variables. Intervention is represented using Pearl’s do-operator to explore causal effects, while counterfactual reasoning involves hypothetical scenarios that address potential outcomes under different circumstances. Structural Causal Models (SCMs), often visualized with Directed Acyclic Graphs (DAGs) where nodes represent variables and edges indicate causal relationships (Pearl, 2009), are widely used to represent the interplay of these components. SCMs provide a mathematical framework for formalizing causal relationships within a system, thus serving as the foundation for causal reasoning and causal learning (Peters et al., 2017). More details about CR tasks, and the comparison between LLMs and traditional methods, are provided in Appendix A.

3 Towards Causal Reasoning in Large Language Models

We separate the roles of language models in CR into two categories. First, LLMs can serve as **causal reasoning engines**, employing methods such as fine-tuning, prompt engineering, external tool integration, and alternative approaches to directly draw causal inferences from input prompts. Second, LLMs can act as **helpers** to traditional methods by extracting causal information and generating causality data to enhance causal analysis across various domains.

3.1 Serving as Causal Reasoning Engines

3.1.1 Injecting Causal Knowledge through Fine-Tuning

Fine-tuning language models has been proven to be a simple and effective technique for transferring learned knowledge from a pre-trained model to

new downstream tasks. Compared to training from scratch, it can significantly reduce training time, and improve performance on new tasks (Lester et al., 2021). For example, Li et al. (2021) used self-supervised learning strategies, such as causal pairs ranking and fine-tuning, to inject causal knowledge into pre-trained models. Additionally, Zheng et al. (2023) proposed Causal Effect Tuning, which is based on causal inference, to learn new knowledge from target data while preserving old knowledge from pre-trained models, thereby building a causal path without storing the pre-training data.

Although scaling language models increases the effectiveness of their prior knowledge in addressing complex causal questions, most training corpora collected from online sources consist of unstructured textual data rather than structured causal annotations. To address this, Jiang et al. (2023) collected a custom dataset containing causal questions and interpreted intents in a structured format for supervised fine-tuning. Building on this, Chen et al. (2024) combined the structural causal model with instruction tuning (Liu et al., 2023a; Wei et al., 2022; Chung et al., 2024) to develop a meta Structural Causal Model, which unified different NLP tasks under a single causal structure and introduced the Structural Instruction Tuning method to learn task-specific causal representations that mimic the causal factors for each task.

Limitations Despite these progresses, LLMs often lack clearly defined supervision for causality, as most ML datasets contain only observational data without interventions, making key causal quantities, such as treatment effects, unknown (Kıcıman et al., 2023b). Furthermore, causality datasets often feature complex, heterogeneous relationships between variables (Huang et al., 2020), which are

more challenging to capture compared to linguistic or perceptual patterns. Thus, developing methods to bridge this gap, such as integrating causal structure learning or leveraging synthetic data for interventions, remains an important area of research.

3.1.2 Embedding Causal Knowledge within Contexts

Prompt-based methods have become crucial in unlocking the internal reasoning capabilities of LLMs across various tasks. Compared with general reasoning, CR tasks, particularly multi-step ones, pose greater challenges due to factors such as the complexity of understanding cause-and-effect relationships, the need for reasoning about interventions and counterfactuals, and the presence of confounding variables (Gandhi et al., 2023; Tan et al., 2023; Ban et al., 2023; Abdali et al., 2023).

To address these challenges, various innovative prompting strategies have been proposed (Liu et al., 2023b; Bagheri et al., 2024). Building on Chain-of-Thought (CoT) prompting (Wei et al., 2023), CausalCoT (Jin et al., 2023a) enables LLMs to reason step-by-step through causal relationships, significantly improving causal inference tasks involving multiple causes and effects. Recent advancements, such as Causal Contextualized Learning (CausalCL) (Wang et al., 2023b), and Zhang et al. (2023) further enhances CR by combining code-like prompting with chain-of-thought prompting, boosting the performance for multihop event reasoning. Multi-modal Causal Reasoning (MCR) (Zang et al., 2023) extends CR to multi-modal settings, integrating visual and textual information for a comprehensive analysis. Counterfactual Prompting Learning (CPL) (He et al., 2022) utilizes counterfactual thinking to guide models in exploring alternative scenarios, thereby improving their understanding of causal mechanisms.

These advanced prompting techniques collectively enhance the ability of LLMs to understand causal relationships, thereby improving precision in CR tasks and enabling models to manage long-term dependencies more effectively (Lu et al., 2022). Furthermore, they reduce biases from training data, enrich the model’s knowledge with external information and commonsense, and do so without substantial increases in computational costs (Yao et al., 2022; Jin et al., 2023a).

Limitations Causal prompt engineering depends on the model’s inherent causality knowledge,

which is often insufficient due to the correlation-driven nature of its modeling approach (Zhang et al., 2024). Crafting effective prompts is labor-intensive, requiring significant domain expertise. Furthermore, the sensitivity to prompt variations can lead to inconsistent outputs, reducing robustness (Perez et al., 2021). Additionally, multi-step causal prompts can increase computational inefficiency, demanding more resources and time. These limitations highlight the necessity of complementary approaches, such as explicit causal modeling or specialized fine-tuning.

3.1.3 External Causal Tool Integration

The motivation for integrating external tools into LLMs is that effective causal reasoning requires more than just internal model representations. By accessing structured knowledge bases and performing specialized computations, external tools enhance LLMs’ capacity for complex causal reasoning, complementing their inherent capabilities. The ability to use these tools can be acquired through methods such as fine-tuning, in-context learning, and prompt manipulation. This integration establishes a crucial link between causal reasoning and language understanding, enabling LLMs to handle complex causal tasks more effectively.

For methods that incorporate tools in prompts, PLAN (Lu et al., 2023) introduced a causal framework for procedural planning that leverages external knowledge bases, such as ConceptNet (Speer et al., 2017). It semantically parses tasks into an entity set to retrieve relevant sub-graphs for improved planning. Additionally, Pawlowski et al. (2023) compared context augmentation and tool augmentation methods: context augmentation uses language models for supplementary operations beyond the causal expert system, whereas tool augmentation incorporates basic Python tools to manipulate outputs from the expert system. Their experiments showed that context-augmented LLMs were more prone to errors compared to tool-augmented ones. For methods involving supervised fine-tuning, dialogue agents such as those by Shuster et al. (2022); Komeili (2021) are trained to effectively use search engines.

Limitations While integrating external causal tools can enhance LLMs, it also introduces several challenges. The effectiveness of these tools depends heavily on the quality and completeness of external resources, which can lead to errors when

inadequate. The alignment between the model and tools is often complex and domain-specific, limiting the model’s generalizability. Additionally, reliance on external tools may slow response times and increase the risk of errors, as observed with context-augmented models (Pawlowski et al., 2023). These challenges indicate that external tools should be used in combination with other approaches to fully address complex CR tasks.

3.1.4 Alternative Approaches

Recent studies have explored various alternative approaches to enhance the causal reasoning capabilities of LLMs, including iterative improvement protocols, multi-agent systems, and rationale-based evaluation. For example, ALLURE (Hasanbeig et al., 2023) uses a systematic auditing and *iterative improvement protocol* with in-context learning examples to refine LLM evaluation. ECHO (Xie et al., 2023) introduces a visio-linguistic dataset based on real-world scenarios to evaluate event causality through human-centric reasoning.

Agentic frameworks have also been employed for causal reasoning in LLMs. CausalGPT (Tang et al., 2023) employs a multi-agent system comprising reasoners and evaluators, with reasoners generating answers and evaluators assessing factual accuracy and causal consistency. Du et al. (2023) presents a debate framework where multiple instances of LLMs debate a causality-related query to reach a consensus. *Selection and inference-based techniques* are another approach. Creswell et al. (2022) alternate between selection and inference steps using pre-trained LLMs to generate interpretable causal reasoning steps, while Creswell and Shanahan (2022) chain reasoning steps to produce valid reasoning traces. *Rationale-based evaluation methods* are also receiving increased attention. Atanasova et al. (2023) use rationale-based methods to conduct tests for evaluating the faithfulness of natural language explanations. Abdali et al. (2023) generate a prior model to address domain knowledge gaps, serving as a heuristic for assessing the informativeness of user feedback. Additionally, Yue et al. (2023) utilize a cascaded model incorporating CoT and Program-of-Thought (PoT) approaches, along with answer sampling and consistency checking, to enhance causal reasoning.

3.2 Enhancing Traditional CR Methods

Beyond serving as causal reasoning engines, LLMs also act as helpers in causal analysis across do-

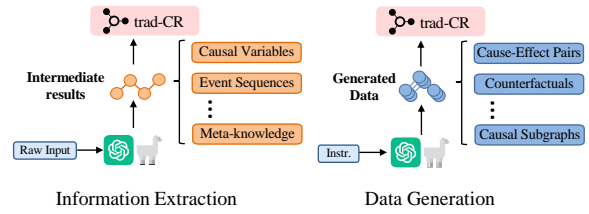


Figure 4: Overview of methods for using LLMs to enhance traditional approaches. (a) Information Extraction: Extracting causal variables and events from text and adjusting for biases. (b) Data Generation: Generating synthetic causal data and forming hypotheses.

main (Figure 4), fulfilling several supportive roles: (1) *Causal Information Extraction*: LLMs can extract causal variables, event sequences, and meta-knowledge from unstructured text, enhancing the quality of causal discovery in real-world applications (Yang et al., 2023; Abdali et al., 2023; Li et al., 2024). Additionally, by leveraging prior knowledge, they identify and adjust for biases and confounders, improving the robustness of causal inferences (Ban et al., 2023). (2) *Causality Data Generation & Augmentation*: LLMs generate synthetic, causally relevant data, such as cause-effect pairs, counterfactuals, and causal subgraphs, thereby enriching existing datasets and supporting traditional CR methods (Staliūnaitė et al., 2021). Furthermore, they form causal hypotheses from observed data patterns, generate plausible counterfactuals, and reason through complex interventions, guiding CR methods to explore less obvious causal pathways and enabling effective outcome evaluation (Gat et al., 2023). These supportive roles illustrate the potential of LLMs to complement traditional causal analysis methods by bridging gaps in data, knowledge, and reasoning. Through their ability to extract hidden causal information and generate synthetic data, LLMs enable more robust and comprehensive causal reasoning. More details on these supportive roles are provided in Appendix B.

4 Evaluating Casual Reasoning in Large Language Models

4.1 Experimental Setup

We first categorize the end tasks into three groups: causal discovery, causal inference, and additional causal tasks. For each category, we evaluate recent LLMs with strategies such as zero-shot, few-shot, direct I/O prompting, and Chain-of-Thought (CoT) reasoning. To reduce variance, we evaluate

Model	Causal Discovery				Causal Inference		Additional Causal Tasks		
	COPA	NPDS	e-CARE	Corr2Cause	CLADDER	CRAB	CRASS	MoCa	Tram
Mistral-7B	86.2 \downarrow 3.0	61.0 \downarrow 10.0	63.4 \uparrow 5.0	24.2 \uparrow 4.6	55.0 \downarrow 6.4	35.6 \downarrow 1.4	56.7 \uparrow 15.0	55.6 \uparrow 3.4	75.8 \downarrow 2.2
+ CoT	88.8 \downarrow 7.4	65.0 \downarrow 1.0	69.8 \uparrow 0.8	36.6 \downarrow 3.2	60.0 \downarrow 4.4	39.0 \downarrow 3.6	71.7 \downarrow 1.7	53.5 \uparrow 0.7	76.8 \downarrow 3.4
Mistral-8 \times 7B	75.2 \uparrow 7.8	73.0 \downarrow 22.0	58.6 \uparrow 15.2	40.0 \downarrow 9.0	57.8 \downarrow 2.2	34.6 \downarrow 8.8	58.3 \uparrow 20.0	57.6 \downarrow 2.7	67.6 \uparrow 7.4
+ CoT	82.8 \uparrow 6.4	65.0 \downarrow 3.0	72.0 \uparrow 5.6	33.6 \uparrow 10.0	58.6 \uparrow 4.4	35.6 \downarrow 1.4	71.7 \uparrow 8.3	59.0 \uparrow 2.1	69.6 \uparrow 6.0
Gemma-2-9B	91.2 \downarrow 8.4	67.0 \downarrow 4.0	69.4 \uparrow 5.4	31.0 \downarrow 1.2	56.8 \downarrow 1.6	45.2 \downarrow 2.6	76.7 \downarrow 1.7	57.6 \uparrow 2.1	79.2 \downarrow 0.8
+ CoT	96.0 \uparrow 0.0	77.0 \downarrow 26.0	80.6 \uparrow 1.4	37.6 \uparrow 2.4	56.0 \downarrow 1.0	40.2 \downarrow 2.8	76.7 \uparrow 8.3	63.2 \downarrow 3.5	83.2 \downarrow 1.8
Gemma-2-27B	94.0 \downarrow 6.5	74.0 \downarrow 10.0	70.0 \uparrow 8.0	23.8 \uparrow 12.4	57.6 \uparrow 6.4	47.4 \downarrow 3.8	86.7 \downarrow 6.7	54.9 \uparrow 7.6	78.6 \uparrow 1.6
+ CoT	97.8 \downarrow 0.6	80.0 \downarrow 17.0	83.8 \downarrow 1.0	42.4 \downarrow 1.4	57.2 \uparrow 6.2	45.4 \downarrow 5.6	91.7 \downarrow 5.0	58.3 \uparrow 9.8	84.8 \downarrow 3.2
LLaMA3-8B	81.4 \uparrow 9.4	57.0 \uparrow 6.0	62.8 \uparrow 8.2	36.0 \downarrow 4.0	53.0 \uparrow 1.2	31.4 \uparrow 10.0	61.7 \uparrow 15.0	58.3 \downarrow 1.4	68.6 \uparrow 9.0
+ CoT	89.8 \uparrow 2.4	56.0 \uparrow 1.0	76.0 \downarrow 1.6	27.2 \uparrow 6.0	56.8 \uparrow 1.4	35.6 \uparrow 0.8	76.7 \uparrow 1.6	51.4 \downarrow 2.8	78.2 \uparrow 1.2
LLaMA3-70B	98.0 \uparrow 0.8	62.0 \uparrow 4.0	76.6 \uparrow 4.8	27.2 \uparrow 4.8	59.2 \uparrow 12.2	33.6 \uparrow 5.0	88.3 \uparrow 0.0	54.9 \uparrow 5.5	81.2 \downarrow 2.0
+ CoT	98.2 \uparrow 0.4	76.0 \downarrow 12.0	82.8 \uparrow 2.0	34.8 \uparrow 3.2	55.0 \uparrow 20.8	43.0 \uparrow 1.0	88.3 \uparrow 6.7	59.0 \uparrow 2.1	83.4 \downarrow 2.4
LLaMA3.1-8B	82.4 \uparrow 12.2	48.0 \downarrow 12.0	61.2 \uparrow 14.6	37.2 \downarrow 14.6	53.6 \downarrow 3.8	34.6 \uparrow 8.8	66.7 \uparrow 16.6	53.5 \downarrow 9.1	68.6 \uparrow 9.2
+ CoT	88.2 \uparrow 6.8	63.0 \downarrow 3.0	77.8 \uparrow 1.4	31.2 \downarrow 4.6	48.4 \uparrow 4.4	35.4 \downarrow 4.0	70.0 \uparrow 5.0	53.5 \downarrow 0.7	79.6 \downarrow 2.2
LLaMA3.1-70B	98.2 \uparrow 0.8	63.0 \uparrow 3.0	78.2 \uparrow 1.4	23.2 \uparrow 3.2	59.2 \uparrow 15.2	37.2 \uparrow 3.4	88.3 \uparrow 6.7	56.9 \uparrow 1.4	80.4 \downarrow 1.0
+ CoT	98.8 \uparrow 0.8	68.0 \downarrow 5.0	81.4 \uparrow 4.4	23.6 \uparrow 13.2	60.6 \uparrow 10.2	44.6 \downarrow 3.6	96.7 \downarrow 3.4	55.6 \uparrow 2.7	83.8 \downarrow 3.0
DeepSeek-V3*	98.0 \uparrow 1.0	62.0 \uparrow 7.0	83.0 \uparrow 3.0	35.0 \uparrow 6.0	62.0 \uparrow 10.0	50.0 \uparrow 5.0	95.0 \uparrow 1.7	63.0 \uparrow 2.0	87.0 \uparrow 1.0
+ CoT*	100.0 \uparrow 0.0	64.0 \uparrow 7.0	84.0 \uparrow 8.0	37.0 \uparrow 9.0	76.0 \uparrow 8.0	56.0 \downarrow 6.0	96.7 \uparrow 0.0	69.0 \downarrow 1.0	86.0 \uparrow 2.0
DeepSeek-R1*	100.0	61.0	84.0	38.0	92.0	51.0	95.0	61.0	82.0
GPT-3.5-turbo	89.8 \uparrow 3.4	54.0 \downarrow 5.0	72.0 \downarrow 2.2	41.8 \uparrow 8.4	52.4 \uparrow 6.8	32.0 \uparrow 6.0	75.0 \uparrow 10.0	56.2 \downarrow 1.3	76.6 \downarrow 0.6
+ CoT	95.2 \uparrow 1.2	70.0 \downarrow 9.0	82.8 \downarrow 4.2	34.6 \downarrow 0.6	54.8 \uparrow 4.4	32.6 \uparrow 11.6	85.0 \uparrow 5.0	56.9 \uparrow 2.1	82.4 \downarrow 1.0
GPT-4o*	100.0 \uparrow 0.0	56.0 \uparrow 11.0	85.0 \uparrow 1.0	47.0 \uparrow 2.0	61.0 \uparrow 18.0	33.0 \uparrow 4.0	95.0 \uparrow 3.3	61.0 \downarrow 4.0	84.0 \uparrow 1.0
+ CoT*	99.0 \uparrow 1.0	73.0 \downarrow 1.0	90.0 \downarrow 3.0	49.0 \uparrow 4.0	69.0 \uparrow 10.0	42.0 \uparrow 3.0	100.0 \downarrow 1.7	61.0 \downarrow 1.0	87.0 \uparrow 1.0
Claude 3.5*	98.0 \uparrow 1.0	63.0 \uparrow 5.0	72.0 \uparrow 12.0	45.0 \uparrow 2.0	70.0 \uparrow 13.0	57.0 \downarrow 1.0	95.0 \uparrow 3.3	42.0 \uparrow 20.0	84.0 \uparrow 0.0
+ CoT*	99.0 \uparrow 1.0	71.0 \uparrow 0.0	86.0 \uparrow 4.0	35.0 \uparrow 8.0	72.0 \uparrow 9.0	58.0 \uparrow 2.0	98.3 \downarrow 3.3	59.0 \uparrow 5.0	88.0 \downarrow 3.0
o1*	100.0	56.0	80.0	59.0	89.0	40.0	96.7	60.0	84.0
o3-mini	99.8	56.0	79.6	60.6	92.2	32.2	93.3	63.9	81.4
Human	95.8	97.7	92.0	94.5	94.8	98.7	98.2	92.0	98.8

Table 2: Overall performance comparison of models across causal tasks. Performance is measured using pass@1 accuracy, with the best results in each category in bold. Zero-shot scores are used as the baseline, while the green and blue arrows indicate the change brought by few-shot learning. *Model is evaluated on 100 random samples due to cost and time constraints.

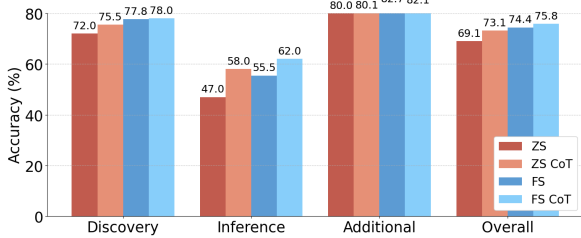


Figure 5: Performance gap of zero-shot (ZS) and few-shot (FS) learning, with and without CoT prompting. The results are from the best performing model (GPT-4o).

each question 5 times using greedy decoding and compute the average accuracy (pass@1) for that question. We then calculate the overall average accuracy across all problems. Specifically, we use COPA (Roemmele et al., 2011), NPDS (Tu et al., 2019), e-CARE (Du et al., 2022), and Corr2Cause (Jin et al., 2023b) for causal discovery; CLADDER (Jin et al., 2023a) and CRAB (Romanou et al., 2023) for causal inference; and CRASS (Frohberg and Binder, 2022), MoCa (Nie et al., 2023), and Tram (Wang and Zhao, 2023) for additional causal tasks. Further details on task definitions, dataset statistics and examples, and implementation details are provided in Appendix C and Appendix D.

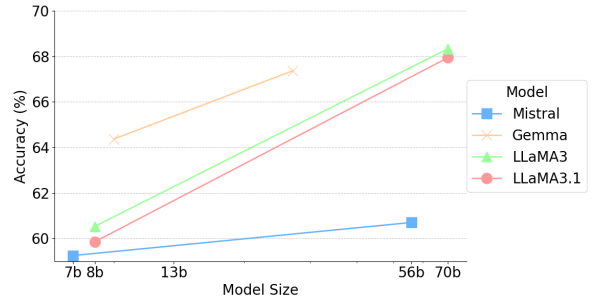


Figure 6: Scaling effect of model size and overall causal reasoning performance. The x-axis (model size) is shown in the log scale. Mistral 8 \times 7b (MoE) is treated as equivalent to 56b.

4.2 Main Results

Overall performance is shown in Table 2. We summarize the key observations as follows:

LLMs significantly underperform humans in many CR Tasks. Our results reveal a notable performance gap between LLMs and humans across many causal reasoning tasks, as illustrated in the Radar chart (Figure 8, Appendix D). Despite recent advancements, LLMs continue to face challenges in fully grasping and applying causal reasoning, often relying on surface-level patterns and struggling with multi-step reasoning. These limitations

Method	Causal Discovery				Causal Inference		Additional Causal Tasks		
	COPA	NPDS	e-CARE	Corr2Cause	CLADDER	CRAB	CRASS	MoCa	Tram
GPT-3.5-turbo	89.8	54.0	72.0	41.8	52.4	32.0	75.0	56.2	76.6
+ CoT	95.2	70.0	82.8	34.6	54.8	32.6	85.0	56.9	82.4
+ causalCoT*	49.0	28.0	37.0	35.0	45.0	19.0	52.0	45.0	48.0
+ FunctionCall*	78.0	58.0	63.0	35.0	57.0	44.0	68.3	53.0	75.0
GPT-4o*	100.0	56.0	85.0	47.0	61.0	33.0	95.0	61.0	84.0
+ CoT*	99.0	73.0	90.0	49.0	69.0	42.0	100.0	61.0	87.0
+ causalCoT*	94.0	58.0	75.0	47.0	75.0	36.0	93.0	56.0	78.0
+ FunctionCall*	100.0	73.0	73.0	49.0	76.0	42.0	98.3	61.0	84.0
Human	95.8	97.7	92.0	94.5	94.8	98.7	98.2	92.0	98.8

Table 3: Performance comparison between Direct IO, CoT, causalCOT prompting and functionCall (Python script) on GPT-3.5-turbo and GPT-4o across various causal reasoning datasets. Performance is measured using pass@1 accuracy, with the best results in each category in bold. *Model is evaluated on 100 random samples due to cost and time constraints.

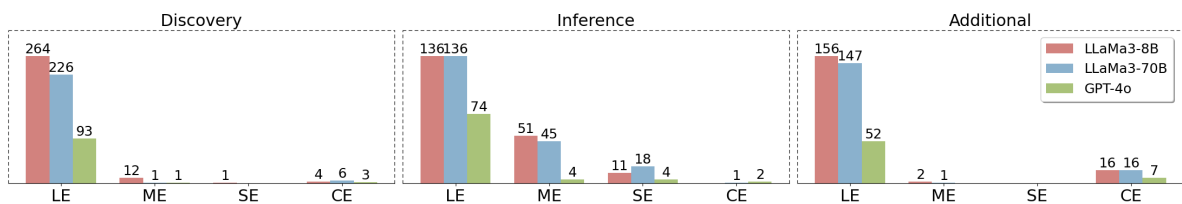


Figure 7: Error analysis across different types of causal reasoning tasks. For each task type, multiple datasets are analyzed with 100 randomly selected incorrect test samples. The y-axis shows the number of incorrect samples per error type: LE (Logical Errors), PE (Probabilistic Errors), SE (Statistical Errors), and CE (Contextual Errors). Detailed definition of each error type is given in Appendix D.

reduce their effectiveness in complex scenarios. Future work should focus on enhancing LLMs’ ability to handle intricate causal structures, potentially through hybrid methods that integrate symbolic reasoning with statistical approaches. Addressing these challenges will be critical for bridging the gap between LLMs and human performance in causal reasoning.

Performance disparity exists between discovery and inference tasks due to contextual and multi-variable reasoning demands. Our results reveal distinct performance trends across causal tasks. All models achieve high accuracy in certain discovery tasks, such as COPA and e-CARE, which require identifying potential causal relationships from observational data. However, their performance drops significantly in inference tasks, particularly when assessing causal strengths from long contexts (CRAB). Discovery tasks involving complex multi-variable interactions, such as Corr2Cause, also present significant challenges. This disparity underscores a fundamental limitation: while models are proficient at identifying potential causal relation candidates, they struggle with evaluating causal strengths and navigating in-

tricate variable dependencies.

Deep thinking models exhibit targeted gains in structured causal reasoning. Our analysis reveals that deep thinking models (DeepSeek-R1, o1, o3-mini) demonstrate particularly strong performance in structured causal reasoning tasks, especially in CLADDER where they achieve remarkable accuracy (92.0% and 92.2% respectively). This suggests that extended thinking contexts may be particularly beneficial for formal causal inference and graph-based reasoning. However, this advantage does not consistently extend to other causal tasks - performance on benchmarks such as COPA, NPDS, and e-CARE remains comparable to or worse than that of standard models. This trend indicates that deep thinking approaches primarily enhance structured causal reasoning rather than general causal understanding.

4.3 Analysis and Discussion

We further analyze the impact of various factors on accuracy below.

Zero-shot or few-shot under CoT prompting. Figure 5 shows that both zero-shot and few-shot models benefit from CoT prompting, with few-shot

CoT achieving the highest overall performance. While few-shot consistently outperforms zero-shot across all tasks, the models struggle the most with inference tasks, where accuracy remains lower than in discovery and additional tasks. These findings highlight the effectiveness of combining example-based learning with structured reasoning but also underscore the need for further improvements, especially in inference tasks, to enhance model performance across complex causal scenarios.

Model size vs. Causal reasoning performance. Figure 6 demonstrates that causal reasoning performance improves consistently with increasing model size across all evaluated models. Both LLaMA3 and LLaMA3.1 achieve significant gains, with Gemma also exhibiting strong improvements. In contrast, Mistral shows only modest improvements as model size increases, likely due to its Mixture of Experts (MoE) architecture. These results underscore the importance of both model size and architecture in advancing causal reasoning capabilities.

Error analysis. Through systematic analysis of model outputs, we identify four common types of errors in causal reasoning.

Logical errors occur when the model makes flawed cause-effect conclusions, such as confusing correlation with causation or failing to distinguish between necessary and sufficient conditions. For example, when asked whether someone caused a car to start in a scenario with multiple people each bringing a working battery, the model incorrectly attributes causation based solely on temporal sequence rather than considering the true causal relationship.

Probabilistic errors arise when the model incorrectly handles numerical reasoning about likelihood and chance, including mishandling of conditional probabilities and confounding variables. In probability comparison tasks, the model often makes computational errors or misinterprets statistical dependencies, leading to incorrect conclusions about relative likelihoods.

Statistical errors manifest when the model fails to properly account for complex relationships between multiple factors, such as indirect effects and selection bias. For instance, in scenarios involving multiple causal paths, the model often overlooks intermediate factors or assumes direct relationships where more complex dependencies exist.

Contextual errors occur when the model fails to incorporate relevant background information or

subtle details that should affect its causal reasoning. This often results in oversimplified conclusions that ignore important contextual factors or conditions that would alter the causal relationship.

Our analysis (Figure 7) shows that logical errors are the most prevalent type of failure. Potential approaches to address these challenges include enhanced representation of causal relationships in the model architecture, systematic consistency checks for reasoning processes, and diversified training data covering a broader range of causal scenarios. Detailed examples of each error type and analysis of these approaches are provided in Appendix D.

Comparison with enhancement methods. We evaluated two methods for improving causal reasoning: CausalCoT (Jin et al., 2023a) and Function Call. CausalCoT provides structured steps for causal inference with domain-specific guidance, while Function Call leverages Python scripts for precise mathematical computations. Our results show that effectiveness varies by task type. For formal reasoning tasks involving causal graph analysis, CausalCoT outperforms other approaches due to its structured inference steps (75.0% on CLADDER). In contrast, Function Call excels in numerical tasks, particularly on CLADDER (GPT-4o: 76.0% vs 69.0% with vanilla CoT), and performs comparably well on other computation-heavy datasets. However, for commonsense causal reasoning tasks like COPA and e-CARE, both methods underperform: CausalCoT shows substantial drops (40.8% and 35% respectively for GPT-3.5-turbo), while Function Call provides minimal gains (COPA: 78.0% vs 95.2% with vanilla CoT), suggesting that vanilla CoT remains the most effective for these tasks. These findings highlight the importance of task-aware method selection: Function Call for quantitative analysis, CausalCoT for structured causal inference with graphs, and vanilla CoT for commonsense causal reasoning.

5 Findings and Reflections

LLMs demonstrate shallow causal reasoning skills. Evaluation metrics often emphasize task-specific accuracy, leaving open questions about the depth of the model’s causal reasoning (Zečević et al., 2023). For instance, Hong et al. (2023) investigated how LLMs process event B in a script-based story, which causally depends on a prior event A. While large models performed better than smaller ones in the $\neg A \rightarrow B$ condition, all models

(even fine-tuned with code and instructions) failed to replicate human behavior when the cause was omitted (nil), indicating that LLMs lack sufficient representation of commonsense CR.

Integration/Collaboration enhances the performance. Tasks that challenge a single LLM can often be addressed by integration with traditional methods, external tools, or multiple LLMs (Tang et al., 2023; Du et al., 2023). Recent approaches, such as neuro-symbolic methods (Yang et al., 2023) and iterative processes (Ban et al., 2023), integrate LLMs with statistical analysis and causal structure learning, improving the quality of causal inferences. The integration with causal analysis tools increases LLMs' access to causal knowledge across domains (Kıcıman et al., 2023b). As intermediaries between specialized tools and natural language, LLMs make causal reasoning more intuitive and user-friendly. This collaboration fosters tools for comprehensively exploring causal relationships in complex systems.

Diverse and standardized causal reasoning benchmarks are needed. Comprehensive evaluation of the causal reasoning abilities of LLMs requires more diverse and standardized benchmarks. Existing benchmarks tend to be relatively simple and domain-specific. Future benchmarks should encompass a broader spectrum of causal reasoning challenges, including multi-step tasks with increasing complexity, while maintaining consistent causal structures. Such datasets would enable LLMs to explore a wider range of causal reasoning scenarios and better adapt to real-world applications (Yang et al., 2024).

6 Future Directions

Deep understanding of existing theories. Integrating causal reasoning capabilities into LLMs presents unique challenges, especially when requiring adherence to established theories. These theories, such as the potential outcomes framework, graph-based causal inference methods, and structural equation modeling, often rely on foundational assumptions, including the stable unit treatment value assumption, ignorability/unconfoundedness, and positivity (Liu et al., 2024b). Considering these assumptions necessitates innovative strategies, such as incorporating specific post-training data and implementing filtering mechanisms to ensure the model's alignment with the theories.

Data efficiency for end-to-end learning. LLMs

have the potential to perform end-to-end causal reasoning, but their effectiveness depends on substantial data to learn and infer causal relationships. A major challenge is the scarcity of high-quality causality data, particularly for counterfactual reasoning, where alternative scenarios are often unavailable or difficult to generate. Improving data efficiency is therefore a critical research direction. Techniques such as analogical reasoning, imaginative simulation, and inferential methods can maximize the utility of existing data by generating synthetic datasets, filling data gaps, or extrapolating from limited information. These approaches enhance the model's capacity to learn causal relationships with constrained training data.

Causality-based internal mechanisms for LLMs.

Incorporating causal reasoning into LLMs internally is supposed to enhance their accuracy and depth in complex problem-solving scenarios (Kıcıman et al., 2023b). Thus, it is crucial to find feasible implementation ways so that LLMs generate content grounded in causality. One strategy could be the introduction of a dual-network architecture, where one is similar to current model and the other focuses on learning causal relations. Both networks are jointly trained to perform various tasks. Another approach involves modifying attention mechanisms to prioritize causal relevance. By embedding internal causal mechanisms, we believe LLMs can enhance their reasoning, address fairness and safety, and better handle multi-modal inputs through improved understanding of causal relations, leading to more accurate and transparent decisions and reduced biases.

7 Conclusion

LLMs equipped with these advanced technologies demonstrate strong potential in causal reasoning, but further research is needed to fully realize and refine their capabilities. Addressing challenges such as managing intricate causal structures is essential for continued development. Enhancing interpretability and transparency through causal reasoning is crucial for building trust and helping users rely on model inferences. Exploring innovative methods will be key to overcoming current limitations and unlocking the full potential of LLMs in complex, real-world scenarios. Advancing these areas will make LLMs more reliable and effective tools for causal analysis across domains.

Limitations

In this survey, we review advancements in causal reasoning with LLMs, though several limitations should be noted. Our focus is primarily on text-based causal reasoning, with limited exploration of multi-modal or visual tasks, which restricts the generalizability of our findings across modalities. Although we evaluate mainstream LLMs and methods on popular benchmarks, we do not cover all enhancement techniques, particularly those requiring post-training or tailored to specific tasks and domains. Additionally, understanding the fundamental principles of causal reasoning in LLMs is crucial for ensuring trust and accountability in real-world applications, yet this area remains challenging and underexplored. Further research is warranted to address these gaps.

References

- S Abdali, A Parikh, S Lim, and E Kiciman. 2023. Extracting self-consistent causal insights from users feedback with llms and in-context learning. *arXiv preprint arXiv:2312.06820*.
- Anonymous. 2024. [A unified causal view of instruction tuning](#). *OpenReview preprint*.
- Pepa Atanasova, Oana-Maria Camburu, Christina Lioma, Thomas Lukasiewicz, Jakob Grue Simonsen, and Isabelle Augenstein. 2023. Faithfulness tests for natural language explanations. *arXiv preprint arXiv:2305.18029*.
- Abdolmahdi Bagheri, Matin Alinejad, Kevin Bello, and Alireza Akhondi-Asl. 2024. [C2p: Featuring large language models with causal reasoning](#). *arXiv preprint arXiv:2407.18069*.
- T Ban, L Chen, D Lyu, X Wang, and H Chen. 2023. Causal structure learning supervised by large language model. *arXiv preprint arXiv:2311.11689*.
- Guangsheng Bao et al. 2024. [Llms with chain-of-thought are non-causal reasoners](#). *arXiv preprint arXiv:2402.16048*.
- Jeroen Berrevoets, Krzysztof Kacprzyk, Zhaozhi Qian, and Mihaela van der Schaar. 2023. Causal deep learning. *arXiv preprint arXiv:2303.02186*.
- Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Scott Wen-tau Yih, and Yejin Choi. 2019. Abductive commonsense reasoning. *arXiv preprint arXiv:1908.05739*.
- Alexander Bondarenko, Magdalena Wolska, Stefan Heindorf, Lukas Blübaum, Axel-Cyrille Ngonga Ngomo, Benno Stein, Pavel Braslavski, Matthias Hagen, and Martin Potthast. 2022. Causalqa: a benchmark for causal question answering. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3296–3308.
- Emanuele Cavenaghi, Alessio Zanga, Fabio Stella, and Markus Zanker. 2024. [Towards a causal decision-making framework for recommender systems](#). *ACM Trans. Recomm. Syst.*, 2(2).
- Viviana Celli. 2022. Causal mediation analysis in economics: Objectives, assumptions, models. *Journal of Economic Surveys*, 36(1):214–234.
- Mandar S Chaudhary, Stephen Ranshous, and Nagiza F Samatova. 2018. A community-driven graph partitioning method for constraint-based causal discovery. In *Complex Networks & Their Applications VI: Proceedings of Complex Networks 2017 (The Sixth International Conference on Complex Networks and Their Applications)*, pages 253–264. Springer.
- Lu Chen, Wei Huang, Ruqing Zhang, et al. 2024. [A unified causal view of instruction tuning](#). *arXiv preprint arXiv:2402.06220*.
- Zhipeng Chen, Kun Zhou, et al. 2023. Chatcot: Tool-augmented chain-of-thought reasoning on chat-based large language models. *arXiv preprint arXiv:2305.14323*.
- David Maxwell Chickering. 2002. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554.
- Hyung Won Chung et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.
- Gregory F Cooper. 1997. A simple constraint-based algorithm for efficiently mining observational databases for causal relationships. *Data Mining and Knowledge Discovery*, 1:203–224.
- Antonia Creswell and Murray Shanahan. 2022. Faithful reasoning using large language models. *arXiv preprint arXiv:2208.14271*.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. 2022. Selection-inference: Exploiting large language models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*.
- Li Du, Xiao Ding, Kai Xiong, Ting Liu, and Bing Qin. 2022. e-care: a new dataset for exploring explainable causal reasoning. *arXiv preprint arXiv:2205.05849*.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*.

- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Amir Feder, Nadav Oved, Uri Shalit, and Roi Reichart. 2021. Causalm: Causal model explanation through counterfactual language models. *Computational Linguistics*, 47(2):333–386.
- Carlos Fernández-Loría and Foster Provost. 2022. Causal decision making and causal effect estimation are not the same. . . and why it matters. *INFORMS Journal on Data Science*, 1(1):4–16.
- Jörg Frohberg and Frank Binder. 2022. **CRASS: A novel data set and benchmark to test counterfactual reasoning of large language models**. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 2126–2140, Marseille, France. European Language Resources Association.
- K Gandhi, JP Fränken, T Gerstenberg, and ND Goodman. 2023. Understanding social reasoning in language models with language models. *arXiv preprint arXiv:2306.15448*.
- Yair Gat, Nitay Calderon, Amir Feder, et al. 2023. Faithful explanations of black-box nlp models using llm-generated counterfactuals. *arXiv preprint arXiv:2310.00603*.
- Clark Glymour, Kun Zhang, and Peter Spirtes. 2019. Review of causal discovery methods based on graphical models. *Frontiers in genetics*, 10:524.
- Sander Greenland, Judea Pearl, and James M. Robins. 1999. **Confounding and Collapsibility in Causal Inference**. *Statistical Science*, 14(1):29 – 46.
- Hosein Hasanbeig, Hiteshi Sharma, Leo Betthausser, Felipe Vieira Frujeri, and Ida Momennejad. 2023. Allure: auditing and improving llm-based evaluation of text using iterative in-context-learning. *arXiv e-prints*, pages arXiv–2309.
- Xuehai He, Diji Yang, Weixi Feng, Tsu-Jui Fu, Arjun Akula, Varun Jampani, Pradyumna Narayana, Sugato Basu, William Yang Wang, and Xin Eric Wang. 2022. Cpl: Counterfactual prompt learning for vision and language models. *arXiv preprint arXiv:2210.10362*.
- David Heckerman, Christopher Meek, and Gregory Cooper. 2006. A bayesian approach to causal discovery. *Innovations in Machine Learning: Theory and Applications*, pages 1–28.
- Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2010. **SemEval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals**. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 33–38, Uppsala, Sweden. Association for Computational Linguistics.
- Xudong Hong et al. 2023. Do large language models and humans have similar behaviors in causal inference with script knowledge? *arXiv preprint arXiv:2311.07311*.
- Biwei Huang, Kun Zhang, Yizhu Lin, Bernhard Schölkopf, and Clark Glymour. 2018. Generalized score functions for causal discovery. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1551–1560.
- Biwei Huang, Kun Zhang, Jiji Zhang, Joseph Ramsey, Ruben Sanchez-Romero, Clark Glymour, and Bernhard Schölkopf. 2020. **Causal discovery from heterogeneous/nonstationary data with independent changes**. *Preprint*, arXiv:1903.01672.
- Jiaxin Huang, Shixiang Shane Gu, Le Hou, et al. 2022. **Large language models can self-improve**. *arXiv preprint arXiv:2210.11610*.
- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Cosmos qa: Machine reading comprehension with contextual commonsense reasoning. *arXiv preprint arXiv:1909.00277*.
- Haitao Jiang, Lin Ge, Yuhe Gao, Jianian Wang, and Rui Song. 2023. Large language model for causal decision making. *arXiv preprint arXiv:2312.17122*.
- Zhijing Jin, Yuen Chen, Felix Leeb, Luigi Gresele, Ojasv Kamal, LYU Zhiheng, Kevin Blin, Fernando Gonzalez Adauto, Max Kleiman-Weiner, Mrinmaya Sachan, et al. 2023a. **Cladder: Assessing causal reasoning in language models**. In *Thirty-seventh conference on neural information processing systems*.
- Zhijing Jin, Jiarui Liu, Zhiheng Lyu, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona Diab, and Bernhard Schölkopf. 2023b. Can large language models infer causation from correlation? *arXiv preprint arXiv:2306.05836*.
- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 252–262.
- Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. 2023a. Causal reasoning and large language models: Opening a new frontier for causality. *arXiv preprint arXiv:2305.00050*.
- Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. 2023b. Causal reasoning and large language models: Opening a new frontier for causality. *arXiv preprint arXiv:2305.00050*.
- Mojtaba Komeili. 2021. Internet-augmented dialogue generation. *arXiv preprint arXiv:2107.07566*.

- Robert Kowalski. 1979. [Algorithm = logic + control](#). *Commun. ACM*, 22(7):424–436.
- Jack Kuipers and Giusi Moffa. 2022. The interventional bayesian gaussian equivalent score for bayesian causal inference with unknown soft interventions. *arXiv preprint arXiv:2205.02602*.
- Kai Lagemann, Christian Lagemann, Bernd Taschler, and Sach Mukherjee. 2023. Deep learning of causal structures in high dimensions under data limitations. *Nature Machine Intelligence*, 5(11):1306–1316.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*.
- Peiwen Li, Xin Wang, et al. 2024. Llm-enhanced causal discovery in temporal domain from interventional data. *arXiv preprint arXiv:2404.14786*.
- Zhongyang Li, Xiao Ding, Kuo Liao, Bing Qin, and Ting Liu. 2021. Causalbert: Injecting causal knowledge into pre-trained models with minimal supervision. *arXiv preprint arXiv:2107.09852*.
- Zhongyang Li, Xiao Ding, Ting Liu, J. Edward Hu, and Benjamin Van Durme. 2020. [Guided generation of cause and effect](#). In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages 3629–3636. International Joint Conferences on Artificial Intelligence Organization. Main track.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*.
- Jintao Liu, Zequn Zhang, Zhi Guo, Li Jin, Xiaoyu Li, Kaiwen Wei, and Xian Sun. 2023b. Kept: Knowledge enhanced prompt tuning for event causality identification. *Knowledge-Based Systems*, 259:110064.
- Xiao Liu, Zirui Wu, et al. 2024a. [Are llms capable of data-based statistical and causal reasoning? benchmarking advanced quantitative reasoning with data](#). *arXiv preprint arXiv:2402.17644*.
- Xiaoyu Liu, Paiheng Xu, Junda Wu, Jiaxin Yuan, Yifan Yang, Yuhang Zhou, Fuxiao Liu, Tianrui Guan, Haoliang Wang, Tong Yu, et al. 2024b. Large language models and causal inference in collaboration: A comprehensive survey. *arXiv preprint arXiv:2403.09606*.
- Xiaoyu Liu, Paiheng Xu, Junda Wu, et al. 2024c. [Large language models and causal inference in collaboration: A comprehensive survey](#). *arXiv preprint arXiv:2403.09606*.
- Y Lu, W Feng, W Zhu, W Xu, XE Wang, M Eckstein, and WY Wang. 2022. [Neuro-symbolic causal language planning with commonsense prompting](#). *arXiv e-prints*.
- Yujie Lu, Weixi Feng, Wanrong Zhu, Wenda Xu, Xin Eric Wang, Miguel Eckstein, and William Yang Wang. 2023. [Neuro-symbolic procedural planning with commonsense prompting](#). In *The Eleventh International Conference on Learning Representations*.
- Zhiyi Luo, Yuchen Sha, Kenny Q Zhu, Seung-won Hwang, and Zhongyuan Wang. 2016. Commonsense causal reasoning between short texts. In *Fifteenth international conference on the principles of knowledge representation and reasoning*.
- Syed S Mahmooda, Daniel Levy, RS Vasan, and TJ Wang. 2014. The framingham heart study and the epidemiology of cardiovascular diseases: A historical perspective. *Lancet*, 383(9921):999–1008.
- Alexander Marx and Jilles Vreeken. 2017. Telling cause from effect using mdl-based local and global regression. In *2017 IEEE international conference on data mining (ICDM)*, pages 307–316. IEEE.
- Valentyn Melnychuk, Dennis Frauen, and Stefan Feuerriegel. 2022. Causal transformer for estimating counterfactual outcomes. In *International Conference on Machine Learning*, pages 15293–15329. PMLR.
- Joris M Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. 2016a. Distinguishing cause from effect using observational data: methods and benchmarks. *Journal of Machine Learning Research*, 17(32):1–102.
- Joris M Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. 2016b. Distinguishing cause from effect using observational data: methods and benchmarks. *The Journal of Machine Learning Research*, 17(1):1103–1204.
- Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. [A corpus and cloze evaluation for deeper understanding of commonsense stories](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849, San Diego, California. Association for Computational Linguistics.
- Nasrin Mostafazadeh, Aditya Kalyanpur, Lori Moon, David Buchanan, Lauren Berkowitz, Or Biran, and Jennifer Chu-Carroll. 2020. Glucose: Generalized and contextualized story explanations. *arXiv preprint arXiv:2009.07758*.
- Allen Nie et al. 2023. Moca: Measuring human-language model alignment on causal and moral judgment tasks. *Advances in Neural Information Processing Systems*, 36:78360–78393.
- Qiang Ning, Zhili Feng, Hao Wu, and Dan Roth. 2019. Joint reasoning for temporal and causal relations. *arXiv preprint arXiv:1906.04941*.

- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Nick Pawlowski, James Vaughan, Joel Jennings, and Cheng Zhang. 2023. Answering causal questions with augmented llms.
- Judea Pearl. 2009. *Causality*. Cambridge university press.
- Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. 2016. *Causal inference in statistics: A primer*. John Wiley & Sons.
- Judea Pearl and Dana Mackenzie. 2018. *The book of why: the new science of cause and effect*. Basic books.
- Judea Pearl et al. 2000. Models, reasoning and inference. *Cambridge, UK: CambridgeUniversityPress*, 19(2):3.
- Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021. [True few-shot learning with language models](#). *Preprint*, arXiv:2105.11447.
- Ronan Perry, Julius Von Kügelgen, and Bernhard Schölkopf. 2022. Causal discovery in heterogeneous environments under the sparse mechanism shift hypothesis. *Advances in Neural Information Processing Systems*, 35:10904–10917.
- Jonas Peters and Peter Bühlmann. 2014. Identifiability of gaussian structural equation models with equal error variances. *Biometrika*, 101(1):219–228.
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. 2017. *Elements of causal inference: foundations and learning algorithms*. The MIT Press.
- Mattia Proserpi, Yi Guo, Matt Sperrin, James S Koopman, Jae S Min, Xing He, Shannan Rich, Mo Wang, Iain E Buchan, and Jiang Bian. 2020. Causal inference and counterfactual prediction in machine learning for actionable healthcare. *Nature Machine Intelligence*, 2(7):369–375.
- Lianhui Qin, Antoine Bosselut, Ari Holtzman, Chandra Bhagavatula, Elizabeth Clark, and Yejin Choi. 2019. Counterfactual story reasoning and generation. *arXiv preprint arXiv:1909.04076*.
- Joseph Ramsey and Bryan Andrews. 2018. Fask with interventional knowledge recovers edges from the sachs model. *arXiv preprint arXiv:1805.03108*.
- Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *2011 AAAI Spring Symposium Series*.
- Raanan Y. Rohekar, Yaniv Gurwicz, and Shami Nisimov. 2023. [Causal interpretation of self-attention in pre-trained transformers](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 31450–31465. Curran Associates, Inc.
- Angelika Romanou, Syrielle Montariol, Debjit Paul, Leo Laugier, Karl Aberer, and Antoine Bosselut. 2023. Crab: Assessing the strength of causal relationships between real-world events. *arXiv preprint arXiv:2311.04284*.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019a. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 3027–3035.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le-Bras, and Yejin Choi. 2019b. Socialliqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*.
- Timo Schick et al. 2024. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36.
- Ladan Shams and Ulrik Beierholm. 2022. Bayesian causal inference: A unifying neuroscience theory. *Neuroscience & Biobehavioral Reviews*, 137:104619.
- Kurt Shuster et al. 2022. Blenderbot 3: A deployed conversational agent that continually learns to responsibly engage. *arXiv preprint arXiv:2208.03188*.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. 31(1).
- Peter Spirtes, Clark N Glymour, and Richard Scheines. 2000. *Causation, prediction, and search*. MIT press.
- Peter Spirtes and Kun Zhang. 2016. Causal discovery and inference: concepts and recent methodological advances. In *Applied informatics*, volume 3, pages 1–28. SpringerOpen.
- Ieva Staliūnaitė, Philip John Gorinski, and Ignacio Iacobacci. 2021. Improving commonsense causal reasoning by adversarial training and data augmentation. *arXiv preprint arXiv:2101.04966*.
- Natasa Tagasovska, Valérie Chavez-Demoulin, and Thibault Vatter. 2020. [Distinguishing cause from effect using quantiles: Bivariate quantile causal discovery](#). In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 9311–9323. PMLR.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2018. Commonsenseqa: A question answering challenge targeting commonsense knowledge. *arXiv preprint arXiv:1811.00937*.
- Xin Tan, Bawei Zou, and Ai Ti Aw. 2023. Evidence-based interpretable open-domain fact-checking with large language models. *arXiv preprint arXiv:2312.05834*.

- Ziyi Tang, Ruilin Wang, Weixing Chen, Keze Wang, Yang Liu, Tianshui Chen, and Liang Lin. 2023. Towards causalgpt: A multi-agent approach for faithful knowledge reasoning via promoting causal consistency in llms. *arXiv preprint arXiv:2308.11914*.
- Michail Tsagris, Giorgos Borboudakis, Vincenzo Lagani, and Ioannis Tsamardinos. 2018. Constraint-based causal discovery with mixed data. *International journal of data science and analytics*, 6:19–30.
- Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. 2006. The max-min hill-climbing bayesian network structure learning algorithm. *Machine learning*, 65:31–78.
- Ruibo Tu, Kun Zhang, Bo Bertilson, Hedvig Kjellstrom, and Cheng Zhang. 2019. Neuropathic pain diagnosis simulator for causal discovery algorithm evaluation. *Advances in Neural Information Processing Systems*, 32.
- Vijay K Vemuri. 2015. Causal inference for statistics, social, and biomedical sciences: an introduction by guido w. imbens and donald b. rubin. *Taylor & Francis*.
- Guangya Wan, Yuqi Wu, Mengxuan Hu, Zhixuan Chu, and Sheng Li. 2024. Bridging causal discovery and large language models: A comprehensive survey of integrative approaches and future directions. *arXiv preprint arXiv:2402.11068*.
- Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, et al. 2023a. Survey on factuality in large language models: Knowledge, retrieval and domain-specificity. *arXiv preprint arXiv:2310.07521*.
- Hecong Wang, Erqian Xu, Pinxin Liu, Zijian Meng, and Zhen Bai. 2023b. Contextualized multi-step commonsense reasoning through context extension.
- Yuqing Wang and Yun Zhao. 2023. Tram: Benchmarking temporal reasoning for large language models. *arXiv preprint arXiv:2310.00835*.
- Zhaowei Wang, Quyet V Do, Hongming Zhang, Jiayao Zhang, Weiqi Wang, Tianqing Fang, Yangqiu Song, Ginny Y Wong, and Simon See. 2023c. Cola: Contextualized commonsense causal reasoning from the causal inference perspective. *arXiv preprint arXiv:2305.05191*.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, et al. 2022. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, et al. 2023. Chain-of-thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.
- Pengzhou Wu and Kenji Fukumizu. 2020. Causal mosaic: Cause-effect inference via nonlinear ica and ensemble method. In *International Conference on Artificial Intelligence and Statistics*, pages 1157–1167. PMLR.
- Yuxi Xie, Guanzhen Li, and Min-Yen Kan. 2023. Echo: A visio-linguistic dataset for event causality inference via human-centric reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4064–4085.
- Linying Yang, Vik Shirvaikar, Oscar Clivio, and Fabian Falck. 2024. A critical review of causal reasoning benchmarks for large language models. *arXiv preprint arXiv:2407.08029*.
- Sen Yang, Xin Li, Leyang Cui, Lidong Bing, and Wai Lam. 2023. Neuro-symbolic integration brings causal and reliable reasoning proofs. *arXiv preprint arXiv:2311.09802*.
- Xiaoyu Yang, Stephen Obadinma, Huasha Zhao, Qiong Zhang, Stan Matwin, and Xiaodan Zhu. 2020. SemEval-2020 task 5: Counterfactual recognition. In *Proceedings of the 14th International Workshop on Semantic Evaluation (SemEval-2020)*, Barcelona, Spain.
- Liuyi Yao, Zhixuan Chu, Sheng Li, Yaliang Li, Jing Gao, and Aidong Zhang. 2021. A survey on causal inference. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15(5):1–46.
- S Yao, J Zhao, D Yu, N Du, I Shafran, K Narasimhan, and Y Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.
- Murong Yue, Jie Zhao, Min Zhang, Liang Du, and Ziyu Yao. 2023. Large language model cascades with mixture of thoughts representations for cost-efficient reasoning. *arXiv preprint arXiv:2310.03094*.
- Chuanqi Zang, Hanqing Wang, Mingtao Pei, and Wei Liang. 2023. Discovering the real association: Multimodal causal reasoning in video question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19027–19036.
- Alessio Zanga, Elif Ozkirimli, and Fabio Stella. 2022. A survey on causal discovery: theory and practice. *International Journal of Approximate Reasoning*, 151:101–129.
- Matej Zečević, Moritz Willig, Devendra Singh Dhami, and Kristian Kersting. 2023. Causal parrots: Large language models may talk causality but are not causal. *arXiv preprint arXiv:2308.13067*.
- Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. Swag: A large-scale adversarial dataset for grounded commonsense inference. *arXiv preprint arXiv:1808.05326*.

- Chenyang Zhang, Haibo Tong, Bin Zhang, and Dongyu Zhang. 2024. [Probing causality manipulation of large language models](#). *Preprint*, arXiv:2408.14380.
- Kun Zhang and Aapo Hyvarinen. 2012. On the identifiability of the post-nonlinear causal model. *arXiv preprint arXiv:1205.2599*.
- Li Zhang, Hainiu Xu, Yue Yang, Shuyan Zhou, Weiqiu You, Manni Arora, and Chris Callison-Burch. 2023. Causal reasoning of entities and events in procedural texts. *arXiv preprint arXiv:2301.10896*.
- Junhao Zheng, Qianli Ma, Shengjie Qiu, Yue Wu, Peitian Ma, Junlong Liu, Huawen Feng, Xichen Shang, and Haibin Chen. 2023. Preserving commonsense knowledge from pre-trained language models via causal inference. *arXiv preprint arXiv:2306.10790*.
- Shengyu Zhu, Ignavier Ng, and Zhitang Chen. 2019. Causal discovery with reinforcement learning. *arXiv preprint arXiv:1906.04477*.

A Casual Reasoning Tasks

Causal discovery and inference (Glymour et al., 2019; Pearl et al., 2000; Spirtes et al., 2000; Spirtes and Zhang, 2016; Pearl et al., 2016; Perry et al., 2022; Yao et al., 2021) are two related fields within various domains such as statistics and machine learning that aim to understand the relationships between variables in a system and identify cause-and-effect connections. (1) *Causal discovery* (Glymour et al., 2019; Zhu et al., 2019; Heckerman et al., 2006; Huang et al., 2018) aims to identify causal relationships from data, especially when the causal structure is unclear. It involves constructing models, such as causal graphs, to elucidate cause-and-effect. This is crucial as real-world constraints often make controlled experiments to establish causality impractical, unethical, or impossible (Spirtes et al., 2000; Vemuri, 2015; Peters et al., 2017). (2) *Causal inference* (Pearl, 2009; Yao et al., 2021; Shams and Beierholm, 2022; Pearl et al., 2016; Mooij et al., 2016b; Kuipers and Moffa, 2022; Peters et al., 2017) involves adding data or interventions to a causal graph from causal discovery to test variable cause-effect relationships. It aims to determine the direct impact of a variable or intervention on another. The methods used for causal discovery may include constraint-based algorithms (Tsagris et al., 2018; Cooper, 1997; Chaudhary et al., 2018), score-based algorithms (Huang et al., 2018; Peters and Bühlmann, 2014; Chickering, 2002; Heckerman et al., 2006), hybrid algorithms (Tsamardinos et al., 2006) and so on (Zanga et al., 2022; Glymour et al., 2019). Constraint-based algorithms reconstruct causal graphs using conditional independence statements from statistical tests. Score-based methods optimize the fitness of the graph using score metrics from observed data. Hybrid approaches integrate both to enhance causal inference.

Although there has been much research based on traditional methods, they have several limitations in data requirements, model scalability, and model robustness that LLMs can address more effectively. (1) *Data Quality*: **Traditional Methods** require well-developed, high-quality datasets, making them vulnerable to missing data and measurement errors. They rely on assumptions like linearity, stationarity, and causal sufficiency, which may not hold in real-world scenarios. **LLMs** capture non-linear interactions and are suitable for various types of data. They can augment sparse or

incomplete data through in-context learning, generating plausible hypotheses, and filling gaps. (2) *Scalability*: **Traditional Methods** can be computationally intensive and less scalable with increasing variables, often requiring costly intervention data. **LLMs** efficiently process large datasets and automate causal discovery, making them more scalable and adaptable to complex datasets. (3) *Robustness*: **Traditional Methods** simulate counterfactuals and predict intervention effects without explicit data, but their reliance on strict assumptions can limit robustness in practical applications. **LLMs** apply methods such as fine-tuning to dynamically integrate vast background knowledge, providing a richer context for causal reasoning and making them more robust to data quality issues.

B Enhancing Traditional CR Methods

B.1 Causal Information Extraction

LLMs can directly support causal reasoning by serving as knowledge bases or extracting causal information. The effectiveness of problem-solving algorithms can be improved by adopting a modular approach that separates the logic and control components (Kowalski, 1979). The logic component represents the essential knowledge for solving the problem, while the control component manages the problem-solving strategy. Decoupling these components enhances flexibility, allowing updates to the control strategy without altering the underlying knowledge base.

Yang et al. (2023) proposed a neuro-symbolic approach with two components: an LLM-based symbolic representation generator that translates natural language into formal knowledge for symbolic inference, and an LLM-free symbolic inference engine. Similarly, Abdali et al. (2023) leveraged LLMs' reasoning capabilities to extract causal variables and event sequences, providing a heuristic for assessing feedback informativeness in the absence of domain knowledge. Additionally, RealTCD (Li et al., 2024) introduced LLM-guided meta-initialization to extract meta-knowledge from textual information, improving the discovery process in real industrial settings.

Another notable contribution is the iterative LLM-supervised causal structure learning (ILS-CSL) framework introduced by Ban et al. (2023), which iteratively refines causal Directed Acyclic Graphs through feedback from LLMs. This framework highlights the efficiency and robustness

of LLMs in improving causal discovery. More broadly, the ability of LLMs to extract causal information can reduce the dependency on purely data-driven methods and unlock new opportunities for hybrid CR approaches. As LLMs evolve, they are likely to play an increasingly important role in causal discovery and inference by connecting unstructured text with structured causal models, making CR more robust and adaptable across domains.

B.2 Causality Data Generation & Augmentation

LLMs store a vast amount of knowledge within their weights, and several works have explored how to extract causal-related knowledge to support traditional causal reasoning models. For instance, Li et al. (2020) introduced a method where, given an input cause sentence, a causal subgraph is retrieved from an event causality network, and the most probable effect event is selected to guide the generation of an effect sentence. As LLMs scale, their capacity to generate causal alternatives has improved. For example, Staliūnaitė et al. (2021) used GPT-2 to generate alternatives and re-write filtered augmentation data for causal relations. Similarly, Gat et al. (2023) showcased how LLMs can generate counterfactual approximations efficiently, focusing on model-agnostic explanations that enhance the interpretability and faithfulness of CR methods.

However, despite the impressive knowledge stored within LLMs, they remain prone to hallucinations, often generating incorrect or misleading facts (Wang et al., 2023a). This issue is particularly problematic in causal reasoning, where precision and factual correctness are essential for drawing valid inferences and supporting decision-making processes. Hallucinated outputs can lead to incorrect causal links, undermining the reliability of LLM-generated causal insights. Moreover, the lack of transparency in how LLMs generate such outputs further complicates efforts to identify and correct errors, posing challenges to researchers and practitioners alike.

Addressing these challenges requires strategies to ensure more controllable and precise causal data generation. High-quality causal data is essential not only for improving LLM accuracy but also for fostering trust in their outputs, particularly in real-world applications like healthcare, economics, and public policy. Future research should focus on robust methods to mitigate hallucinations and enhance the controllability of LLM-generated causal

information, bridging the gap between language model capabilities and practical causality needs.

C Tasks & Benchmarks Overview

We first categorize the various end tasks into three groups: causal discovery, causal inference, and additional causal tasks. Causal discovery involves identifying potential causal relationships between variables within large, often unannotated, datasets. This is crucial for understanding causal dynamics in complex systems. Causal inference refers to quantifying the impact of a treatment (or intervention) on an outcome, determining how changes in one variable causally influence another, typically using observational data. Beyond these two primary tasks, there are additional causal tasks such as causal attribution, counterfactual reasoning, and causal explanation. These tasks extend the scope of causal analysis, often integrating causal reasoning with other reasoning tasks.

We show the statistics of existing benchmarks for causal reasoning tasks (Table 4) with those used in our experiments in bold. We also provide their inclusion of the three levels of causality: association, intervention, and counterfactual. From this data, we draw several key insights: (1) *Opportunity for Improvement*: There is a need for more datasets incorporating **intervention** and **counterfactual** reasoning to comprehensively evaluate causal reasoning. (2) *Balanced Data*: Greater emphasis on **real-world** scenarios and underrepresented domains, such as health and social contexts, would enhance the robustness of causal benchmarks. (3) *Scaling Challenges*: Although large datasets exist, developing datasets for complex causal tasks such as interventions remains difficult, limiting their representation. (4) *Future Work*: Advancing causal reasoning in LLMs will require benchmarks that balance real and synthetic data, cover all three levels of causality, and span diverse domains. In addition, we show examples from each dataset used in our experiments in Figure 9.

D Implementation Details

We begin by discussing the evaluation metrics used across different benchmarks. TRAM uses accuracy to assess causal reasoning tasks. NPDS evaluates performance with both F1 score and causal accuracy, while Corr2Cause employs F1 as the primary metric, alongside precision, recall, and accuracy. CRAB focuses on Pairwise and Graded

Dataset	Size (Unit)	Domain	Source	Association	Intervention	Counterfactual
Causality Discovery						
SemEval Tasks 8 (2010)	1.3K (P)	Mixed	Real	✓		
COPA (2011)	1K (Q)	Dailylife	Real	✓		
CEPairs (2016a)	108 (P)	Mixed	Real	✓		
CausalNet (2016)	62M (P)	Dailylife	Synthetic	✓		
MultiRC (2018)	9.8K (Q)	Mixed	Real	✓		
NPDS (2019)	770 (R)	Health	Synthetic	✓		
TCR Benchmark (2019)	25 (D)	News	Real	✓		
SemEval Tasks 5 (2020)	20K (S)	Mixed	Real	✓		✓
CausalBank (2020)	314M (P)	Mixed	Synthetic	✓		
GLUCOSE (2020)	4.8K (S)	Mixed	Real	✓		
CausaLM (2021)	9.8K (S)	Ethnic & Gender	Real	✓		✓
e-CARE (2022)	21K (Q)	Mixed	Real	✓		✓
CausalQA (2022)	1.1M (Q)	Mixed	Real	✓		
Corr2Cause (2023b)	200K (S)	Mixed	Synthetic	✓		
Causal Inference						
ROC Stories (2016)	100K (S)	Mixed	Real	✓		
CSR (2019)	16K (S)	Mixed	Real	✓	✓	✓
ART (2019)	20K (S)	Mixed	Real	✓		
CausalQA (2022)	1.1M (Q)	Mixed	Real	✓		
COPEs (2023c)	1.3K (P)	Mixed	Real	✓	✓	
CLADDER (2023a)	10K (S)	Mixed	Synthetic	✓	✓	✓
CRAB (2023)	2.7K (P)	Dailylife	Real	✓		
Additional Causal Tasks						
SWAG (2018)	113K (Q)	Mixed	Real	✓		
CommonsenseQA (2018)	12K (Q)	Mixed	Real	✓		
ATOMIC (2019a)	24K (S)	Mixed	Real	✓		
Cosmos QA (2019)	35K (Q)	Mixed	Real	✓		✓
SocialIQA (2019b)	38K (Q)	Social	Real	✓		
CRASS (2022)	275 (Q)	Mixed	Real	✓		✓
MoCa (2023)	144 (S)	Mixed	Real	✓		
Tram (2023)	1.2K (Q)	Dailylife	Real	✓		

Table 4: Benchmarks for causal reasoning tasks. Datasets used in our experiments are in bold. Sizes are reported in various units: P (causal pairs), R (causal relations), S (samples), Q (questions), and D (documents). We also provide their inclusion of the three levels of causality: association, intervention, and counterfactual.

Causality tasks, using Macro F1 for its 4-class Pairwise Causality Multiple Choice Question (MCQ) task. MOCA compares LLM responses with human participants using agreement (Agg), Area Under the Curve (AUC), Mean Absolute Error (MAE), and Cross-Entropy (CE), though we directly assess LLMs on causal tasks without relying on agreement metrics. CRASS uses top-1 accuracy, while e-CARE assesses MCQs with accuracy, average-BLEU (n=4), ROUGE-L, Perplexity, and Causal Explanation Quality (CEQ) for explanations. Finally, COPA uses accuracy as the primary metric.

Across these benchmarks, metrics such as F1, precision, recall, AUC, and BLEU assess different aspects of performance, including classification quality, alignment with human responses, and explanation quality. We select pass@1 accuracy as our primary metric because it is straightforward and widely used for evaluating multiple-choice and

prediction-based causal tasks.

For evaluation, we randomly select 500 samples from each dataset’s test set using a random seed of 42, or the maximum available size if fewer samples exist. Due to cost and rate constraints, models like GPT-4o and Claude-3.5 are tested on 100 samples. We utilize a unified prompt pattern (Figure 9) comprising three parts: (1) **Identity Recognition**, such as "You are a highly intelligent question-answering bot with profound knowledge of causal inference"; (2) **Question-Context**, which describes the specific question to be addressed; and (3) **Answer Format Control**, which guides the model’s response. The latter instructs the model to either provide a direct answer with "Give me the answer directly" for Direct I/O or to begin with reasoning or evidence, concluding with the final result marked by "#####" for CoT reasoning, ensuring the output is easy to parse. In the few-shot learning setup, we

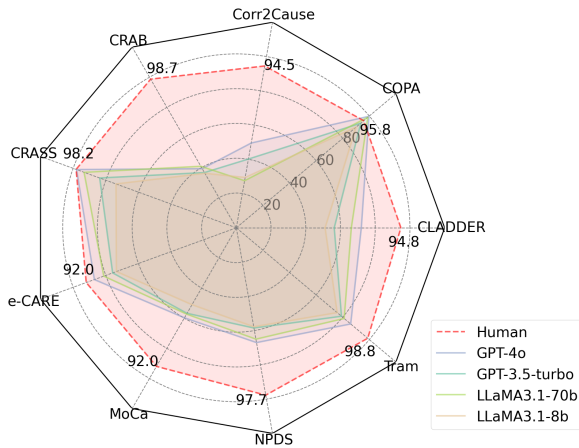


Figure 8: Overview of human and LLM performance on causal reasoning tasks, with annotated human scores.

use two fixed examples for most datasets and one for the CRAB dataset, which has a longer context. Direct I/O prompts request direct answers, while CoT prompts involve generating CoT examples with GPT-4o (Figure 10), followed by manual validation for accuracy. Pass @1 Accuracy is calculated using a string parsing function to extract answers from outputs tagged with "Answer" or "#####" Any model failing to follow the specified format is treated as a failure case.

We evaluate a diverse set of language models, including Mistral AI’s instructional models, Google’s Gemma models, Meta’s LLaMA 3 and 3.1 series, DeepSeek’s V3 and R1 models, OpenAI’s GPT-3.5-turbo, GPT-4o, o1 and o3-mini, and Anthropic’s Claude-3.5, ensuring comprehensive coverage of various architectures and capabilities. For open-source models, we used the API provided by the Replicate and Together platform, while proprietary models were accessed through their respective company APIs. The specific model versions used in our experiments are listed below:

- mistralai/Mistral-7B-Instruct-v0.3
- mistralai/Mixtral-8x7B-Instruct-v0.1
- google/gemma-2-9b-it
- google/gemma-2-27b-it
- meta-llama/Meta-Llama-3-8B-Instruct
- meta-llama/Meta-Llama-3-70B-Instruct-Turbo
- meta-llama/Meta-Llama-3.1-8B-Instruct

- meta-llama/Meta-Llama-3.1-70B-Instruct-Turbo
- Deepseek-V3
- Deepseek-R1
- gpt-3.5-turbo-0125
- gpt-4o-2024-08-06
- claude-3-5-sonnet-20240620
- o1-2024-12-17
- o3-mini-2025-01-31

To assess human performance on each dataset, we first reviewed the original papers for the datasets and found that some (Du et al., 2022; Froberg and Binder, 2022; Wang and Zhao, 2023) already provided human evaluation results. For the remaining datasets, we asked three graduate students to solve 100 randomly selected problems from each dataset. Prior to the evaluation, the participants underwent training to develop a strong foundation in causal reasoning. They then solved the problems independently, and their average accuracy is calculated as the final result.

For our error analysis, we first categorize errors into four primary types. **Statistical errors** (SE) arise from issues in data collection, measurement, or analysis, distorting relationships between variables. Sub-classes include: *confounding*, where an extraneous variable influences both independent and dependent variables, creating misleading associations; *selection bias*, where the sample is unrepresentative of the target population, leading to biased estimates; and *measurement error*, where inaccurate data collection results in flawed assessments. **Logical errors** (LE) involve flaws in reasoning, such as *directionality errors* (misidentifying causal direction), *inconsistency* (contradictory reasoning), and *over-generalization* (drawing broad conclusions from specific cases without proper support). **Probabilistic errors** (PE) stem from incorrect handling of probabilities and statistical inferences, including *Bayesian inference errors* (misapplication of priors or updates) and *probability comparison errors* (misinterpreting conditional probabilities or neglecting base rates). Finally, **contextual errors** (CE) occur when subtle details, such as shifts in tone or minor actions, are overlooked, leading to flawed causal inferences. We select three representative models (LLaMA3-8B, LLaMA3-70B, and

GPT-4o) for error analysis. For each dataset, 100 incorrect test samples are randomly chosen and manually analyzed. The results and analysis are presented in Section 4.3.

Detailed Error Analysis This appendix provides detailed examples and analysis of the error types discussed in Section 4.3. For each type, we present representative cases that illustrate common failure modes.

- Logical Errors - Case Study

Context: Billy and Suzy inherited an unusual type of hybrid car that has two special car batteries called Bartlett batteries. The car won't start unless it has at least one Bartlett battery. Since all the car needed was at least one Bartlett battery, they were able to start the car.

Question: Did Billy cause the car to start?

Model's Answer: A delightful question! Let's break down the stages: Billy brought a battery, Suzy also brought a battery, and either battery alone would start the car. Therefore, Billy's action was sufficient to cause the car to start. Answer: Yes

Error Analysis: The model exhibits several reasoning flaws. First, it conflates sufficiency with causation - while Billy's action was sufficient, it wasn't necessary. Second, it overlooks redundant causation by failing to recognize that Suzy's battery would have started the car regardless. Third, it relies on temporal sequence, implicitly assuming first action implies causation.

- Statistical Errors - Case Study

Context: A study tracked 1000 patients' recovery rates. In the data, 70% of patients who exercised regularly recovered. Among those who didn't exercise, 60% still recovered. The exercise group had more younger patients (avg age 45) than non-exercise group (avg age 65).

Question: Does exercise improve recovery rates?

Model's Answer: Let's analyze the data: The exercise group shows a 70% recovery rate while the non-exercise group shows a 60% recovery rate. There's a 10% improvement in the exercise group. Answer: Yes, exercise clearly improves recovery rates.

Error Analysis: The model makes several statistical reasoning errors. It fails to account for selection bias in group composition and ignores age as a potential confounding variable. The model draws causal conclusions from raw correlations without controlling for population differences, missing the

need for age-stratified analysis or matched comparisons.

- Probabilistic Errors - Case Study

Context: Imagine a self-contained, hypothetical world with only the following probabilities. Is the chance of kwox smaller when observing yupt?

Model's Answer: Let's calculate step by step: $P(kwox|yupt) = 0.47 * 0.57 + 0.65 * 0.43 = 0.5141$. Since $P(kwox) \approx 0.5141 < 0.5$, Answer: No.

Error Analysis: The model demonstrates several probabilistic reasoning failures. It incorrectly applies the conditional probability formula and shows confusion between joint and conditional probabilities. Most critically, it makes a comparison with the wrong baseline probability, leading to an invalid conclusion about relative likelihoods.

- Contextual Errors - Case Study

Context: In a small company with high employee turnover, all employees who left in the last quarter had completed their major projects. The company conducted exit interviews where employees cited various reasons for leaving including better opportunities, relocation, and personal reasons.

Question: Does project completion cause employee turnover?

Model's Answer: Looking at the data: All departing employees had completed their projects, and there is a perfect correlation between project completion and leaving. Answer: Yes, project completion appears to cause employee turnover.

Error Analysis: The model shows several contextual understanding failures. It ignores the broader context of employee career paths and job market dynamics, while failing to consider natural project completion as part of the normal work cycle. The model overlooks explicitly stated alternative reasons for leaving in the exit interviews. Furthermore, it misses crucial temporal aspects, namely that employees typically complete projects before leaving, regardless of their reason for departure.

E Comparison with Traditional Causal Reasoning Methods

While our primary focus is on LLM-based approaches, we present comparison results with traditional causal reasoning methods using the Tübingen cause-effect pairs dataset, a well-established benchmark for causal direction identification, from [Kıcıman et al. \(2023a\)](#). As shown in Table 5, traditional covariance-based methods like Slope ([Marx](#)

Model	Acc.
Slope (Marx and Vreeken, 2017)	0.75
bQCD (Tagasovska et al., 2020)	0.68
PNL-MLP (Zhang and Hyvarinen, 2012)	0.75
Mosaic (Wu and Fukumizu, 2020)	0.83
GPT-3.5-turbo	0.81
GPT-3.5-turbo (causal agent)	0.86
GPT-3.5-turbo (single prompt)	0.89
GPT-4 (single prompt)	0.96

Table 5: Performance comparison between traditional causal discovery methods and LLM-based approaches on the Tübingen cause-effect pairs benchmark. Traditional methods rely on statistical analysis of observational data, while LLM-based approaches leverage natural language understanding.

and Vreeken, 2017), bQCD (Tagasovska et al., 2020), and PNL-MLP (Zhang and Hyvarinen, 2012) achieve accuracies ranging from 0.68 to 0.75, while Mosaic (Wu and Fukumizu, 2020) reaches 0.83. In comparison, Kıcıman et al. (2023a)’s prompting strategies with LLMs demonstrate progressively improving performance: from 0.81 with basic GPT-3.5-turbo (without any special prompting), improving to 0.86 with causal agent prompting (using the system message "You are a helpful assistant for causal reasoning"), and further reaching 0.89 with single prompt engineering (reformulating the task as a direct question about causal direction with step-by-step reasoning requirements). Their experiments with GPT-4 using the same single prompt strategy ultimately achieved 0.96 accuracy. These results suggest that well-designed prompting techniques can effectively leverage LLMs’ inherent understanding of causal relationships, outperforming traditional statistical approaches on this benchmark dataset.

CLADDER

Question: {Imagine a self-contained, hypothetical world with only the following conditions, and without any unmentioned factors or causal relationships: Full moon has a direct effect on sprinkler and weather. Weather has a direct effect on ground. Sprinkler has a direct effect on ground. For situations where there is no full moon, the probability of wet ground is 38%. For situations where there is a full moon, the probability of wet ground is 43%. Does full moon negatively affect ground through weather and sprinkler?}

[Direct I/O] Give me the answer directly. Answer: Yes or Answer: No. Do not use any other format.

[CoT] Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'. The answer format is ####YES or ####NO.

COPA

Premise: {The man's voice sounded hoarse.}

Question: What was the { What was the cause?}?

A. {He had a cold.}

B. {He quit smoking.}

[Direct I/O] Give me the answer directly. Answer: A or Answer: B. Do not use any other format.

[CoT] Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'. The answer format is ####A or ####B.

CRAB

Context:

Article A: {International pressure may be the only way to stop the Brazilian government from taking a 201csuicide201d path in the Amazon, one of the country2019s most respected scientists has said, as the world2019s biggest rainforest continues...}

Article B: {Hollywood star Leonardo DiCaprio has announced a \$5 million pledge to help reduce the number of fires in the Amazon rainforest. It comes amid global outcry over fires in the Brazilian Amazon, with celebrities calling for action and pledging to donate to efforts in order to tackle the issue.DiCaprio announced on Sunday that Earth Alliance, an organization dedicated to environmental philanthropy that he co-chairs, was pledging \$5 million...}

Event 1: {The EU has recently signed a trade deal with South America.}

Event 2: {Leonardo DiCaprio pledges \$5 million to help reduce the number of fires in the Amazon rainforest.}

How much did event 1 cause event 2 to happen?

[A] High causality: Event 1 is definitely responsible for Event 2.

[B] Medium causality: Event 1 might have been responsible for Event 2.

[C] Low causality: The context gives a little indication that there is a connection between the two events, but background info might suggest a low causal connection.

[D] No causality: Events are somehow related but definitely NOT causally related.

[Direct I/O] Answer: A, or Answer: B, or Answer: C, or Answer: D. Do not use any other format.

[CoT] Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'. The answer format is ####A or ####B or ####C or ####D.

CRASS

Question: {A woman opens a treasure chest. What would have happened if the woman had not opened the treasure chest?}

Options:

A){The treasure chest would have been open.},

B){That is not possible.},

C){The treasure chest would have remained closed.},

D){I don't know.}

[Direct I/O] Answer: A, or Answer: B, or Answer: C, or Answer: D. Do not use any other format.

[CoT] Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'. The answer format is [A] or [B] or [C] or [D].

E-CARE

Premise:{Tom wants to go to space.}

A. {He aspired to be an astronaut.}

B. {He aspired to be a calligrapher.}

Which of the following is more likely to be true?

[Direct I/O] Answer: A or Answer: B. Do not use any other format.

[CoT] Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'. The answer format is ####A or ####B.

MOCA
<p>Story: {Sara has an electric guitar that she plays through an amplifier. The guitar will play through the amplifier if either the amplifier is set to channel A or the power mixer is turned on. Today, the amplifier is set to channel A. Sara checks to see if the power mixer is on. She sees that the power mixer is turned on, and she does not turn off the power mixer or adjust it at all. Because the guitar will play through the amplifier if either the amplifier is set to channel A or the power mixer is turned on, Sara's guitar plays through the amplifier.}</p> <p>Question: {Did the guitar play through the amplifier because Sara did not turn off the power mixer?}</p> <p>[Direct I/O] Answer: Yes or Answer: No. Do not use any other format.</p> <p>[CoT] Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'. The answer format is ####YES or ####NO.</p>
NPDS
<p>Question: {L Achillodynι causes R Br00f6stbesv00e4r. ("R" and "L" refer to the right and left sides of the body, respectively). Answer with true or false.}</p> <p>[Direct I/O] Answer: True or Answer: False. Do not use any other format.</p> <p>[CoT] Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'. The answer format is ####TRUE or ####FALSE.</p>
TRAM
<p>Premise:{As evening approaches, households engage in their routine culinary activities.}</p> <p>Question:{What's the more plausible RESULT?}</p> <p>A. {People celebrated the start of the new year.}</p> <p>B. {People prepared their dinner.}</p> <p>[Direct I/O] Answer: A or Answer: B. Do not use any other format.</p> <p>[CoT] Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'. The answer format is ####A or ####B.</p>
CORR2CAUSE
<p>Question: {Suppose there is a closed system of 5 variables, A, B, C, D and E. All the statistical relations among these 5 variables are as follows: A correlates with B. A correlates with C. A correlates with D. A correlates with E. B correlates with C. B correlates with D. B correlates with E. C correlates with D. C correlates with E. D correlates with E. However, C and E are independent given A and B. C and E are independent given A, B and D. D and E are independent given A and B. D and E are independent given A, B and C. Determine the truth value the following statement: A influences B through some mediator(s).}</p> <p>Determine the truth value the following statement: {A influences B through some mediator(s).}</p> <p>[Direct I/O] Answer: neutral or Answer: contradiction or Answer: entailment. Do not use any other format.</p> <p>[CoT] Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'. The answer format is ####NEUTRAL or ####CONTRADICTION or ####ENTAILMENT.</p>

Figure 9: Prompts for evaluating LLMs on different casual reasoning tasks.

COPA COT EXAMPLE

Premise: The man's voice sounded hoarse.

Question: What was the cause?

- A. He had a cold.
- B. He quit smoking.

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'.

A hoarse voice is often a symptom of a cold, as it can result from inflammation of the vocal cords caused by infection. Quitting smoking, while beneficial for overall health, generally leads to improved vocal quality over time rather than immediate hoarseness.

A

NPDS COT EXAMPLE

Question: L Achillodynyn causes R Bröstbesvär. 'R' and 'L' refer to the right and left sides of the body, respectively. Answer with true or false.

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'.

Achillodynia refers to pain in the Achilles tendon, while bröstbesvär refers to chest pain or chest-related issues. These are two distinct anatomical regions: the Achilles tendon is located in the lower leg, and chest pain pertains to the upper torso. There is no direct physiological connection between Achilles tendon issues and chest pain. Therefore, the claim that left-sided Achillodynia causes right-sided Bröstbesvär lacks anatomical or medical evidence.

False

TRAM COT EXAMPLE

As evening approaches, households engage in their routine culinary activities.

What's the more plausible RESULT?

- A. People celebrated the start of the new year.
- B. People prepared their dinner.

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'.

As evening approaches, households typically engage in their routine activities, which often include preparing meals. This situation aligns more naturally with the daily routine of preparing dinner rather than an annual event like celebrating the new year. Celebrating the new year is a rare and specific occasion, while preparing dinner is a common, everyday activity. The evidence suggests that dinner preparation is a more plausible result in this context.

B

E-CARE COT EXAMPLE

Tom wants to go to space.

- A. He aspired to be an astronaut.
- B. He aspired to be a calligrapher.

Which of the following is more likely to be true?

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'.

Since Tom wants to go to space, the likelihood of this aspiration aligning with a related field or interest is greater. Aspiring to be an astronaut is directly connected to space exploration, and people who want to go to space typically have a background or interest in relevant fields such as science, engineering, or aeronautics.

Aspiring to be a calligrapher, while a valuable skill, does not naturally correlate with a desire to go to space, as it is unrelated to space exploration or the skill sets typically required for such an endeavor.

A

CLADDER COT EXAMPLE

Question: Imagine a self-contained, hypothetical world with only the following conditions, and without any unmentioned factors or causal relationships: Full moon has a direct effect on sprinkler and weather. Weather has a direct effect on ground. Sprinkler has a direct effect on ground. For situations where there is no full moon, the probability of wet ground is 38%. For situations where there is a full moon, the probability of wet ground is 43%. Does full moon negatively affect ground through weather and sprinkler?

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'.

To determine whether the full moon negatively affects the ground through weather and sprinkler, we need to consider the chain of causal relationships in this hypothetical world. The full moon has a direct effect on both the sprinkler and the weather, which in turn directly affect the ground.

Here are the key details:

- When there is ****no full moon****, the probability of wet ground is 38%.
- When there ****is a full moon****, the probability of wet ground is 43%.

If the full moon negatively affected the ground through the weather and sprinkler, we would expect the probability of wet ground to decrease in the presence of a full moon. However, the probability of wet ground increases (from 38% to 43%) when there is a full moon.

Thus, based on this evidence, the full moon does not negatively affect the ground through the weather and sprinkler. Instead, it seems to have a positive or neutral effect.

No

CRASS COT EXAMPLE

Question: A girl reveals a thing. What would have happened if the girl had concealed the thing?

- A) The thing would have been revealed.
- B) That is not possible.
- C) Everything would have been fine.
- D) The thing would have not been visible.

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'.

To answer this question from a causal inference perspective, we must examine the hypothetical scenario where the girl did not reveal the thing (i.e., she concealed it). The question is asking us to compare the actual world (where the thing was revealed) to the counterfactual world (where the thing would have been concealed).

- **Option A:** "The thing would have been revealed" directly contradicts the hypothetical situation because the girl concealing the thing would prevent it from being revealed.
- **Option B:** "That is not possible" might refer to the idea that the girl cannot conceal the thing in some absolute sense, but no evidence suggests that such a scenario is impossible. Thus, this is an unlikely answer.
- **Option C:** "Everything would have been fine" introduces an assumption about the broader consequences of the concealment, which is speculative and not necessarily implied by the question.
- **Option D:** "The thing would have not been visible" aligns with the hypothetical situation where concealment prevents the thing from being revealed. This directly follows from the premise of the question.

Based on this reasoning, the most plausible answer is that if the girl had concealed the thing, it would not have been visible.

D

MOCA COT EXAMPLE

Story: Sara has an electric guitar that she plays through an amplifier. The guitar will play through the amplifier if either the amplifier is set to channel A or the power mixer is turned on. Today, the amplifier is set to channel A. Sara checks to see if the power mixer is on. She sees that the power mixer is turned on, and she does not turn off the power mixer or adjust it at all. Because the guitar will play through the amplifier if either the amplifier is set to channel A or the power mixer is turned on, Sara's guitar plays through the amplifier.

Question: Did the guitar play through the amplifier because Sara did not turn off the power mixer?

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'.

The guitar will play through the amplifier if either the amplifier is set to channel A or the power mixer is turned on. In this case, we know that the amplifier is already set to channel A, which alone is sufficient for the guitar to play through the amplifier, regardless of the power mixer's status.

Even though Sara did not turn off the power mixer, the amplifier being set to channel A was already enough for the guitar to work. Therefore, the guitar would have played through the amplifier even if the power mixer had been off or turned on. Sara's action of not turning off the power mixer had no effect on the outcome.

No

CORR2CAUSE COT EXAMPLE

Question: Suppose there is a closed system of 5 variables, A, B, C, D, and E. All the statistical relations among these 5 variables are as follows: A correlates with B. A correlates with C. A correlates with D. A correlates with E. B correlates with C. B correlates with D. B correlates with E. C correlates with D. C correlates with E. D correlates with E. However, C and E are independent given A and B. C and E are independent given A, B, and D. D and E are independent given A and B. D and E are independent given A, B, and C.

Determine the truth value of the following statement: A influences B through some mediator(s).

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'.

The statement asks whether A influences B through some mediator(s). Let's reason through this based on the given information about correlations and conditional independencies.

1. Correlation Relationships: A correlates with B, C, D, and E. Similarly, B correlates with C, D, and E. This suggests that A and B are connected, but it doesn't yet clarify whether the relationship is direct or mediated.

2. Conditional Independencies:

- C and E are independent given A and B, and also independent given A, B, and D. This tells us that once we know A and B, C and E do not provide additional information about each other.

- D and E are independent given A and B, and also independent given A, B, and C. This implies that once A and B are known, D and E do not add any further influence over each other.

These independencies suggest that the relationships between variables are conditional, possibly implying a complex interaction but without directly stating that a mediator is necessary between A and B.

3. No Direct Evidence of Mediation: The problem specifies conditional independencies between some variables given others, but there is no clear indication of a mediating variable between A and B. The provided conditions do not imply a mediation relationship but rather highlight correlations and conditional independencies.

Given that no mediating variable is directly implicated between A and B from the given information, the truth value of the statement "A influences B through some mediator(s)" cannot be definitively concluded as true (entailment) or false (contradiction). Instead, the information is insufficient to prove or disprove mediation.

neutral

CRAB COT EXAMPLE

Context:

Article A:

International pressure may be the only way to stop the Brazilian government from taking a “suicide” path in the Amazon, one of the country’s most respected scientists has said, as the world’s biggest rainforest continues to be ravaged by thousands of deliberate fires.

The large number of conflagrations – set illegally to clear and prepare land for crops, cattle and property speculation – has prompted the state of Amazonas to declare an emergency, created giant smoke clouds that have drifted hundreds of miles, and sparked international concerns about the destruction of an essential carbon sink.

“Our house is burning,” tweeted the French president, Emmanuel Macron, who called for emergency talks on the subject at this week’s G7 summit. But the response to the crisis has been mixed: while Norway and Germany have halted donations to the Brazilian government’s Amazon fund, the EU has recently signed a trade deal with South America, and the UK spent this week focusing on post-Brexit business with Brazil.

Scientists say the ongoing destruction will have dire consequences for Brazil and the world.

Article B:

Hollywood star Leonardo DiCaprio has announced a \$5 million pledge to help reduce the number of fires in the Amazon rainforest.

It comes amid global outcry over fires in the Brazilian Amazon, with celebrities calling for action and pledging to donate to efforts in order to tackle the issue.

DiCaprio announced on Sunday that Earth Alliance, an organization dedicated to environmental philanthropy that he co-chairs, was pledging \$5 million in funding to protect the Amazon.

“The governments around the world, including Brazil, need to work together to make sure this doesn’t continue,” he said.

Event 1: The EU has recently signed a trade deal with South America.

Event 2: Leonardo DiCaprio pledges \$5 million to help reduce the number of fires in the Amazon rainforest.

How much did event 1 cause event 2 to happen?

[A] High causality: Event 1 is definitely responsible for Event 2.

[B] Medium causality: Event 1 might have been responsible for Event 2.

[C] Low causality: The context gives a little indication that there is a connection between the two events, but background info might suggest a low causal connection.

[D] No causality: Events are somehow related but definitely NOT causally related.

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by '####'.

Event 1, the EU signing a trade deal with South America, and Event 2, Leonardo DiCaprio pledging \$5 million to help reduce fires in the Amazon, are both related to the broader context of the Amazon rainforest fires and global reactions. However, the two events appear to stem from different sources of influence and motivation.

Event 1 focuses on trade and economic relations between the EU and South America, whereas Event 2 involves an individual actor (DiCaprio) taking philanthropic action. The trade deal seems to be driven by economic priorities, while DiCaprio’s pledge is motivated by environmental concerns and public outcry. There is no clear evidence or connection that suggests the trade deal directly caused DiCaprio to pledge his donation. The timing of both events might be coincidental, as they both respond to the ongoing Amazon crisis, but they do not seem to be causally related.

Thus, the causality between the two events is likely very low, if not non-existent.

D

Figure 10: The CoT examples generated by GPT-4o.

CASUAL COT EXAMPLE

Question: Imagine a self-contained, hypothetical world with only the following conditions, and without any unmentioned factors or causal relationships: Full moon has a direct effect on sprinkler and weather. Weather has a direct effect on ground. Sprinkler has a direct effect on ground. For situations where there is no full moon, the probability of wet ground is 38%. For situations where there is a full moon, the probability of wet ground is 43%. Does full moon negatively affect ground through weather and sprinkler?

Guidance: Address the question by following the steps below:

Step 1) Extract the causal graph: Identify the causal graph that depicts the relationships in the scenario. The diagram should simply consist of edges denoted in "var1-> var2" format, separated by commas.

Step 2) Determine the query type: Identify the type of query implied by the main question. Choices include "marginal probability", "conditional probability", "explaining away effect", "backdoor adjustment set", "average treatment effect", "collider bias", "normal counterfactual question", "average treatment effect on treated", "natural direct effect" or "natural indirect effect". Your answer should only be a term from the list above, enclosed in quotation marks.

Step 3) Formalize the query: Translate the query into its formal mathematical expression based on its type, utilizing the "do(·)" notation or counterfactual notations as needed.

Step 4) Gather all relevant data: Extract all the available data. Your answer should contain nothing but marginal probabilities and conditional probabilities in the form "P(...)=..." or "P(...|...)=...", each probability being separated by a semi-colon. Stick to the previously mentioned denotations for the variables.

Step 5) Deduce the estimand using causal inference: Given all the information above, deduce the estimand using skills such as do-calculus, counterfactual prediction, and the basics of probabilities. Answer step by step.

Step 6) Calculate the estimand: Insert the relevant data in Step 4 into the estimand, perform basic arithmetic calculations, and derive the final answer. There is an identifiable answer. Answer step by step.

Begin your response with reasoning or evidence to support your explanation, then return me the final result marked by ####.

Let's solve the question step-by-step as per the provided guidance:

Step 1: Extract the causal graph The relationships described are:

- Full moon affects sprinkler and weather. - Weather affects ground. - Sprinkler affects ground.

The causal graph is:

“ FullMoon -> Weather, FullMoon -> Sprinkler, Weather -> Ground, Sprinkler -> Ground “

Step 2: Determine the query type The query asks whether the full moon negatively affects the ground through weather and sprinkler. This involves understanding the total effect of Full Moon on the probability of wet ground, considering direct and indirect causal pathways.

The correct term is "natural indirect effect".

Step 3: Formalize the query We want to assess the effect of Full Moon (F) on Ground (G) mediated through Weather (W) and Sprinkler (S). This can be expressed as:

$$NIE = E[G | do(F = 1)] - E[G | do(F = 0)]$$

Where: - G = Ground - F = 1 (Full Moon present), F = 0 (Full Moon absent) - "do" notation removes any confounding to isolate causal effects.

Step 4: Gather all relevant data From the problem, we are given:

1. $P(G = 1 | F = 0) = 0.38$ (No full moon → wet ground). 2. $P(G = 1 | F = 1) = 0.43$ (Full moon → wet ground).

Step 5: Deduce the estimand using causal inference To determine whether Full Moon negatively affects Ground through Weather and Sprinkler, we compare the probabilities $P(G = 1 | F = 1)$ and $P(G = 1 | F = 0)$. Specifically:

- If $P(G = 1 | F = 1) > P(G = 1 | F = 0)$, then the Full Moon positively affects Ground. - If $P(G = 1 | F = 1) < P(G = 1 | F = 0)$, then the Full Moon negatively affects Ground.

Here, the difference:

$$\Delta = P(G = 1 | F = 1) - P(G = 1 | F = 0)$$

Step 6: Calculate the estimand Substitute the values given:

$$\Delta = P(G = 1 | F = 1) - P(G = 1 | F = 0)$$

$$\Delta = 0.43 - 0.38$$

$$\Delta = 0.05$$

Since $\Delta > 0$, the Full Moon positively affects the probability of wet ground through Weather and Sprinkler. Therefore, there is no negative effect.

Final Result

No

Figure 11: The causalCoT prompt used in our evaluation.

FUNCTIONAL CALL EXAMPLE

Question:

Imagine a self-contained, hypothetical world with only the following conditions, and without any unmentioned factors or causal relationships:

Smoking has a direct effect on effort and freckles. Effort has a direct effect on freckles. For nonsmokers who are lazy, the probability of freckles is 92%. For nonsmokers who are hard-working, the probability of freckles is 50%. For smokers who are lazy, the probability of freckles is 80%. For smokers who are hard-working, the probability of freckles is 36%.

For nonsmokers, the probability of being hard-working is 17%. For smokers, the probability of being hard-working is 74%. If we disregard the mediation effect through effort, would smoking negatively affect freckles?

Begin your response with reasoning or evidence to support your explanation, then return the final result marked by ####. The answer format is #### Yes or #### No.

Code interpreter request:

Given probabilities:

`p_f_nonsmoker_lazy = 0.92`

`p_f_nonsmoker_hardworking = 0.50`

`p_f_smoker_lazy = 0.80`

`p_f_smoker_hardworking = 0.36`

Given probability of being hardworking:

`p_hardworking_nonsmoker = 0.17`

`p_hardworking_smoker = 0.74`

Derive probabilities of being lazy:

`p_lazy_nonsmoker = 1 - p_hardworking_nonsmoker`

`p_lazy_smoker = 1 - p_hardworking_smoker`

Calculate weighted probabilities:

`p_freckles_nonsmoker = (p_f_nonsmoker_lazy × p_lazy_nonsmoker) + (p_f_nonsmoker_hardworking × p_hardworking_nonsmoker)`

`p_freckles_smoker = (p_f_smoker_lazy × p_lazy_smoker) + (p_f_smoker_hardworking × p_hardworking_smoker)`

`p_freckles_nonsmoker, p_freckles_smoker`

Code interpreter response: Code executed successfully.

Code interpreter request:

Output the results

`p_freckles_nonsmoker, p_freckles_smoker`

Code interpreter response: (0.8486, 0.4744)

Final response:

The computed probabilities of having freckles, when disregarding the mediation effect through effort, are:

$$P(\text{freckles} \mid \text{nonsmoker}) = 84.86\%$$

$$P(\text{freckles} \mid \text{smoker}) = 47.44\%$$

Conclusion By disregarding the mediation effect through effort, smoking appears to have a negative effect on the probability of developing freckles. Smokers have a lower probability of having freckles compared to nonsmokers.

####No

Figure 12: A Function Call example from CLADDER dataset evaluation.