ON THE ELIGIBILITY OF LLMS FOR COUNTERFACTUAL REASONING: A DECOMPOSITIONAL STUDY

Anonymous authorsPaper under double-blind review

ABSTRACT

Counterfactual reasoning has emerged as a crucial technique for generalizing the reasoning capabilities of large language models (LLMs). By generating and analyzing counterfactual scenarios, researchers can assess the adaptability and reliability of model decision-making. Although prior work has shown that LLMs often struggle with counterfactual reasoning, it remains unclear which factors most significantly impede their performance across different tasks and modalities. In this paper, we propose a decompositional strategy that breaks down the counterfactual generation from causality construction to the reasoning over counterfactual interventions. To support decompositional analysis, we investigate 11 datasets spanning diverse tasks, including natural language understanding, mathematics, programming, and vision-language tasks. Through extensive evaluations, we characterize LLM behavior across each decompositional stage and identify how modality type and intermediate reasoning influence performance. By establishing a structured framework for analyzing counterfactual reasoning, this work contributes to the development of more reliable LLM-based reasoning systems and informs future elicitation strategies.

1 Introduction

Large language models (LLMs) have exhibited remarkable proficiency across a diverse range of tasks, including natural language understanding (Devlin et al., 2019; Kuang et al., 2025) and multimodal reasoning (Hu et al., 2017; Lu et al., 2022; Yang et al., 2023). Despite these advancements, concerns persist regarding their reasoning and generalization capabilities. A particularly challenging aspect of model evaluation is **Counterfactual Reasoning**, i.e., the ability to adjust responses when presented with modified premises (Pearl & Mackenzie, 2018) (e.g., *What is the outcome in a hypothetical condition?*). Investigating the counterfactual reasoning of LLMs provides an interpretable step to understand their adaptability under hypothetical alterations to input conditions (Gat et al., 2024; Huang et al., 2024).

Prior studies have demonstrated that LLMs often struggle with counterfactual reasoning and frequently fail to maintain logical consistency or adjust to context shifts (Li et al., 2023; Nguyen et al., 2024; Wang et al., 2024). While these works highlight notable performance gaps, they lack a standardized framework for systematically analyzing and understanding counterfactual behaviors in LLMs. Consequently, it remains unclear what factors most significantly impact LLM performance in counterfactual scenarios. Furthermore, counterfactual reasoning has often been evaluated in a direct and monolithic manner, primarily by introducing interventions and assessing model responses (Li et al., 2023), without grounding the analysis in the underlying causal structure that gives rise to such interventions. This overlooks the foundational role of causal modeling. Specifically, the identification of causal variables and their dependencies are essential for understanding counterfactuals.

To address these gaps, we outline a **Decompositional Strategy** that breaks down the analysis of counterfactual reasoning into distinct stages. Our approach departs from prior work that focuses solely on counterfactual generation. Instead, we begin by examining the causal structure of factual conditions, which serves as the necessary foundation for valid counterfactual reasoning.

As illustrated in Figure 1, our methodology is outlined into four stages. First, we assess (i) whether LLMs can accurately identify the four variable groups critical to causal reasoning: Exposure, Covariate, Mediator, and Outcome. Next, we evaluate (ii) whether LLMs can correctly construct a

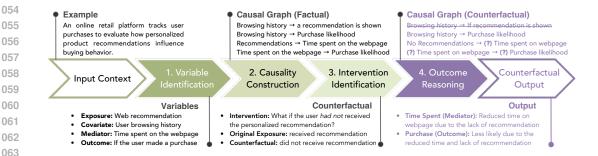


Figure 1: A workflow and illustrative example that decomposes LLM-based counterfactual reasoning into four stages: (1) identifying causal variables (e.g., whether web recommendation is shown), (2) constructing the causal graph (e.g., browsing history \rightarrow a recommendation is shown), (3) specifying the counterfactual intervention (e.g., no recommendation shown), and (4) reasoning about the counterfactual outcome (e.g., less likely to purchase a product online).

corresponding causal graph in the form of a directed acyclic graph (DAG). Building on this causality modeling, we then study LLMs' counterfactual reasoning abilities by evaluating (iii) whether they can identify the correct intervened variable (i.e., the Exposure), and (iv) whether they can accurately infer the counterfactual mediators and final outcomes by reasoning over the updated causal graph.

To support our decompositional study, we construct a benchmark by collecting and curating 11 counterfactual datasets across diverse tasks, including natural language understanding, mathematics, programming, and vision-language reasoning. We curate each dataset by extracting factual and counterfactual variables, identifying causal elements, and constructing corresponding causal graphs as reference structures for evaluation purpose. In experiments, we test the performance of leading LLMs across each decompositional stage to analyze their sufficiency in handling individual reasoning components. Based on the observed performance among these decomposed evaluations, we propose targeted improvements, such as integrating modality-specific function-calling interfaces within a toolaugmented learning paradigm, to address critical reasoning bottlenecks. Additionally, we evaluate the impact of different elicitation (prompting) strategies, including Chain-of-Thought (CoT) (Wei et al., 2022), Chain-of-Thought with Self-Consistency (CoT-SC) (Wang et al., 2022), and Tree-of-Thought (ToT) (Yao et al., 2023) reasoning. Collectively, our evaluations provide a solid step for understanding and enhancing LLMs in complex reasoning tasks and imaginative scenarios ¹.

In summary, we make the following contributions:

- **Decompositional Framework**—We propose a decompositional strategy that spans from causal modeling to counterfactual reasoning, enabling a systematic evaluation of LLMs' capabilities in understanding and performing counterfactual tasks.
- Benchmark Construction—We construct a comprehensive evaluation benchmark by curating causal structures and counterfactual instances across multiple domains. This benchmark standardizes decompositional evaluations and supports consistent analysis across tasks and modalities.
- Evaluation and Improvement Strategy—We evaluate leading LLMs under diverse tasks. By identifying LLMs' capabilities in specific decompositional stage, we propose actionable strategies to improve LLMs' counterfactual adaptability.

2 Related Work

Counterfactual Reasoning. A fundamental component of causal inference is to examine hypothetical scenarios, which addresses the question: *What would have occurred had a particular factor or decision differed?* (Pearl & Mackenzie, 2018) This method facilitates causal analysis by comparing observed outcomes with those projected under alternative conditions. Empirical research has

¹Codes available at: https://anonymous.4open.science/r/Counterfactual_NeurIPS_ 2025-D8E6/.

 demonstrated the broad applicability of counterfactual reasoning across multiple domains, including healthcare, business, and fairness (Gvozdenović et al., 2021; Kyrimi et al., 2025; Gow et al., 2016; Kasirzadeh & Smart, 2021; Koonce et al., 2011). This stands in contrast to many contemporary AI systems, which predominantly rely on statistical correlations while lacking robust capacities for abstract reasoning and causal inference (Jiao et al., 2024).

Counterfactual Reasoning in AI and NLP. Counterfactual reasoning has emerged as a powerful framework for enhancing model interpretability and causal understanding in AI and NLP. In medical AI, SyncTwin (Qian et al., 2021) proposed a counterfactual estimation framework that constructs synthetic patient data to predict potential outcomes under alternative treatments. In NLP, a Counterfactual Reasoning Model (CRM) (Feng et al., 2021) is developed using LLMs to generate contrastive samples, improving sentiment analysis and inference tasks. There are also order-faithfulness metrics (Gat et al., 2024) to evaluate causal explanations in black-box models. These contributions demonstrate the versatility of counterfactual methods in improving model transparency and reliability across domains.

LLMs and Elicitation. Recent advances have significantly enhanced LLMs' reasoning capabilities through several elicitation (prompting) approaches. The introduction of Chain-of-Thought (CoT) prompting (Wei et al., 2022) and its extensions, the Self-Consistency CoT (Wang et al., 2022) and Tree-of-Thought (ToT) (Yao et al., 2023), have enabled more structured and reliable multi-step reasoning. These innovations collectively represent a paradigm shift from simple pattern recognition to deliberate, verifiable reasoning in LLMs.

3 METHODOLOGY: DECOMPOSING COUNTERFACTUAL REASONING

This section presents our methodology for decomposing counterfactual reasoning. We begin by introducing foundational concepts in causality and counterfactual reasoning (Section 3.1). Subsequently, we detail the evaluation tasks used to assess models (Section 3.2) and our corresponding construction of benchmarks over multimodal datasets (Section 3.3).

3.1 PRELIMINARY: FROM CAUSALITY TO COUNTERFACTUAL REASONING

Causality. Causality depicts the dependencies about how one variable influences another, i.e., the underlying causal effects. There are four types of variables commonly used in causal analysis: exposure, covariate, mediator, and outcome. Specifically: (1) Exposure (or treatment, intervention, denoted X) refers to the action or condition imposed on a system; (2) Outcome (Y) denotes the resulting response or effect influenced by the exposure; (3) Covariate (Z) is the pre-treatment variable that may influence both X and Y; (4) Mediator (M) lies on the causal pathway from X to Y, representing intermediate mechanisms through which the exposure exerts its influence.

Example 1 Consider a dataset that records students' academic performance in the presence of a tutoring tool. Here, the exposure X indicates whether a student used the tool. The outcome Y corresponds to the student's final exam score. The covariate Z may include socioeconomic factors (e.g., parental income), which could influence both the use of tutoring and academic outcomes. The mediator M is the number of hours the student spends studying per week.

Causal Graph. The relationships among exposure, covariate(s), mediator(s), and outcome(s) can be formally represented using a directed acyclic graph (DAG), commonly named a *causal graph* (Pearl & Mackenzie, 2018) that captures causal relationships, where the exposure X influences the outcome Y both directly and indirectly through a mediator M. Covariate Z may affect X, M, and Y, as illustrated in Figure 2.

Example 2 The corresponding causal graph, illustrated in Figure 2-(a), would include: an arrow from $X \to M$ (e.g., tutoring influences study time), $M \to Y$ (study time influences exam performance), $X \to Y$ (direct effect of tutoring on scores), $Z \to X$ and $Z \to Y$ (e.g., socioeconomic status affects both tutoring usage and academic performance).

Counterfactual Reasoning. Counterfactual reasoning aims to answer:

163

164

165

167

169 170

171 172

173

174 175

176

177

178

179

181

182

183

185

186

187

188

189

190 191 192

193

196

197

199 200

201

202

203

204

205

206

207

208

210

211

212 213 214

215



Figure 2: (a) Causal graph structure and (b)(c) A factual/counterfactual example.

Given an observed instance (X = x, Z = z, M = m, Y = y), what would the outcome Y be if the exposure X were set to a different value x', while keeping the covariate Z fixed?

The observed instance (x, z, m, y) is also known as the **factual** case. In contrast, counterfactual reasoning seeks to determine the outcomes under an alternate intervention, that is, when X = x'.

In causal graph, we assume a two-stage causal mechanism: (1) Mediator Function: $M = f_M(X, Z)$ and (2) Outcome Function: $Y = f_Y(X, M, Z)$. Then, the counterfactual outcome under an alternative exposure x' can be computed via: $Y_{x'} = f_Y(x', f_M(x', z), z)$, where we simulate a new mediator value $M_{x'} = f_M(x', z)$ based on the counterfactual exposure x' and observed covariate z. Then, we predict the counterfactual outcome $Y_{x'}$ using the counterfactual exposure x', the simulated mediator $M_{x'}$, and the same covariate z.

Example 3 As exemplified in Figure 2-(b), a student with z = LOW-INCOME socioeconomic status did not receive tutoring (x=0), studied 5 hours per week (m=5), and scored 78 on the exam (y = 78).

We now ask: What would the student's score have been if they had received tutoring (x'=1)? We compute: (i) Simulated study time: $m' = f_M(x' = 1, z = \text{LOW-INCOME}) = 9$ and (ii) counterfactual score: $Y_{x'=1} = f_Y(x'=1, m'=9, z=\text{LOW-INCOME}) = 85.$

We conclude: The tutoring would have increased the student's score from 78 to 85 (Figure 2-(c)).

3.2 **DECOMPOSITIONAL EVALUATION TASK**

Counterfactual reasoning is often described as a structured chain of analysis from identifying variables to modeling causal relations, specifying interventions, and simulating outcomes (Pearl, 2009; Bareinboim et al., 2022). Recent work in ML and LLMs further emphasizes the importance of disentangling these steps to evaluate reasoning capacity (Kiciman et al., 2023; Chi et al., 2024). Motivated by the need for decomposition, we design four evaluation tasks that reflect the full pipeline of counterfactual reasoning, with each task targeting a distinct capability in counterfactual analysis.

- Task I: Causal Variable Identification. Given inputs containing factual information, a model is required to identify the values of the causal variables (X, Z, M, Y). This step serves as the foundation for subsequent causal modeling and counterfactual reasoning.
- Task II: Causal Graph Construction. Given the identified variables, the model is tasked with constructing a DAG that captures the causal relationships among them. This step evaluates the model's ability to discover causal dependencies.
- Task III: Counterfactual Identification. Given a counterfactual query (e.g., "What if variable X had been different?"), the LLM must identify the new value of (i.e., the intervention). This task evaluates whether the model can detect intervention in the counterfactual condition.
- **Task IV: Outcome Reasoning.** Based on the constructed causal graph and identified intervention, the model is prompted to predict the counterfactual outcome. This step measures whether the model can simulate the hypothetical scenario while respecting the underlying causal mechanisms.

3.3 BENCHMARKING COUNTERFACTUALS

Next, we introduce the datasets we leverage for the decompositional evaluations:

Table 1: Summary of counterfactual benchmarks including data source, use case, presence of causal variables (●: present, ●: partially present), the definition of counterfactual condition, included modalities, and number of instances. Concrete examples are shown in Appendix A.

| Data | Use Case | Ca | usal | Varia | ble | Counterfactual | Modality | Num |
|--------------------------------------|---------------------------|----------------|------|-------|-----|-------------------------------------|-------------|--------|
| Data | Use Case | \overline{X} | Z | M | Y | Condition | Modality | Nulli |
| CRASS (Frohberg & Binder, 2021) | Question answering | • | • | • | • | "What if" condition | Text | 274 |
| CLOMO (Huang et al., 2023) | Text logic parsing | • | • | • | • | New premise for textual statement | Text | 1,100 |
| RNN-Typology (Ravfogel et al., 2019) | Text syntax parsing | • | • | • | • | New syntactic structure of sentence | Text | 584 |
| CVQA-Bool (Zhang et al., 2024b) | Question answering | • | • | • | • | Hypothetical behavioral pattern | Text,Image | 1,130 |
| CVQA-Count (Zhang et al., 2024b) | Numerical reasoning | • | • | • | • | Hypothetical numerical pattern | Text,Image | 2,011 |
| COCO (Le et al., 2023) | Text-image matching | • | • | • | • | "What if" condition | Text,Image | 17,410 |
| Arithmetic (Wu et al., 2024) | Mathematical reasoning | • | • | • | • | Change number base | Symbol | 6,000 |
| MalAlgoQA (Sonkar et al., 2024) | Question Answering | • | • | • | • | "What if" condition | Text,Symbol | 807 |
| HumanEval-Exe (Chen et al., 2021) | Code Execution simulation | • | • | • | • | Hypothetical coding criterion | Text,Code | 981 |
| Open-Critic (Vezora, 2024) | Code generation | • | • | • | • | Hypothetical descriptive functions | Text,Code | 8,910 |
| Code-Preference (Vezora, 2024) | Code summarization | • | • | • | • | Hypothetical code structures | Text,Code | 9,389 |

Data Sources and Use Cases. As shown in Table 1, we collect a diverse set of datasets to ensure broad coverage across various NLP tasks and modalities. The included use cases are: (1) Question Answering, evaluated using CRASS (Frohberg & Binder, 2021), CVQA-Bool (Zhang et al., 2024b), and MalAlgoQA (Sonkar et al., 2024), which involve answering general-purpose textual or visually grounded questions; (2) Text Parsing, using CLOMO (Huang et al., 2023) for logical structure reconstruction and RNN-Typolog (Ravfogel et al., 2019) for syntactic structure understanding; (3) Reasoning Tasks, with CVQA-Count (Zhang et al., 2024b) for numerical reasoning and Arithmetic (Wu et al., 2024) for symbolic arithmetic computation; (4) Multimodal Matching, represented by the COCO dataset (Le et al., 2023) for image-text alignment; (5) Code-based Tasks, including HumanEval-Exe (Chen et al., 2021) for execution simulation, Open-Critic (Vezora, 2024) for generation, and Code-Preference (Vezora, 2024) for summarization. These datasets are therefore intentionally positioned as a prerequisite stage to support safe and informed downstream application.

These datasets span four modalities, natural language text, images, mathematical symbols, and code, that encompass diverse definitions of counterfactual interventions tailored to each task. Collectively, they support a comprehensive multimodal evaluation of LLMs' abilities to reason under varied counterfactual settings and data types.

Our Preprocessing. To support our decompositional evaluations, we curate those datasets to augment each instance with three additional aspects of information relevant to Tasks I–III (Section 3.2). Specifically, we begin by identifying and annotating the causal variables (X,Z,M,Y) from the original data, questions, or descriptions. Using these annotations, we construct a DAG to represent the underlying causal structure of each data instance, which enrich original instances with causal and counterfactual structures. A running example in Figure 2.

Preprocessing Feasibility. Notably, all datasets are built upon "what-if" conditions or hypothetical scenarios (as outlined in Table 1) and intervention-style narratives, thus naturally supporting counterfactual interventions. For each instance, we parse and extract the intervened variables and, guided by the previously constructed DAG, annotate the corresponding counterfactual outcomes and construct a matched counterfactual graph. We provide the curated instances in Appendix A.

4 EXPERIMENT

We aim to empirically answer two research questions: $\mathbf{RQ_1}$: How well do LLMs perform when their counterfactual reasoning is decomposed into distinct reasoning tasks? $\mathbf{RQ_2}$: What auxiliary techniques can improve LLMs' counterfactual reasoning? We defer the experimental settings into Appendix B.1 and additional results into Appendix B.2.

LLMs. We evaluate reasoning-centric and multimodal LLMs due to the nature of counterfactual reasoning tasks. Specifically, we leverage GPT-5, GPT-04-mini-high, Qwen3-VL-235B-A22B-Thinking, Llama-4-Scout-17Bt, Llama-4-Maverick-17B-128E, Gemini2.5-Pro, DeepSeek-VL.

Metrics. We use the F1 score for Tasks I, II, and IV, as they involve multiple instances (e.g., M, Z, or graph edges) that require set-level evaluation. For Task III, which typically involves a single intervention on X, we use accuracy to assess whether the LLM correctly identifies the intervened X.

Table 2: LLMs' performance in causal variable identification, we report means of F1 across all instances for each variable. Each value is scaled to 100%. The standard deviation is in Table 7.

| Dataset | GP | T-5 | GP | GPT-o4 | | ren3 | Llan | na4-S | Llama4-M | | Gemini2.5 | | DeepSeek | |
|-----------------|-------|-------|-------|-----------------|--------|-------------|------------|---------|----------|-------|-----------|-------|----------|-------|
| Dataset | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 |
| | | | | $v_1 = \lambda$ | (Expo | sure), v | $_{2}=Z$ (| Covaria | ite) | | | | | |
| CRASS | 92.3 | 91.1 | 91.0 | 89.2 | 87.3 | 85.4 | 88.5 | 86.9 | 90.6 | 89.1 | 88.6 | 86.2 | 89.5 | 87.1 |
| CLOMO | 89.8 | 87.6 | 88.1 | 85.6 | 83.5 | 81.9 | 87.1 | 85.3 | 88.8 | 87.0 | 84.3 | 82.8 | 86.4 | 84.2 |
| RNN-Topo | 87.9 | 85.4 | 85.9 | 83.7 | 80.4 | 78.6 | 84.3 | 82.6 | 85.7 | 84.2 | 81.7 | 80.1 | 83.8 | 82.0 |
| CVQA-Bool | 79.4 | 76.2 | 79.8 | 76.5 | 72.3 | 70.5 | 77.9 | 75.8 | 79.1 | 76.9 | 68.5 | 66.9 | 70.7 | 68.3 |
| CVQA-Count | 74.7 | 72.3 | 74.3 | 72.9 | 68.9 | 67.2 | 73.6 | 71.9 | 74.4 | 72.5 | 65.8 | 63.7 | 67.4 | 65.1 |
| COCO | 72.8 | 70.2 | 73.2 | 71.1 | 67.2 | 65.4 | 72.5 | 70.8 | 73.6 | 71.7 | 62.6 | 60.9 | 65.9 | 63.2 |
| Arithmetic | 88.2 | 86.5 | 84.9 | 82.8 | 75.7 | 73.8 | 80.3 | 78.6 | 81.6 | 79.8 | 76.9 | 74.5 | 78.3 | 76.1 |
| MalAlgoQA | 84.1 | 81.3 | 81.5 | 78.9 | 72.9 | 70.6 | 79.5 | 77.1 | 80.6 | 78.2 | 73.5 | 71.2 | 75.9 | 73.4 |
| HumanEval-Exe | 69.3 | 66.9 | 71.4 | 69.2 | 63.7 | 61.9 | 67.8 | 65.7 | 68.9 | 66.7 | 59.6 | 57.3 | 62.1 | 59.8 |
| Open-Critic | 71.7 | 69.4 | 70.1 | 67.3 | 61.8 | 59.7 | 66.5 | 64.7 | 67.6 | 65.8 | 57.3 | 55.9 | 60.4 | 58.1 |
| Code-Preference | 49.6 | 68.4 | 80.2 | 69.0 | 72.9 | 60.5 | 73.6 | 61.9 | 75.2 | 63.4 | 68.4 | 66.5 | 61.2 | 59.3 |
| | | | | $v_1 = \Lambda$ | I (Med | iator), v | $_2 = Y$ | (Outcon | ne) | | | | | |
| CRASS | 87.4 | 91.7 | 84.1 | 89.3 | 72.8 | 79.9 | 81.2 | 87.4 | 82.4 | 88.6 | 74.1 | 81.6 | 76.2 | 83.5 |
| CLOMO | 83.1 | 89.4 | 81.3 | 87.9 | 68.9 | 77.4 | 77.2 | 84.9 | 78.6 | 86.0 | 71.2 | 78.9 | 73.5 | 80.7 |
| RNN-Topo | 81.7 | 86.3 | 79.4 | 85.6 | 67.1 | 75.8 | 75.3 | 83.1 | 76.6 | 84.4 | 69.3 | 76.4 | 71.5 | 78.9 |
| CVQA-Bool | 73.5 | 79.3 | 73.9 | 80.1 | 59.7 | 66.9 | 71.8 | 78.3 | 72.9 | 79.4 | 57.3 | 63.2 | 60.5 | 68.4 |
| CVQA-Count | 69.6 | 75.2 | 70.1 | 76.2 | 56.8 | 63.6 | 68.9 | 74.8 | 70.2 | 75.9 | 54.7 | 60.4 | 57.9 | 65.1 |
| COCO | 67.3 | 73.4 | 68.0 | 74.4 | 54.2 | 61.8 | 66.2 | 72.1 | 67.4 | 73.5 | 52.8 | 58.3 | 55.7 | 62.6 |
| Arithmetic | 82.1 | 85.6 | 82.3 | 73.4 | 63.6 | 71.9 | 79.1 | 85.1 | 80.3 | 86.0 | 62.8 | 73.5 | 65.7 | 74.9 |
| MalAlgoQA | 79.2 | 83.4 | 76.8 | 80.4 | 61.5 | 69.3 | 76.6 | 81.2 | 77.8 | 82.4 | 60.2 | 70.5 | 63.8 | 72.3 |
| HumanEval-Exe | 63.2 | 67.4 | 66.0 | 70.3 | 51.9 | 59.7 | 61.9 | 66.3 | 63.0 | 67.5 | 49.7 | 57.0 | 53.6 | 60.5 |
| Open-Critic | 66.3 | 70.6 | 64.1 | 68.9 | 49.7 | 57.2 | 64.7 | 69.0 | 65.6 | 70.0 | 47.3 | 54.8 | 51.3 | 58.7 |
| Code-Preference | 65.9 | 75.3 | 66.3 | 77.0 | 50.4 | 78.6 | 64.5 | 79.3 | 65.7 | 80.4 | 48.2 | 76.1 | 52.5 | 59.8 |

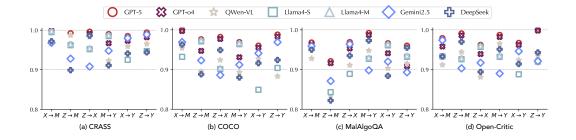


Figure 3: Evaluation on causal graph construction. We evaluate F1 score to balance (i) whether the constructed edges under one category (e.g., $X \to M$) is correctly constructed if the (X, Z, M, Y) are already given. Additional results for all other datasets at Figure 5.

4.1 LLM PERFORMANCE ON DECOMPOSITIONAL TASKS (RQ₁)

Setting. We evaluate LLMs independently on each decompositional task. For each task, we explicitly provide the ground-truth outputs from the preceding tasks to isolate and measure LLMs' capabilities specific to that task. For example, when assessing models' ability to construct causal graphs, we supply the original inputs along with the ground-truth causal variables.

Task I: Causal Variable Identification. Table 2 presents the performance of LLMs on identifying causal variables (X, Z, M, Y). We observe that model performance is strongly influenced by the modality of the dataset. Specifically, datasets involving more complex modalities (e.g., images, mathematical symbols, codes) tend to reduce LLM accuracy (e.g., <0.7 F1 on Open-Critic), even when variables like X, Z, Y are explicitly present in the context.

Interestingly, even within the text modality, LLMs show notable difficulty in identifying the implicit mediator M, which often requires reasoning about the underlying causal pathways connecting X, Z, and Y. This suggests that the challenge lies not only in the complexity of the input modality but also in the abstractness and inferential nature of the variable type itself. Together, these findings highlight

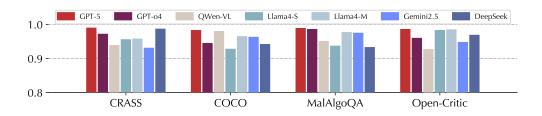


Figure 4: Evaluation of LLMs' accuracy in identifying the correct intervention (i.e., the counterfactual value of X). Additional results for all other datasets are provided in Figure 6.

Table 3: LLM performance (F1 mean) in reasoning the counterfactual mediator (M') and outcome (Y'). Standard deviation is in Table 8.

|). Standard deviation is in Tuble 6. | | | | | | | | | | | | | | |
|--------------------------------------|------|------|------|------|------|------|------|-------|------|------|------|--------|------|------|
| Dataset | GP | T-5 | GP | Г-о4 | Qw | en3 | Llan | na4-S | Llam | a4-M | Gem | ini2.5 | Deep | Seek |
| Dataset | M' | Y' | M' | Y' | M' | Y' | M' | Y' | M' | Y' | M' | Y' | M' | Y' |
| CRASS | 92.1 | 88.0 | 90.5 | 86.2 | 80.5 | 73.9 | 70.1 | 63.5 | 84.9 | 79.5 | 81.7 | 75.2 | 82.9 | 77.1 |
| CLOMO | 90.2 | 85.3 | 88.7 | 83.9 | 77.8 | 71.6 | 67.2 | 60.9 | 82.9 | 77.2 | 79.3 | 72.8 | 80.5 | 74.3 |
| RNN-Topo | 88.9 | 83.4 | 87.9 | 81.6 | 75.6 | 69.4 | 65.2 | 58.7 | 80.5 | 75.0 | 77.1 | 70.6 | 78.3 | 72.0 |
| CVQA-Bool | 81.2 | 74.5 | 77.1 | 70.2 | 65.4 | 58.6 | 54.8 | 48.1 | 70.9 | 64.3 | 63.2 | 56.8 | 66.7 | 59.8 |
| CVQA-Count | 79.2 | 72.0 | 76.2 | 69.3 | 62.2 | 55.7 | 51.7 | 45.2 | 67.5 | 61.0 | 60.1 | 53.8 | 63.5 | 56.9 |
| COCO | 77.8 | 70.1 | 75.3 | 66.9 | 60.1 | 53.4 | 49.3 | 42.7 | 65.4 | 58.7 | 57.8 | 51.5 | 61.3 | 54.6 |
| Arithmetic | 87.8 | 82.7 | 85.8 | 80.9 | 69.8 | 63.2 | 59.1 | 52.4 | 76.3 | 70.6 | 72.1 | 65.4 | 74.0 | 67.5 |
| MalAlgoQA | 85.1 | 79.6 | 83.6 | 77.8 | 67.5 | 60.9 | 57.4 | 50.7 | 74.0 | 68.2 | 69.6 | 62.9 | 71.8 | 65.1 |
| HumanEval-Exe | 75.7 | 71.5 | 73.4 | 66.5 | 58.2 | 51.5 | 47.7 | 41.2 | 63.6 | 56.9 | 55.8 | 49.4 | 59.4 | 52.7 |
| Open-Critic | 75.3 | 69.4 | 73.8 | 67.5 | 56.0 | 49.4 | 45.8 | 39.2 | 61.5 | 54.7 | 53.7 | 47.3 | 57.2 | 50.6 |
| Code-Preference | 77.0 | 71.0 | 74.4 | 66.8 | 57.1 | 50.4 | 46.6 | 40.0 | 62.7 | 55.9 | 54.7 | 48.3 | 58.3 | 51.6 |

the need for improved methods that enhance LLMs' capacity to handle both cross-modal complexity and deeper causal reasoning.

Task II: Causal Graph Construction. As described in our experimental setup, we isolate each decompositional step by providing the ground-truth outputs of preceding steps as inputs. For causal graph construction, we supply the identified variables X, Z, M, Y and prompt LLMs to construct the corresponding counterfactual graph. The results are presented in Figures 3 and 5. Notably, the overall performance mostly exceeds 0.9 F1 scores, indicating that LLMs can accurately construct graph edges. Moreover, the impact of dataset modality and variable types (e.g., explicit Z vs. implicit M) appears to be minimal in this step. We attribute this to the rule-based nature of causal graph construction: since causal graph structures are well-defined (as shown in Figure 2), it is relatively less challenging for LLMs to apply construction rules and generate the correct causal relationships.

Insights from causality modeling.

The major challenge in causal modeling lies in causal variable identification, where (1) LLMs are highly sensitive to the complexity and structure of the input modality, and (2) implicit variables (i.e., the mediator M) reveal a critical gap in LLMs' causal reasoning capabilities.

Task III: Counterfactual Identification. Next, we evaluate LLMs' capability in identifying interventions, i.e., determining the counterfactual values of X (i.e., X'). As shown in Figures 4 and 6, the experimental results indicate that LLMs are generally effective at recognizing the counterfactual values of X across most datasets and modalities. This demonstrates that LLMs have a solid grasp of pinpointing intervention points within the context. However, note that the task remains relatively isolated and does not challenge the model's ability to propagate the effects of the intervention through downstream variables (e.g., M', Y'), which we address in Task IV.

Task IV: Outcome Reasoning. In this final task, we evaluate LLMs' ability to infer the mediator (M') and outcome (Y') under a counterfactual intervention. As shown in Table 3, LLMs consistently exhibit insufficient performance in inferring these implicit variables across all datasets. Notably, since both M' and Y' are implicit under the counterfactual condition (whereas only M is implicit in the factual condition), this result suggests that LLMs lack sufficient capacity to reason over causal chains, even when the underlying structure is explicitly provided.

Table 4: Improvement in LLM performance (comparing with Table 2) for identifying explicit causal variables. Results are reported on six representative datasets spanning all major modalities.

| Dataset | | GPT-o4 | | | Qwen3 | | | Llama4-S | ; | | Gemini2.5 | | | |
|---|---------------------------------------|---------------------------------------|---------------------------------------|---|--|--|--|--|---|--|---------------------------------------|--|--|--|
| Dataset | Dataset $X Z Y$ | | Y | X | Z | Y | X | Z | Y | X | Z | Y | | |
| CRASS CLOMO CVQA-Count COCO MalAlgoQA | +6.0 +4.7 +17.7 +5.5 +4.2 | +6.6 +5.1 +15.9 +3.1 +5.9 | +5.5 +7.9 +18.2 +4.4 +2.6 | +10.8 +12.7 +21.9 +8.9 +9.6 | +9.4 +11.3 +21.4 +7.6 +8.1 | +9.7 +12.1 +22.7 +8.2 +6.8 | +15.2 +21.0 +32.0 +7.2 +12.5 | +11.5 +14.4 +26.1 +6.0 +10.2 | +11.9 +16.8 +24.1 +9.7 +6.5 | +6.6 +6.2 +18.1 +6.9 +12.4 | +5.3 +9.4 +15.7 +3.8 +8.5 | +7.1 +15.5 +19.3 +5.2 +8.3 | | |
| Open-Critic | +12.8 | +14.5 | +7.0 | +15.4 | +10.9 | +8.3 | +17.4 | +5.1 | +1.2 | +7.3 | +6.1 | +9.8 | | |

Insights from counterfactual reasoning.

Regardless of whether the setting is factual or counterfactual, the primary challenge lies not in identifying causal variables and performing causal reasoning. In particular, the complex input modality and the implicit nature of mediation hinder effective reasoning through causal pathways.

4.2 AUXILIARY TECHNIQUES TO IMPROVE COUNTERFACTUAL REASONING (RQ₂)

Given the insights from previous evaluations, we aim to correspondingly address the limitations arising from multimodal complexity and intermediate reasoning, we propose augmenting LLMs with two auxiliary techniques: (i) tool-augmented execution and (ii) advanced elicitation strategies.

4.2.1 TOOL-AUGMENTED EXECUTION IN EXPLICIT VARIABLE IDENTIFICATION

Settings. To enhance LLM performance in identifying explicit variables (X, Z, Y) across different modalities, we adopt a tool-augmented approach, where the LLM dynamically calls additional specialized tools to assist in entity identification, identical to a named entity recognition (NER) paradigm. We leverage several pretrained models tailored to the multimodality of datasets: (1) **Text-based NER:** We use BERT-BASE-NER (Devlin et al., 2018; Tjong Kim Sang & De Meulder, 2003) to identify candidate entities in text-based data, including mathematical symbols.(2) **Vision-based NER:** We employ GROUNDING-DINO-BASE (Liu et al., 2023) to detect all relevant objects in images and generate focused regions by masking out irrelevant backgrounds. (3) **Code analysis:** We adopt GRAPHCODEBERT (Guo et al., 2020) to extract functions, variables, and control structures for programming tasks.

After identifying candidate entities across each modality, we prompt the LLM to refine and filter the final set of explicit variables (X, Z, Y) according to the formal definitions provided in Section 3.1.

Experimental Results. We randomly select three representative LLMs and conduct experiments across multiple datasets. As shown in Table 4, tool-augmented execution consistently improves LLM performance in identifying explicit causal variables (X, Z, Y) across all modalities. For example, by leveraging GRAPHCODEBERT to parse code structures and forward the results to Llama, which gains a clearer understanding of programming logic and achieves an F1 improvement up to 0.189. Similarly, in the vision modality, the object detector GROUNDING-DINO-BASE assists by generating a set of candidate visual objects, which GPT-o4 can then contextualize and compose into a coherent factual variable. For instance, detected objects like "Woman," "Knife," and "Apple" can be effectively integrated by GPT-o4 into the causal expression: "A woman cuts an apple with a knife."

These results demonstrate that tool-augmented learning effectively mitigates modality-specific bottlenecks by offloading low-level entity recognition to specialized models, allowing the LLM to focus on higher-level reasoning. Looking forward, there is potential to explore alternative tool configurations that may yield comparable or even superior performance. Additionally, future work may explore multi-agent frameworks with specialized agents to collaboratively handle different variable categories.

4.2.2 ADVANCED ELICITATION STRATEGIES FOR REASONING OVER IMPLICIT VARIABLES

Settings. To enhance LLM reasoning over implicit variables, particularly factual M and counterfactual M', Y', we implement advanced elicitation strategies that guide the model through more

Table 5: Improvement of LLM performance (F1 score) in reasoning implicit variables.

| Dataset Elicitation | | | GPT-o4 | | | Qwen3 | | I | Jama4-S | | Gemini2.5 | | |
|---------------------|----------------------|-----------------------|-------------------------|------------------------|----------------------|----------------------|-------------------------|-----------------------|-------------------------------|------------------------|-----------------------|------------------------|----------------------|
| Dataset | Elicitation | M | M' | Y' | M | M' | Y' | M | M' | Y' | M | M' | Y' |
| CRASS | CoT CoT-SC ToT | +5.6 +5.3 +6.8 | +4.0 +5.0 +5.2 | +3.1 +5.5 +4.1 | +7.4 +8.6 +8.9 | +5.6 +7.3 +6.7 | +4.1 +6.0 +5.0 | +7.9 +10.5 +8.9 | +5.8 +8.2 +7.4 | +4.2 +5.7 +4.7 | +6.8 +9.4 +8.2 | +4.5 +6.8 +5.9 | +3.2 +5.0 +4.3 |
| CLOMO | CoT CoT-SC ToT | +6.4 +7.0 +10.1 | +4.9 +5.2 +3.8 | $+3.6 \\ +4.2 \\ +2.9$ | +8.3 +9.6 +8.7 | +6.5 +7.6 +5.2 | +4.8 +5.4 +3.5 | +8.2 +9.9 +7.2 | +6.9 +7.7 +5.1 | +5.0 +5.5 +3.9 | +7.6 +8.9 +12.4 | +5.3 +6.5 +4.2 | +3.9 +4.7 +3.0 |
| CVQA-Count | CoT CoT-SC ToT | +5.7 +4.1 +4.9 | $^{+4.4}_{+4.2}_{+3.8}$ | $+3.5 \\ +2.4 \\ +3.1$ | +7.9 +7.2 +6.9 | +6.1 +6.6 +5.3 | $^{+4.4}_{+4.8}_{+4.0}$ | +8.2 +5.9 +6.5 | $^{+6.4}_{+6.8}$ $^{+5.2}$ | +4.5 +5.7 +4.0 | +7.3 +12.2 +6.1 | +5.4 +7.6 +4.5 | +3.8 +4.5 +3.1 |
| COCO | CoT CoT-SC ToT | +3.8 +5.4 +4.5 | +2.9 +4.1 +3.5 | $+2.1 \\ +3.1 \\ +2.7$ | +6.0 +7.2 +6.8 | +4.6 +5.7 +5.1 | +3.2 +4.3 +3.6 | +5.4 +7.6 +6.2 | +4.0 +5.9 +5.0 | $+2.6 \\ +3.8 \\ +3.5$ | +4.9 +6.9 +5.6 | $+3.4 \\ +5.1 \\ +4.3$ | +2.3 +3.6 +3.0 |
| MalAlgoQA | CoT CoT-SC ToT | +4.6 +5.5 +6.2 | +3.5 +4.3 +4.8 | +2.7 +3.4 +3.8 | +7.0 +7.7 +8.5 | +5.5 +6.1 +6.7 | +4.0 +4.6 +5.0 | +6.6 +7.9 +8.6 | +5.0 +6.3 +7.1 | $+3.3 \\ +4.1 \\ +4.6$ | +5.8 +7.2 +8.1 | +4.0 +5.2 +6.0 | +2.7 +3.7 +4.3 |
| Open-Critic | CoT CoT-SC ToT | +3.5 +4.2 +5.0 | +2.7 +3.3 +3.8 | +2.0 +2.6 +2.9 | +5.3 +6.1 +6.7 | +4.2 +4.9 +5.2 | +3.0 +3.6 +3.8 | +5.0 +5.8 +6.1 | +3.7 +4.5 +5.3 | +2.5 +3.1 +3.6 | +4.5 +5.0 +6.3 | +3.1 +3.5 +4.7 | +2.1 +2.7 +3.4 |

structured reasoning. Specifically, we apply Chain-of-Thought (CoT) (Wei et al., 2022), CoT with Self-Consistency (CoT-SC) (Wang et al., 2022), and Tree-of-Thought (ToT) (Yao et al., 2023):

- CoT: given pre-determined explicit variables, LLMs are encouraged to infer intermediate variables step-by-step.
- CoT-SC: LLMs are prompted to generate multiple reasoning paths and select the final answer via majority voting or consensus.
- ToT: LLMs are prompted to explore multiple parallel reasoning paths in a branching structure and evaluate candidate outputs based on intermediate criteria.

In implementation, both CoT-SC and ToT are executed with k=5 sampled reasoning paths. ToT further evaluates candidate outputs by scoring their textual similarity using BERTScore (Zhang et al., 2019), specifically assessing how well the intermediate results align with the original task statement.

Experimental Results. Our experimental results in Table 5 show that advanced elicitation strategies generally lead to improved performance in reasoning over implicit variables (M, M', Y') despite factual or counterfactual cases. However, we also observe that more complex prompting strategies (e.g., CoT-SC and ToT) can sometimes perform slightly worse than simpler approaches (e.g., CoT). While these advanced methods encourage more exhaustive exploration of reasoning paths, they may also induce overthinking behavior in LLMs, leading the model to introduce unnecessary causal links or misinterpret the underlying problem structure. For instance, consider the following context "A person is running a marathon and collapses." with expected mediator "Dehydration". While CoT and CoT-SC strategies correctly identify the mediator, ToT leads LLMs to overanalyze and identify "lack of training" or "overexertion" as the mediator. These choices, although related, are not directly supported by the input data and reflect an over-extension of the reasoning process.

The over-qualification of elicitation strategies (e.g., ToT) highlights that, while advanced prompting techniques can improve reasoning capabilities, they may also introduce complexities that divert the model from the most straightforward and contextually supported causal pathways. Therefore, it's crucial to balance the reasoning depth while maintaining alignment with the input data.

5 CONCLUSION

This work provides a decompositional framework for evaluating counterfactual reasoning in LLMs. We collect a set of datasets across multimodalities, and then curate them with reference causal variables, structured graphs, and counterfactual intervention. Next, we use our curated dataset to study the reasoning process into distinct stages from causal variable identification to outcome inference. We uncover the qualifications of current LLMs in counterfactual reasoning and where they fall short. Based on experimental insights, we further propose improvements to offer actionable insights for enhancing LLM reasoning, particularly in multimodality and implicit reasoning settings.

ETHICS STATEMENT

This work does not raise ethical concerns. All experiments were conducted on publicly available datasets spanning text, image, code, and symbolic reasoning tasks (e.g., CRASS, CLOMO, COCO, HumanEval). No private, personal, or sensitive information was accessed or processed during this research. The methodology and evaluations strictly comply with the licensing terms and intended academic usage of the benchmark datasets.

REPRODUCIBILITY STATEMENT

To support reproducibility, we have curated and released the complete benchmark construction process, together with evaluation scripts, model prompts, and experimental configurations, through an anonymous GitHub repository released in Introduction section. This repository provides detailed instructions for dataset preprocessing, task decomposition, and model evaluation, enabling independent researchers to reproduce and extend our experiments.

REFERENCES

- Elias Bareinboim, Juan D Correa, Duligur Ibeling, and Thomas Icard. On pearl's hierarchy and the foundations of causal inference. In *Probabilistic and causal inference: the works of judea pearl*, pp. 507–556. 2022.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Haoang Chi, He Li, Wenjing Yang, Feng Liu, Long Lan, Xiaoguang Ren, Tongliang Liu, and Bo Han. Unveiling causal reasoning in large language models: Reality or mirage? *Advances in Neural Information Processing Systems*, 37:96640–96670, 2024.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. URL http://arxiv.org/abs/1810.04805.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.
- Fuli Feng, Jizhi Zhang, Xiangnan He, Hanwang Zhang, and Tat-Seng Chua. Empowering language understanding with counterfactual reasoning. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, pp. 2226–2236, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.196. URL https://aclanthology.org/2021.findings-acl.196/.
- Jörg Frohberg and Frank Binder. Crass: A novel data set and benchmark to test counterfactual reasoning of large language models. *arXiv preprint arXiv:2112.11941*, 2021.
- Yair Ori Gat, Nitay Calderon, Amir Feder, Alexander Chapanin, Amit Sharma, and Roi Reichart. Faithful explanations of black-box nlp models using llm-generated counterfactuals. In *The Twelfth International Conference on Learning Representations*, 2024.
- Ian D Gow, David F Larcker, and Peter C Reiss. Causal inference in accounting research. *Journal of Accounting Research*, 54(2):477–523, 2016.
- Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. Graphcodebert: Pre-training code representations with data flow. *arXiv preprint arXiv:2009.08366*, 2020.

Emilia Gvozdenović, Lucio Malvisi, Elisa Cinconze, Stijn Vansteelandt, Phoebe Nakanwagi, Emmanuel Aris, and Dominique Rosillon. Causal inference concepts applied to three observational studies in the context of vaccine development: from theory to practice. *BMC Medical Research Methodology*, 21:1–10, 2021.

- Ronghang Hu, Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Kate Saenko. Learning to reason: End-to-end module networks for visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pp. 804–813, 2017.
- Yinya Huang, Ruixin Hong, Hongming Zhang, Wei Shao, Zhicheng Yang, Dong Yu, Changshui Zhang, Xiaodan Liang, and Linqi Song. Clomo: Counterfactual logical modification with large language models. *arXiv preprint arXiv:2311.17438*, 2023.
- Yinya Huang, Ruixin Hong, Hongming Zhang, Wei Shao, Zhicheng Yang, Dong Yu, Changshui Zhang, Xiaodan Liang, and Linqi Song. CLOMO: Counterfactual logical modification with large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 11012–11034, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.593. URL https://aclanthology.org/2024.acl-long.593/.
- Licheng Jiao, Yuhan Wang, Xu Liu, Lingling Li, Fang Liu, Wenping Ma, Yuwei Guo, Puhua Chen, Shuyuan Yang, and Biao Hou. Causal inference meets deep learning: A comprehensive survey. *Research*, 7:0467, 2024.
- Atoosa Kasirzadeh and Andrew Smart. The use and misuse of counterfactuals in ethical machine learning. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pp. 228–236, 2021.
- Emre Kiciman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language models: Opening a new frontier for causality. *Transactions on Machine Learning Research*, 2023.
- Lisa Koonce, Karen K Nelson, and Catherine M Shakespeare. Judging the relevance of fair value for financial instruments. *The Accounting Review*, 86(6):2075–2098, 2011.
- Jiayi Kuang, Ying Shen, Jingyou Xie, Haohao Luo, Zhe Xu, Ronghao Li, Yinghui Li, Xianfeng Cheng, Xika Lin, and Yu Han. Natural language understanding and inference with mllm in visual question answering: A survey. *ACM Comput. Surv.*, 57(8), March 2025. ISSN 0360-0300. doi: 10.1145/3711680. URL https://doi.org/10.1145/3711680.
- Evangelia Kyrimi, Somayyeh Mossadegh, Jared M Wohlgemut, Rebecca S Stoner, Nigel RM Tai, and William Marsh. Counterfactual reasoning using causal bayesian networks as a healthcare governance tool. *International Journal of Medical Informatics*, 193:105681, 2025.
- Tiep Le, Vasudev Lal, and Phillip Howard. Coco-counterfactuals: Automatically constructed counterfactual examples for image-text pairs. *Advances in Neural Information Processing Systems*, 36:71195–71221, 2023.
- Jiaxuan Li, Lang Yu, and Allyson Ettinger. Counterfactual reasoning: Testing language models' understanding of hypothetical scenarios. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (*Volume 2: Short Papers*), pp. 804–815, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-short.70. URL https://aclanthology.org/2023.acl-short.70/.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. Grounding dino: Marrying dino with grounded pre-training for open-set object detection, 2023.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022.

Van Bach Nguyen, Paul Youssef, Christin Seifert, and Jörg Schlötterer. LLMs for generating and evaluating counterfactuals: A comprehensive study. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 14809–14824, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.870. URL https://aclanthology.org/2024.findings-emnlp.870/.

Judea Pearl. Causality. Cambridge university press, 2009.

- Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect.* Basic books, 2018.
- Zhaozhi Qian, Yao Zhang, Ioana Bica, Angela Wood, and Mihaela van der Schaar. Synctwin: Treatment effect estimation with longitudinal outcomes. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 3178–3190. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/19485224d128528da1602ca47383f078-Paper.pdf.
- Shauli Ravfogel, Yoav Goldberg, and Tal Linzen. Studying the inductive biases of rnns with synthetic variations of natural languages. *arXiv* preprint arXiv:1903.06400, 2019.
- Shashank Sonkar, Naiming Liu, MyCo Le, and Richard Baraniuk. Malalgoqa: Pedagogical evaluation of counterfactual reasoning in large language models and implications for ai in education. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 15554–15567, 2024.
- Erik F. Tjong Kim Sang and Fien De Meulder. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pp. 142–147, 2003. URL https://www.aclweb.org/anthology/W03-0419.
- Vezora. Code preference pairs. https://huggingface.co/datasets/Vezora/ Code-Preference-Pairs, 2024.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Yongjie Wang, Xiaoqi Qiu, Yu Yue, Xu Guo, Zhiwei Zeng, Yuhong Feng, and Zhiqi Shen. A survey on natural language counterfactual generation. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 4798–4818, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.276. URL https://aclanthology.org/2024.findings-emnlp.276/.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim, Jacob Andreas, and Yoon Kim. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 1819–1862, 2024.
- Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*, 2023.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: deliberate problem solving with large language models. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.

 Chunhui Zhang, Yiren Jian, Zhongyu Ouyang, and Soroush Vosoughi. Working memory identifies reasoning limits in language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 16896–16922, Miami, Florida, USA, November 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.938. URL https://aclanthology.org/2024.emnlp-main.938/.

Letian Zhang, Xiaotong Zhai, Zhongkai Zhao, Yongshuo Zong, Xin Wen, and Bingchen Zhao. What if the tv was off? examining counterfactual reasoning abilities of multi-modal language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21853–21862, 2024b.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.

A COMPLEMENTARY INFORMATION OF CAUSALITY AND CAUSAL GRAPH

This appendix presents a concise overview of each dataset, followed by its causal structure and graphical representation, alongside a concrete example. For each dataset, we identify the four variable types—Exposure, Covariate, Mediator, and Outcome—and distinguish their roles in both factual and counterfactual scenarios, illustrating each with directed edges in the corresponding causal graph. In addition, we include a sample prompt used to generate responses on a simple text-parsing Q&A task dataset.

CRASS Example. The **Counterfactual Reasoning Assessment for Structured Scenarios (CRASS)** dataset is designed to evaluate whether language models can reason about hypothetical alternatives to factual events. Each example in CRASS presents a factual scenario (e.g., "A woman opens a treasure chest") followed by a counterfactual question (e.g., "What would have happened if the woman had not opened the treasure chest?"). Models are asked to select the most logically consistent outcome from multiple-choice options, such as "The treasure chest would have remained closed", which is labeled as correct. The following displays a full example:

```
"input": "A woman opens a treasure chest. What would have happened if
    the woman had not opened the treasure chest?",
"target_scores": {
    "The treasure chest would have been open.": 0,
    "That is not possible.": 0,
    "The treasure chest would have remained closed.": 1,
    "I don't know.": 0
}
```

CRASS Prompt.

```
# Role: Counterfactual Reasoning Analyst
Your task is to identify causal variables and generate causal graphs
    with edges. All outputs must follow strict JSON formatting for
    automated saving.

## Context Format
**Premise:**
{premise}

**Questionized Counterfactual Conditional (QCC):**
{qcc}

## Task Instructions
1. **Variable Identification**
- Extract and classify variables into:
```

```
702
               - 'Exposure (X)': Direct intervention (verb phrase, e.g., "act of
703
                    opening").
704
                 'Outcome (Y) ': Final result (state phrase, e.g., "chest opened"
705
                   ) .
                - 'Covariate (Z)': Pre-existing conditions influencing X/Y.
706
               - 'Mediator (M)': Mechanism linking X-Y (physical/logical process
707
                   ) .
708
709
           2. **Causal Graph Construction**
710

    Build factual/counterfactual edges following these strict rules:

711
           **Factual Graph Requirements:**
712
            "\rightarrow" Must include ALL these relationships:
713
               1. Z - X (Covariate influences Exposure)
714
               2. Z - M (Covariate influences Mediator)
715
               3. Z - Y (Covariate influences Outcome)
               4. X - M (Exposure affects Mediator)
716
               5. M - Y (Mediator affects Outcome)
717
               6. X - Y (Direct exposure effect)
718
719
           **Counterfactual Graph Requirements:**
720
           - Must include ALL these relationships:
               1. Z - M' (Same covariate influence on adjusted Mediator)
721
               2. Z - Y' (Same covariate influence on new Outcome)
722
               3. {\ensuremath{\text{X'}}} - {\ensuremath{\text{M'}}} (Modified exposure affects adjusted Mediator)
723
               4. M' - Y' (Adjusted mediator affects new Outcome)
724
               5. X' - Y' (Direct modified exposure effect)
           - Use prime notation (M', Y') for counterfactual variables
725
           - Maintain Z's original relationships
726
727
```

CRASS Causality & Causal Graph.

```
729
730
          "factual_roles": {
731
            "Exposure": ["act of opening treasure chest"],
732
            "Covariate": ["key possession", "physical capability"],
733
            "Mediator": ["lock mechanism release"],
734
            "Outcome": ["chest opened"]
735
          "counterfactual_roles": {
736
            "Exposure": ["omission of opening action"],
737
            "Covariate": ["key possession", "physical capability"],
            "Mediator": ["lock state preservation"],
739
            "Outcome": ["chest remains closed"]
740
          "causal_graph": {
741
            "factual_edges": [
742
               ["key possession", "act of opening treasure chest"],
743
               ["key possession", "lock mechanism release"], ["key possession", "chest opened"],
744
               ["physical capability", "act of opening treasure chest"], \[
745
               ["physical capability", "lock mechanism release"], ["physical capability", "chest opened"],
746
747
               ["act of opening treasure chest", "lock mechanism release"],
748
               ["lock mechanism release", "chest opened"],
749
               ["act of opening treasure chest", "chest opened"],
750
            "counterfactual_edges": [
751
               ["key possession", "lock state preservation"], ["key possession", "chest remains closed"],
752
753
               ["physical capability", "lock state preservation"],
["physical capability", "chest remains closed"],
754
               ["omission of opening action", "lock state preservation"],
755
               ["lock state preservation", "chest remains closed"],
```

757

758

759

760

761

762

763

764

765

766

767

768

769 770

771

772

773

774

775

776

777

778

779

782

783

784

785

787 788 789

790 791

792

793

794

795

797

800

801

802

803

804 805

806 807

808

```
["omission of opening action", "chest remains closed"],
]
}
```

CLOMO Example. The Counterfactual Logical Modification (CLOMO) dataset is designed to evaluate whether large language models can perform controlled, counterfactual edits to natural language arguments in a logically coherent way. Each example presents a base argument and two premises: Premise 1 has a logical sensitivity with the original argument, while Premise 2 does not. The model is instructed to modify the argument such that Premise 2 has a logical sensitivity with the original argument, while Premise 1 no longer is. For instance, given an argument attributing the rise in gasoline prices fully to government policies, the model must produce a revised version (e.g., changing "fully responsible" to "partly leads") of an argument that shifts logical sensitivity from one premise to another without introducing new claims. The following displays a full example:

```
"instruction": "In the following, you will see an argument and 2
   premises, where Premise 1 provides a necessary assumption to the
   Argument. Please modify the Statements in the Argument until
   Premise 2 provides a necessary assumption to the Argument instead,
    while Premise 1 fails to provide a necessary assumption to the
   Argument. Note that no additional statement should be added. ",
"input": "Argument: Statement1: Consumer advocate : there is no doubt
   that the government is responsible for the increased cost of
   gasoline, because the government's policies have significantly
   increased consumer demand for fuel, and as a result of increasing
   demand, the price of gasoline has risen steadily. Premise1: The
   government can bear responsibility for that which it indirectly
   causes.Premise2: Consumer demand for gasoline cannot increase
   without causing gasoline prices to increase. Please write the
   modified argument below: ",
"output": "Statement1: Consumer advocate : there is no doubt that the
   government partly leads to the increased cost of gasoline, because
    the government's policies have significantly increased consumer
   demand for fuel, and as a result of increasing demand, the price
   of gasoline has risen steadily undoubtedly."
```

CLOMO Causality & Causal Graph.

```
"factual_roles": {
  "Exposure": ["Premise 1 as necessary assumption"],
  "Covariate": [
    "Government's policy impact on demand",
    "Demand-price relationship assumption"
 "Mediator": ["Causal attribution mechanism (direct vs indirect)"],
  "Outcome": ["Full responsibility attribution to government"]
"counterfactual_roles": {
  "Exposure": ["Premise 2 as necessary assumption"],
  "Covariate": [
    "Government's policy impact on demand",
    "Demand-price relationship assumption"
  "Mediator": ["Responsibility attribution modifier (partial vs full)"]
  "Outcome": ["Partial responsibility attribution to government"]
"causal_graph": {
  "factual_edges": [
    ["Government's policy impact on demand", "Premise 1 as necessary
       assumption"],
```

```
810
             ["Government's policy impact on demand", "Causal attribution
811
                mechanism"],
812
             ["Government's policy impact on demand", "Full responsibility
813
                 attribution to government"],
             ["Demand-price relationship assumption", "Premise 1 as necessary
814
                assumption"],
815
             ["Demand-price relationship assumption", "Causal attribution
816
                mechanism"],
817
             ["Demand-price relationship assumption", "Full responsibility
818
                attribution to government"],
             ["Premise 1 as necessary assumption", "Causal attribution mechanism
819
                "],
820
             ["Causal attribution mechanism", "Full responsibility attribution
821
                to government"],
822
             ["Premise 1 as necessary assumption", "Full responsibility
                attribution to government"]
823
824
         "counterfactual_edges": [
825
             ["Government's policy impact on demand", "Responsibility
826
                 attribution modifier"],
827
             ["Government's policy impact on demand", "Partial responsibility
828
                attribution to government"],
             ["Demand-price relationship assumption", "Responsibility
829
                attribution modifier"],
830
             ["Demand-price relationship assumption", "Partial responsibility
831
                attribution to government"],
832
             ["Premise 2 as necessary assumption", "Responsibility attribution
833
                modifier"],
             ["Responsibility attribution modifier", "Partial responsibility
834
                attribution to government"],
             ["Premise 2 as necessary assumption", "Partial responsibility
836
                 attribution to government"]
837
838
         }
839
840
```

RNN-Typology Example. This is a synthetic dataset contains sentence pairs that reflect syntactic alterations to word orders (e.g., converting English from subject-verb-object (SVO) to subject-object-verb (SOV) order). For example, the factual sentence "*Tim saw Lucas*." (SVO) is transformed to its SOV equivalent "*Tim Lucas saw*.".

```
"tim saw lucas.": "tim lucas saw."
```

RNN-Typology Causality & Causal Graph.

841

842

843

844 845

846 847

848 849

850

851

852

853 854

855

856

857

858 859

860

861

862

```
"factual_roles": {
    "Exposure": ["subject-verb-object order"],
    "Covariate": ["syntactic rule", "Lexical items (Tim, saw, Lucas)"],
    "Mediator": ["SOV reordering operation"],
    "Outcome": ["tim saw lucas."]
},
    "counterfactual_roles": {
        "Exposure": ["subject-object-verb order"],
        "Covariate": ["syntactic rule", "Lexical items (Tim, saw, Lucas)"],
        "Mediator": ["SVO restoration operation"],
        "Outcome": ["tim lucas saw."]
},
    "causal_graph": {
        "factual_edges": [
        ["syntactic rule", "subject-verb-object order"],
        ["syntactic rule", "SOV reordering operation"],
        ["syntactic rule", "tim saw lucas."],
        ["Lexical items (Tim, saw, Lucas)", "subject-verb-object order"],
```

```
864
             ["Lexical items (Tim, saw, Lucas)", "SOV reordering operation"],
865
             ["Lexical items (Tim, saw, Lucas)", "tim saw lucas."],
866
             ["subject-verb-object order", "SOV reordering operation"],
867
             ["SOV reordering operation", "tim saw lucas."],
            ["subject-verb-object order", "tim saw lucas."]
868
869
          "counterfactual_edges": [
870
            ["syntactic rule", "SVO restoration operation"], ["syntactic rule", "tim lucas saw."],
871
            ["Lexical items (Tim, saw, Lucas)", "SVO restoration operation"], ["Lexical items (Tim, saw, Lucas)", "tim lucas saw."],
872
873
            ["subject-object-verb order", "SVO restoration operation"],
874
             ["SVO restoration operation", "tim lucas saw."],
875
             ["subject-object-verb order", "tim lucas saw."]
876
877
878
879
```

CVQA-Bool Example. Counterfactual Visual Question Answering(CVQA) is designed to assess the ability of vision-language models to perform counterfactual reasoning over images. Each example presents a factual visual query-answer pair (e.g., "Is there a red sandal here?" \rightarrow yes) grounded in a COCO image, along with a corresponding counterfactual query that modifies a key visual condition (e.g., "Would there be a red sandal here if all shoes were removed?" \rightarrow no). The task requires the model to infer changes in object presence or relationships under hypothetical alterations to the scene. The dataset focuses on a boolean query type. The following displays an example:

| real image | query | answer | new query | new answer | type |
|------------|-----------------------------|--------|---|------------|---------|
| | Is there a red sandal here? | yes | Would there be a red sandal here if all shoes were removed? | no | boolean |

CVQA-Bool Causality & Causal Graph.

880

881

882

883

884

885

893 894 895

896

897

898

899

900

902

903

904

905

906

907

908

909

910

911

912 913

914

915

916

```
"factual_roles": {
  "Exposure": ["presence of red sandal"],
  "Covariate": ["original shoe collection in cart", "visual recognition
       capability"],
  "Mediator": ["sandal-as-shoe categorical inclusion"],
  "Outcome": ["yes"]
"counterfactual_roles": {
  "Exposure": ["removal of all shoes"],
  "Covariate": ["original shoe collection in cart", "visual recognition
       capability"],
  "Mediator": ["sandal-shoe categorical dependency"],
  "Outcome": ["no"]
"causal_graph": {
  "factual_edges": [
    ["original shoe collection in cart", "presence of red sandal"],
    ["original shoe collection in cart", "sandal-as-shoe categorical
        inclusion"],
    ["original shoe collection in cart", "yes"],
    ["visual recognition capability", "presence of red sandal"], ["visual recognition capability", "sandal-as-shoe categorical
        inclusion"],
    ["visual recognition capability", "yes"],
    ["presence of red sandal", "sandal-as-shoe categorical inclusion"],
    ["sandal-as-shoe categorical inclusion", "yes"],
    ["presence of red sandal", "yes"]
 ],
```

CVQA-Count Example. Visual Counterfactual Query Dataset (CVQA) also evaluates whether language models can perform a direct or indirect numerical counterfactual reasoning grounded in visual inputs. Each example consists of a factual visual question (e.g., "How many plates are there?" \rightarrow 1) paired with a corresponding counterfactual query that modifies the quantity in a clearly defined way (e.g., "How many plates would there be if 2 more plates were added?" \rightarrow 3). The model must integrate visual perception (e.g., detecting a single white plate in an image) with numerical logic (e.g., adding 2) to produce the correct answer. The dataset focuses on a counting query type. The following displays an example:

| real image | query | answer | new query | new answer | type |
|------------|---------------------------|--------|---|------------|-----------------|
| | How many plates are there | 1 | How many plates would there be if 2 more plates were added? | 3 | direct counting |

CVQA-Count Causality & Causal Graph.

```
"factual_roles": {
  "Exposure": ["current plate presence (1 unit)"],
  "Covariate": ["original plate count (1)", "visual counting capability
     "],
  "Mediator": ["visual plate detection mechanism"],
  "Outcome": ["1"]
"counterfactual_roles": {
  "Exposure": ["addition of 2 plates"],
  "Covariate": ["original plate count (1)", "visual counting capability
  "Mediator": ["numerical addition operation"],
  "Outcome": ["3"]
"causal_graph": {
  "factual edges": [
    ["original plate count (1)", "current plate presence (1 unit)"],
    ["original plate count (1)", "visual plate detection mechanism"],
    ["original plate count (1)", "1"],
    ["visual counting capability", "current plate presence (1 unit)"],
    ["visual counting capability", "visual plate detection mechanism"],
    ["visual counting capability", "1"],
    ["current plate presence (1 unit)", "visual plate detection
       mechanism"],
    ["visual plate detection mechanism", "1"],
    ["current plate presence (1 unit)", "1"]
  "counterfactual_edges": [
```

```
["original plate count (1)", "numerical addition operation"],
    ["original plate count (1)", "3"],
    ["visual counting capability", "numerical addition operation"],
    ["visual counting capability", "3"],
    ["addition of 2 plates", "numerical addition operation"],
    ["numerical addition operation", "3"],
    ["addition of 2 plates", "3"]
    ]
}
```

COCO Example. Common Objects in Context(COCO) dataset provides automatically constructed counterfactual examples for evaluating multimodal reasoning in image-text pairs. Each instance contains two images and two near-identical captions that differ only in a key noun (e.g., "A big burly grizzly bear is shown with grass in the background" vs. "A big burly grizzly bear is shown with deer in the background"). The dataset is designed to test whether models can detect minimal semantic changes and determine whether the new image visually aligns with the counterfactual caption. The goal is to assess visual-textual consistency and a model's sensitivity to causal or identity-based alterations in structured multimodal contexts. The following displays an example:

| Factual Caption | Image 0 | Counterfactual Caption | Image 1 |
|---|---------|--|---------|
| A big burly grizzly bear is shown with grass in the background. | | A big burly grizzly bear is shown with deer in the background. | |

COCO Causality & Causal Graph.

```
999
1000
          "factual_roles": {
1001
             "Exposure": ["original image (bear with grass)"],
1002
             "Covariate": ["bear presence", "background context"],
1003
             "Mediator": ["grass visual detection"],
1004
             "Outcome": ["caption with 'grass'"]
1005
          "counterfactual_roles": {
1006
             "Exposure": ["modified image (bear with deer)"],
             "Covariate": ["bear presence", "background context"],
1008
             "Mediator": ["deer-grass substitution mechanism"],
1009
             "Outcome": ["caption with 'deer'"]
1010
          "causal_graph": {
1011
             "factual_edges": [
1012
               ["bear presence", "original image (bear with grass)"],
1013
               ["bear presence", "grass visual detection"], ["bear presence", "caption with 'grass'"],
1014
               ["background context", "original image (bear with grass)"], ["background context", "grass visual detection"], ["background context", "caption with 'grass'"],
1015
1016
1017
               ["original image (bear with grass)", "grass visual detection"],
1018
               ["grass visual detection", "caption with 'grass'"],
1019
               ["original image (bear with grass)", "caption with 'grass'"]
1020
             "counterfactual_edges": [
1021
               ["bear presence", "deer-grass substitution mechanism"],
1022
               ["bear presence", "caption with 'deer'"],
1023
               ["background context", "deer-grass substitution mechanism"], ["background context", "caption with 'deer'"],
1024
               ["modified image (bear with deer)", "deer-grass substitution
1025
                   mechanism"],
```

1032

1033

1034

1035

1036

1037

1038

1039 1040

1041

1042

1043 1044

1045

1046

1047

1048 1049 1050

1051

```
["deer-grass substitution mechanism", "caption with 'deer'"],

["modified image (bear with deer)", "caption with 'deer'"]

1028

]
1029
}
1030
}
```

Arithmetic Example. Base-computation Arithmetic dataset evaluates counterfactual numerical reasoning by testing arithmetic operations across multiple numeral systems(e.g. Base-8, 9, 10, 11, 16). Each example pairs a factual base-10 calculation with a counterfactual alternate-base computation (e.g., base-8: $14_8+57_8=73_8$, base-16: $EC_{16}+DD_{16}=1C9_{16}$). The dataset includes inputs (num1, num2), the numeral system (e.g., "8" for octal, "16" for hexadecimal), and the base-specific result (addrst). It assesses models' ability to adapt numeral system transitions and consistency in counterfactual reasoning. The following display a base-8 computation example and a base-16 example:

```
{
   "8": {
     "num1": "14",
     "num2": "57",
     "addrst": "73"
   },
   "16": {
     "num1": "EC",
     "num2": "DD",
     "addrst": "1C9"
   }
}
```

Base-8 Arithmetic Causality & Causal Graph.

```
1052
1053
         "factual_roles": {
           "Exposure": ["10-based system"],
1054
           "Covariate": ["14", "57"],
1055
           "Mediator": ["base-10 arithmetic operation"],
1056
           "Outcome": ["71"]
1057
1058
         "counterfactual_roles": {
           "Exposure": ["8-based system"],
1059
           "Covariate": ["14", "57"],
1060
           "Mediator": [
1061
             "base-8 to base-10 conversion",
1062
              "base-10 sum conversion to base-8"
1063
           1,
           "Outcome": ["73"]
1065
         "causal_graph": {
           "factual_edges": [
              ["14", "10-based system"],
              ["14", "base-10 arithmetic operation"],
1068
             ["14", "71"],
["57", "10-based system"],
1069
1070
              ["57", "base-10 arithmetic operation"],
1071
             ["57", "71"],
1072
             ["10-based system", "base-10 arithmetic operation"],
1073
              ["base-10 arithmetic operation", "71"],
              ["10-based system", "71"]
1074
1075
           "counterfactual_edges": [
1076
             ["14", "base-8 to base-10 conversion"],
1077
             ["14", "73"],
["57", "base-8 to base-10 conversion"],
1078
              ["57", "73"],
1079
              ["8-based system", "base-8 to base-10 conversion"],
```

```
1080
1081
1082
1082
1083
1084
1085
1086
["base-8 to base-10 conversion", "base-10 sum conversion to base-8"
],
["base-10 sum conversion to base-8", "73"],
["8-based system", "73"]
]
1086
]
1087
1088
]
1088
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
1089
]
```

Base-16 Arithmetic Causality & Causal Graph.

1087

1088

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

```
1089
1090
         "factual_roles": {
           "Exposure": ["10-based system"],
1091
           "Covariate": ["EC", "DD"],
1092
           "Mediator": ["N.A."],
1093
           "Outcome": ["N.A."]
1094
1095
         "counterfactual_roles": {
           "Exposure": ["16-based system"],
1096
           "Covariate": ["EC", "DD"],
1097
           "Mediator": [
1098
             "hex-to-decimal conversion",
1099
              "decimal-to-hex reversion"
1100
           "Outcome": ["1C9"]
1101
1102
         "causal_graph": {
1103
           "factual_edges": [
1104
              ["EC", "10-based system"],
              ["DD", "10-based system"]
1105
1106
           "counterfactual_edges": [
1107
              ["EC", "hex-to-decimal conversion"],
1108
              ["hex-to-decimal conversion", "decimal-to-hex reversion"],
1109
             ["EC", "1C9"],
["DD", "hex-to-decimal conversion"],
1110
              ["DD", "1C9"],
1111
              ["16-based system", "hex-to-decimal conversion"],
1112
             ["decimal-to-hex reversion", "1C9"],
1113
              ["16-based system", "1C9"]
1114
           1
1115
         }
1116
1117
```

MalAlgoQA Example (Malformed Algorithmic Question Answering (MalAlgoQA) dataset is designed intentionally including factual and counterfactual rationales between multiple-choice question answering to validate a language model's ability to discern sound reasoning in the presence of rationales(factual or counterfactual). Each question is presented alongside a factual rationale that supports the correct answer (e.g., "Correctly ordered the values from greatest to least: 276, 254, 237, 235." \rightarrow C), and is paired with counterfactual rationales (e.g., "Ordered least to greatest") that correspond to plausible but incorrect or altered answers (e.g., A). Each example is decomposed into factual and counterfactual role pairs, allowing researchers to assess how changes in reasoning paths (rationales) lead to different answer choices. The following display an example of raw data and its decomposed data points:

```
1128
           "Question": "Which list shows the following number in order from
1129
              highest to lowest?",
1130
           "Answer": "C",
1131
           "Choice_A":" 235
                              237
                                    254
                                         276 ",
1132
                                         254 ",
           "Choice_B":" 237
                              276
                                    235
           "Choice_C":" 276
                                         235 ",
1133
                              254
                                   237
           "Choice_D":" 276
                                         237 ",
                              254
                                   235
```

```
1134
           "Rationale_A": "Ordered least to greatest",
1135
           "Rationale_B": "Ordered greatest to least by ones place.",
1136
           "Rationale_C":"Correctly ordered the values from greatest to least:
1137
              276, 254, 237, 235.",
           "Rationale_D": "Switched last 2 numbers."
1138
1139
1140
1141
1142
               "Question": "Which list shows the following number in order from
1143
                  highest to lowest?",
1144
               "Answer": "C",
1145
               "Counterfactual Answer": "A",
               "Choice_A":" 235 237 254 276 ",
1146
               "Choice_B":" 237 276 235 254 ",
1147
               "Choice_C": 276 254 237 235 ",
1148
               "Choice_D": " 276 254 235 237 ",
1149
               "Counterfactual Rationale": "Ordered least to greatest",
               "Rationale_C": "Correctly ordered the values from greatest to
1150
                   least: 276, 254, 237, 235."
1151
           },
1152
1153
           {
1154
               "Question": "Which list shows the following number in order from
1155
                   highest to lowest?",
               "Answer": "C",
1156
               "Counterfactual Answer": "B",
1157
               "Choice_A":" 235 237 254 276 ",
1158
               "Choice_B":" 237
                                 276
                                       235
                                            254 ",
1159
               "Choice_C":" 276
                                            235 ",
                                 254
                                       237
               "Choice_D":" 276 254 235 237 ",
1160
               "Counterfactual Rationale": "Ordered greatest to least by ones
1161
                   place",
1162
               "Rationale_C": "Correctly ordered the values from greatest to
1163
                   least: 276, 254, 237, 235."
1164
           } .
1165
           {
1166
               "Question": "Which list shows the following number in order from
1167
                  highest to lowest?",
1168
               "Answer":"C",
1169
               "Counterfactual Answer": "D",
1170
               "Choice_A":" 235 237 254 276 ",
               "Choice_B":" 237 276 235 254 ",
1171
               "Choice_C":" 276 254 237
                                           235 ",
1172
               "Choice_D":" 276 254 235 237 ",
1173
               "Rationale_C":"Correctly ordered the values from greatest to
1174
                   least: 276, 254, 237, 235.",
               "Counterfactual Rationale": "Switched last 2 numbers."
1175
           }
1176
1177
```

Take the first decomposed sample as an example showing MalAlgoQA Causality & Causal Graph.

1178

```
1180
1181
         "factual_roles": {
           "Exposure": ["Ordered from greatest to least"],
1182
           "Covariate": ["Number set (276, 254, 237, 235)"],
1183
           "Mediator": ["Descending comparison logic"],
1184
           "Outcome": ["Choice C"]
1185
1186
         "counterfactual_roles": {
           "Exposure": ["Ordered least to greatest"],
1187
           "Covariate": ["Number set (276, 254, 237, 235)"],
```

```
1188
            "Mediator": ["Ascending comparison logic"],
1189
            "Outcome": ["Choice A"]
1190
1191
          "causal_graph": {
            "factual_edges": [
1192
              ["Number set (276, 254, 237, 235)", "Ordered from greatest to least
1193
                  "],
1194
              ["Number set (276, 254, 237, 235)", "Descending comparison logic"], ["Number set (276, 254, 237, 235)", "Choice C"],
1195
              ["Ordered from greatest to least", "Descending comparison logic"],
1196
              ["Descending comparison logic", "Choice C"],
1197
              ["Ordered from greatest to least", "Choice C"]
1198
1199
            "counterfactual_edges": [
1200
              ["Number set (276, 254, 237, 235)", "Ordered least to greatest"],
              ["Number set (276, 254, 237, 235)", "Ascending comparison logic"],
1201
              ["Number set (276, 254, 237, 235)", "Choice A"],
1202
              ["Ordered least to greatest", "Ascending comparison logic"],
1203
              ["Ascending comparison logic", "Choice A"], ["Ordered least to greatest", "Choice A"]
1204
1205
            1
         }
1207
```

HumanEval-Exe Example This dataset performs a programming-related task: code execution. It is designed to probe the ability of code-execution language models to perform counterfactual reasoning in the context of program behavior. Each example consists of a function definition and a test case input, and the model is asked to predict the output under a factual assumption (e.g., Python's default 0-based indexing). The example is paired with a counterfactual version of the same test case, where a hypothetical condition is introduced—such as switching to 1-based indexing. The model must then predict the corresponding counterfactual output. For instance, given a function that checks for close floating-point elements in a list, the model is expected to reason whether the list and threshold would yield a different outcome if indexing conventions were altered. The full example is shown as follows:

```
"instruction": "from typing import List\n\n\ndef has_close_elements(
    numbers: List[float], threshold: float) -> bool:\n \"\"\"
    Check if in given list of numbers, are any two numbers closer to
    each other than\n given threshold.\n >>> has_close_elements
    ([1.0, 2.0, 3.0], 0.5)\n False\n >>> has_close_elements([1.
        0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3)\n True\n \"\"\"",
    "input": "[1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.3",
    "output": "True",
    "counterfactual_output": "True"
}
```

HumanEval-Exe Causality & Causal Graph.

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217 1218

1219

1220

1221

1222

1224

1225

```
1229
1230
         "factual_roles": {
           "Exposure": ["0-based indexing"],
1231
           "Covariate": [
1232
             "List values [1.0, 2.0, 3.9, 4.0, 5.0, 2.2]",
1233
             "Threshold 0.3",
1234
             "Pairwise comparison algorithm"
1235
           "Mediator": ["Range iteration logic (0 <= i < j < len(numbers))"],
1236
           "Outcome": ["True"]
1237
         "counterfactual_roles": {
1239
           "Exposure": ["1-based indexing"],
1240
           "Covariate": [
             "List values [1.0, 2.0, 3.9, 4.0, 5.0, 2.2]",
1241
             "Threshold 0.3",
```

1268

1269

1270

1271

1272

1273

1274

1275

```
1242
             "Pairwise comparison algorithm"
1243
           ],
1244
           "Mediator": ["Range iteration logic (1 <= i < j <= len(numbers))"],
1245
           "Outcome": ["True"]
         },
1246
         "causal_graph": {
1247
           "factual_edges": [
1248
             ["List values [...]", "0-based indexing"],
1249
             ["Threshold 0.3", "0-based indexing"],
1250
             ["Pairwise comparison algorithm", "0-based indexing"],
             ["0-based indexing", "Range iteration logic (0 <= i < j < len(
1251
                 numbers))"],
1252
             ["Range iteration logic (0 <= i < j < len(numbers))", "True"],
1253
             ["List values [...]", "True"],
1254
             ["Threshold 0.3", "True"]
1255
           ],
           "counterfactual_edges": [
1256
             ["List values [...]", "1-based indexing"],
1257
             ["Threshold 0.3", "1-based indexing"],
1258
             ["Pairwise comparison algorithm", "1-based indexing"],
             ["1-based indexing", "Range iteration logic (1 <= i < j <= len(
1259
                 numbers))"],
1260
             ["Range iteration logic (1 <= i < j <= len(numbers))", "True"],
1261
             ["List values [...]", "True"],
1262
             ["Threshold 0.3", "True"]
1263
           1
1264
1265
1266
```

Open-Critic Example This dataset performs a programming-related task: code generation. It is a synthetic code editing benchmark designed to train and evaluate large language models on their ability to identify and fix bugs in code. Each example consists of a natural language *task description*, a *correct code* solution that satisfies the task, and a *counterfactual explanation* that describes bugs introduced into a similar but faulty implementation. The objective is for the model to learn to generate or repair the *bugged version of the code* based on the bug descriptions as a counterfactual scenario. For instance, given a correct nested loop that prints all valid digit combinations excluding the number 5 and repetitions, the counterfactual explanation points out altered range values (e.g., using range (100) and range (1)), leading to an incomplete iteration logic. The full example is shown as follows:

```
1276
1277
          "task": "Create a nested loop to print every combination of numbers
1278
              between 0-9, excluding any combination that contains the number
1279
              5. Additionally, exclude any combination that contains a
1280
              repeating digit. Implement the solution without using any built-
              in functions or libraries to check for repeating digits.",
1281
          "correct_code": "for i in range(10): # First digit\n
                                                                   for j in
1282
              range(10): # Second digit\n for k in range(10): # Third digit\
1283
                   # Checking for the conditions\n
                                                     if i != 5 and j != 5 and k
1284
               != 5 and i != j and i != k and j != k:\n print(i, j, k)",
          "correct_explanation": "This code will generate and print every
1285
              combination of three digits between 0-9 that do not contain the
1286
              number 5 and do not have any repeating digits.",
1287
          "bugged_code": "for i in range(10): \n
                                                      for j in range(100):
1288
                                      if i != 5 and j != 5 and k != 5 and i !=
              for k in range(1): \n
1289
               j and i != k and j != k:\n
                                             print(i, j, k)",
          "counterfactual_explanation": "Reading through this code, I notice
1290
              that:\n\n* The second loop is set to 'range(100)' instead of '
1291
              range(10)', which means it will only iterate over the numbers
              from 0 to 99. This will cause the script to miss some
1293
              combinations.\n* The third loop is set to 'range(1)', which means
1294
               it will only iterate over the number 0. This will cause the
1295
              script to only print combinations with one digit, instead of
              three.\n\nThese bugs will prevent the script from generating and
```

```
printing all possible combinations of three digits between 0-9
that do not contain the number 5 and do not have any repeating
digits.\n\nTips for avoiding these mistakes:\n\n* Double-check
the range values in each loop to ensure they are correct.\n* Make
    sure the loops iterate correctly over the desired range of
    values."
}
```

Open-Critic Causality & Causal Graph.

1296

1297

1298

1299

1300

```
1304
1305
           "factual_roles": {
1306
               "Exposure": ["Correct explanation (valid ranges and checks)"],
1307
               "Covariate": [
1308
                 "Task requirements (0-9 digits)",
                 "Exclusion logic (no 5/repeats)",
1309
                 "Nested loop structure"
1310
1311
               "Mediator": ["Proper range initialization (range(10) x3)"],
1312
               "Outcome": ["Correct triple-nested loop code"]
1313
1314
           "counterfactual_roles": {
               "Exposure": ["Counterfactual explanation (invalid ranges)"],
1315
               "Covariate": [
1316
                 "Task requirements (0-9 digits)",
1317
                 "Exclusion logic (no 5/repeats)",
1318
                 "Nested loop structure"
1319
               "Mediator": [
1320
                 "Flawed range parameters (range(100)/range(1))",
                 "Incomplete digit iteration"
1322
               1.
1323
               "Outcome": ["Bugged code with limited iterations"]
1324
               },
           "causal_graph": {
1325
               "factual_edges": [
1326
                 ["Task requirements", "Correct explanation"],
1327
                 ["Exclusion logic", "Correct explanation"],
1328
                 ["Nested loop structure", "Correct explanation"],
                 ["Task requirements", "Proper range initialization"],
1329
                 ["Exclusion logic", "Proper range initialization"],
1330
                 ["Nested loop structure", "Proper range initialization"],
1331
                 ["Correct explanation", "Correct triple-nested loop code"],
1332
                 ["Task requirements", "Correct triple-nested loop code"],
1333
                 ["Exclusion logic", "Correct triple-nested loop code"],
1334
                 ["Correct explanation", "Proper range initialization"],
                 ["Proper range initialization", "Correct triple-nested loop
1335
                     code"],
1336
                 ["Correct explanation", "Correct triple-nested loop code"]
1337
               ],
1338
               "counterfactual_edges": [
                 ["Task requirements", "Flawed range parameters"],
1339
                 ["Exclusion logic", "Flawed range parameters"],
1340
                 ["Nested loop structure", "Flawed range parameters"],
1341
                 ["Task requirements", "Bugged code with limited iterations"],
1342
                 ["Exclusion logic", "Bugged code with limited iterations"],
1343
                 ["Nested loop structure", "Bugged code with limited iterations"
1344
                     ],
                 ["Counterfactual explanation", "Flawed range parameters"],
1345
                 ["Flawed range parameters", "Incomplete digit iteration"],
1346
                 ["Incomplete digit iteration", "Bugged code with limited
1347
                     iterations"],
1348
                 ["Counterfactual explanation", "Bugged code with limited
1349
                     iterations"]
```

1354

1355

1356

1357

1358

1359

1360

1361

1362

1396

1398 1399

1400

1401

1402

1403

```
}
```

Code-Preference Example This dataset performs a programming-related task: code summarization. It contains pairs of duplicate code examples, with the only difference being the bugged code example has the bugged code 'surgically transplanted in' while the corrected code is left the same. Each example consists of a natural language *instruction*, a *correct code* solution that satisfies the instruction, and a *bug explanation* that describes bugs. The objective is for the model to learn to summarize and generate the bug descriptions as a counterfactual scenario. For instance, given a correct nested loop that prints all valid digit combinations excluding the number 5 and repetitions compared with a bugged loop, the bug description generated by the model will be able to point out altered range values (e.g., using range (100) and range (1)) leading to an incomplete iteration logic in its summarized response of bug explanation. The full example is shown as follows:

```
1364
          "Instruction": 'Create a nested loop to print every combination of
1365
              numbers between 0-9, excluding any combination that contains the
              number 5. Additionally, exclude any combination that contains a
1366
              repeating digit. Implement the solution without using any built-
1367
              in functions or libraries to check for repeating digits.',
          "bugged_code": 'What are the problems with this code? '''\npython\
1369
                                          for j in range(100): \n for k in
              nfor i in range(10): \n
              range(1): \n if i != 5 and j != 5 and k != 5 and i != j and i !=
1370
               k and j != k:\n print(i, j, k)\n''',
1371
          "bug_explanation": 'Reading through this code, I notice that:\n\n*
1372
              The second loop is set to 'range(100)' instead of 'range(10)',
1373
              which means it will only iterate over the numbers from 0 to 99.
1374
              This will cause the script to miss some combinations.\n* The
1375
              third loop is set to 'range(1)', which means it will only iterate
               over the number 0. This will cause the script to only print
1376
              combinations with one digit, instead of three.\n\nThese bugs will
1377
               prevent the script from generating and printing all possible
1378
              combinations of three digits between 0-9 that do not contain the
1379
              number 5 and do not have any repeating digits.\n\nTips for
              avoiding these mistakes:\n\n* Double-check the range values in
              each loop to ensure they are correct.\n* Make sure the loops
1381
              iterate correctly over the desired range of values.\n\nHere is
1382
              the corrected code:\n\n'''python\nfor i in range(10): # First
1383
                         for j in range(10): # Second digit\n
                                                                  for k in range
1384
              (10): # Third digit\n # Checking for the conditions\n if i
              != 5 and j != 5 and k != 5 and i != j and i != k and j != k:\n
1385
                 print(i, j, k) \n'''
1386
          "correct_code": 'Here is an example of a nested loop in Python to
1387
              print every combination of numbers between 0-9, excluding any
1388
              combination that contains the number 5 or repeating digits: \n\
1389
              '''python\nfor i in range(10): # First digit\n
                                                                  for j in range
              (10): # Second digit\n for k in range(10): # Third digit\n #
1390
              Checking for the conditions\n if i != 5 and j != 5 and k != 5
1391
                                                     print(i, j, k) \n'''\n
              and i != j and i != k and j != k:\n
1392
              nThis code will generate and print every combination of three
1393
              digits between 0-9 that do not contain the number 5 and do not
1394
              have any repeating digits.'
1395
```

Code-Preferencec Causality & Causal Graph.

```
{
   "factual_roles": {
      "Exposure": ["Correct code (range(10) loops)"],
      "Covariate": [
      "Task requirements (0-9 digits)",
      "Exclusion logic (no 5/repeats)",
      "Nested loop structure"
```

```
1404
            ],
1405
            "Mediator": ["Proper range initialization (range(10) x3)"],
1406
           "Outcome": ["Correct explanation (valid ranges and checks)"]
1407
         "counterfactual_roles": {
1408
            "Exposure": ["Bugged code (range(100)/range(1))"],
1409
            "Covariate": [
1410
              "Task requirements (0-9 digits)",
1411
              "Exclusion logic (no 5/repeats)",
1412
              "Nested loop structure"
1413
            "Mediator": ["Flawed range parameters", "Incomplete digit iteration"]
1414
1415
            "Outcome": ["Bug explanation (incorrect ranges analysis)"]
1416
          "causal_graph": {
1417
            "factual_edges": [
1418
              ["Task requirements (0-9 digits)", "Correct code (range(10) loops)"
1419
1420
              ["Exclusion logic (no 5/repeats)", "Correct code (range(10) loops)"
1421
              ["Nested loop structure", "Correct code (range(10) loops)"],
1422
              ["Task requirements (0-9 \text{ digits})", "Proper range initialization (
1423
                  range(10) x3)"],
1424
              ["Exclusion logic (no 5/repeats)", "Proper range initialization (
1425
                  range(10) x3)"],
1426
              ["Nested loop structure", "Proper range initialization (range(10)
1427
                  x3)"],
              ["Task requirements (0-9 digits)", "Correct explanation (valid
1428
                  ranges and checks) "],
1429
              ["Exclusion logic (no 5/repeats)", "Correct explanation (valid
1430
                  ranges and checks) "],
1431
              ["Correct code (range(10) loops)", "Proper range initialization (
                  range(10) x3)"],
1432
              ["Proper range initialization (range(10) x3)", "Correct explanation
1433
                   (valid ranges and checks) "],
1434
              ["Correct code (range(10) loops)", "Correct explanation (valid
1435
                  ranges and checks) "]
1436
            "counterfactual_edges": [
1437
              ["Task requirements (0-9 digits)", "Flawed range parameters"], ["Task requirements (0-9 digits)", "Incomplete digit iteration"], ["Task requirements (0-9 digits)", "Bug explanation (incorrect
1438
1439
1440
                  ranges analysis)"],
              ["Exclusion logic (no 5/repeats)", "Flawed range parameters"], ["Exclusion logic (no 5/repeats)", "Incomplete digit iteration"],
1441
1442
              ["Exclusion logic (no 5/repeats)", "Bug explanation (incorrect
1443
                  ranges analysis)"],
1444
              ["Nested loop structure", "Flawed range parameters"],
              ["Nested loop structure", "Incomplete digit iteration"],
1445
1446
              ["Nested loop structure", "Bug explanation (incorrect ranges
                  analysis)"],
1447
              ["Bugged code (range(100)/range(1))", "Flawed range parameters"],
1448
              ["Flawed range parameters", "Incomplete digit iteration"],
1449
              ["Incomplete digit iteration", "Bug explanation (incorrect ranges
1450
                  analysis)"],
1451
              ["Bugged code (range(100)/range(1))", "Bug explanation (incorrect
                  ranges analysis)"]
1452
1453
         }
1454
1455
```

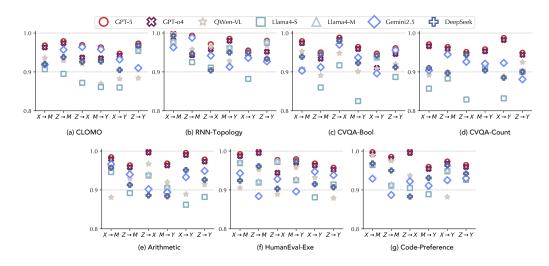


Figure 5: Additional evaluation on causal graph construction, complementing to Figure 3.

B COMPLEMENTARY EXPERIMENT

B.1 SETTING

This section lists the experimental settings used in this study.

Value **Description** Hyperparameter 0.7 Controls randomness in generation Temperature Top-p (nucleus sampling) 0.95 Probability mass for sampling Max tokens Maximum number of tokens to generate Stop sequences ["\n", "Q:"] Used to truncate responses CoT, CoT-SC, ToT Prompting strategy used in Section 4.2.2 Prompt format Tool-calling API Enabled (Selective) Used in tool-augmented experiments

Table 6: LLM query hyperparameters used during all experiments.

Computational Resources. All experiments were conducted on a high-performance computing server equipped with six NVIDIA RTX 6000 Ada Generation GPUs, each with 49 GB of dedicated VRAM. The system utilized CUDA version 12.8 and NVIDIA driver version 570.124.06. These GPUs supported parallel execution of model querying, evaluation, and tool-augmented tasks across our benchmark datasets. The hardware configuration ensured sufficient memory bandwidth and processing capability to accommodate large-scale inference, particularly for multimodal tasks and multi-sample prompting strategies such as CoT-SC and ToT. No resource-related constraints were encountered during experimentation.

B.2 ADDITIONAL RESULT AND DISCUSSION

This part presents additional experimental results that complement the main evaluation in the body of the paper. These supplementary findings, together with what has been presented in previous sections, offer comprehensive insights into LLMs capabilities across different decompositional tasks.

To deepen our understanding of how LLMs handle counterfactual reasoning, we analyze representative datasets that cover text, multimodal, symbolic, and code-based modalities. We aim to uncover impediments at each decompositional stage. Below, we highlight dataset-level characteristics that either enable or hinder performance.

Textual Question-Answering and Logic Parsing. Datasets built on natural language (e.g., textual QA and logical modification tasks) generally facilitate the recognition of explicit causal variables such as exposures and outcomes. Models can easily link a stated intervention ("if tutoring was not

Table 7: LLMs' performance (F1 standard deviations, scaled to 100%) in causal variable identification (reordered; adjusted variability).

| Dataset | GP | T-5 | GP. | Г-о4 | Qw | en3 | Llan | na4-S | Llam | a4-M | Gem | ini2.5 | Deep | Seek |
|-----------------|-------|-----------|--------|----------|-----------|---------|-------|-------|-------|-------|-------|--------|-------|-------|
| Dataset | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 | v_1 | v_2 |
| | | $v_1 = 0$ | X (Exp | osure), | $v_2 = Z$ | (Covari | iate) | | | | | | | |
| CRASS | 4.6 | 4.9 | 9.1 | 9.9 | 5.6 | 6.2 | 8.3 | 7.9 | 3.5 | 4.4 | 4.5 | 5.3 | 4.8 | 5.0 |
| CLOMO | 4.8 | 5.1 | 9.3 | 10.0 | 6.8 | 7.2 | 9.6 | 8.9 | 5.0 | 5.5 | 5.9 | 6.5 | 5.3 | 5.7 |
| RNN-Topo | 5.0 | 5.3 | 9.5 | 10.2 | 7.3 | 7.8 | 10.2 | 9.5 | 5.4 | 5.9 | 6.4 | 6.9 | 5.8 | 6.2 |
| CVQA-Bool | 5.4 | 5.7 | 9.8 | 10.4 | 9.7 | 10.2 | 12.4 | 13.1 | 8.9 | 9.4 | 10.8 | 11.5 | 9.4 | 10.1 |
| CVQA-Count | 5.6 | 5.9 | 9.9 | 10.5 | 10.6 | 11.3 | 13.7 | 14.2 | 9.8 | 10.5 | 11.5 | 12.3 | 10.3 | 11.2 |
| COCO | 5.8 | 6.1 | 10.0 | 10.6 | 11.9 | 12.6 | 15.2 | 15.8 | 11.1 | 11.8 | 12.9 | 13.7 | 11.6 | 12.5 |
| Arithmetic | 5.2 | 5.5 | 9.4 | 10.1 | 8.9 | 9.5 | 11.3 | 12.0 | 7.5 | 8.1 | 8.4 | 9.1 | 7.9 | 8.6 |
| MalAlgoQA | 5.3 | 5.6 | 9.6 | 10.3 | 9.5 | 10.1 | 12.2 | 12.9 | 8.2 | 8.7 | 9.0 | 9.7 | 8.5 | 9.2 |
| HumanEval-Exe | 6.0 | 6.3 | 10.2 | 10.8 | 12.8 | 13.6 | 16.3 | 17.1 | 11.7 | 12.4 | 13.9 | 14.7 | 12.5 | 13.4 |
| Open-Critic | 6.1 | 6.4 | 10.3 | 10.9 | 13.5 | 14.2 | 17.1 | 17.8 | 12.6 | 13.3 | 14.6 | 15.3 | 13.7 | 14.5 |
| Code-Preference | 6.0 | 6.2 | 10.1 | 10.7 | 13.1 | 13.9 | 16.7 | 17.4 | 12.2 | 12.9 | 14.3 | 15.0 | 13.1 | 13.9 |
| | | $v_1 = 0$ | M (Me | diator), | $v_2 = Y$ | (Outco | me) | | | | | | | |
| CRASS | 4.9 | 5.2 | 9.2 | 10.0 | 9.6 | 7.4 | 12.1 | 9.7 | 8.5 | 6.1 | 9.1 | 6.8 | 8.7 | 6.2 |
| CLOMO | 5.0 | 5.3 | 9.4 | 10.2 | 10.5 | 8.1 | 12.9 | 10.3 | 9.2 | 6.6 | 9.8 | 7.5 | 9.3 | 6.8 |
| RNN-Topo | 5.1 | 5.4 | 9.6 | 10.3 | 11.1 | 8.5 | 13.6 | 10.8 | 9.6 | 7.0 | 10.5 | 7.9 | 9.8 | 7.2 |
| CVQA-Bool | 5.7 | 6.0 | 10.0 | 10.6 | 13.8 | 11.3 | 16.9 | 14.2 | 12.4 | 10.0 | 14.1 | 12.5 | 13.2 | 10.7 |
| CVQA-Count | 5.9 | 6.2 | 10.1 | 10.7 | 14.6 | 12.1 | 17.8 | 15.1 | 13.3 | 10.8 | 15.0 | 13.3 | 14.1 | 11.6 |
| COCO | 6.0 | 6.3 | 10.2 | 10.8 | 15.3 | 12.8 | 18.5 | 15.9 | 14.0 | 11.5 | 15.7 | 14.1 | 14.8 | 12.4 |
| Arithmetic | 5.4 | 5.7 | 9.5 | 10.2 | 12.5 | 9.8 | 15.4 | 12.5 | 11.0 | 8.6 | 13.1 | 9.2 | 11.3 | 8.7 |
| MalAlgoQA | 5.6 | 5.9 | 9.7 | 10.4 | 13.1 | 10.4 | 16.1 | 13.2 | 11.6 | 9.2 | 13.8 | 10.0 | 12.0 | 9.3 |
| HumanEval-Exe | 6.2 | 6.5 | 10.4 | 11.0 | 16.2 | 13.6 | 19.3 | 16.7 | 14.8 | 12.3 | 16.5 | 14.8 | 15.5 | 13.1 |
| Open-Critic | 6.3 | 6.6 | 10.5 | 11.1 | 16.9 | 14.4 | 8.1 | 17.5 | 15.5 | 13.0 | 17.2 | 15.5 | 16.1 | 13.9 |
| Code-Preference | 6.1 | 6.4 | 10.3 | 10.9 | 16.5 | 14.0 | 19.7 | 17.0 | 15.2 | 12.6 | 16.9 | 15.1 | 15.8 | 13.5 |

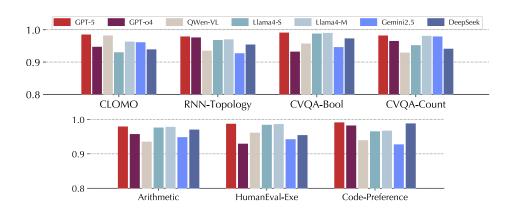


Figure 6: Additional evaluation of LLMs' intervention identification, complementing to Figure 4.

provided") to its corresponding variable. However, mediators are often described abstractly (through latent constructs like "trust," "motivation," or "belief states") rather than explicit entities. This makes Task I disproportionately difficult for mediators compared to exposures or outcomes. Errors at this stage then propagate to Task IV, where the model must simulate how such mediators would affect outcomes under intervention. Even when causal graph construction (Task II) is near-perfect due to the structured nature of logical relations, the absence of explicit mediators limits downstream reasoning fidelity.

Vision-Language Counterfactuals. Multimodal datasets combining images with text pose unique challenges. When asked to identify causal variables, LLMs must ground textual descriptions in visual objects. For example, distinguishing "presence of a ball" (exposure) from "action of kicking" (mediator) requires fine-grained alignment of object attributes with causal semantics. This grounding step introduces errors in Task I, especially when visual scenes are cluttered or ambiguous. Even when interventions (Task III) are identified correctly, outcome reasoning (Task IV) suffers because

Table 8: LLM performance (F1 standard deviation) in reasoning the counterfactual mediator (M') and outcome (Y').

| Dataset | GPT-5 | | GPT-o4 | | Qw | en3 | Llan | na4-S | Llam | a4-M | Gemini2.5 | | DeepSeek | |
|-----------------|-------|-----|--------|-----|------|------|------|-------|------|------|-----------|------|----------|------|
| Dataset | M' | Y' | M' | Y' | M' | Y' | M' | Y' | M' | Y' | M' | Y' | M' | Y' |
| CRASS | 4.2 | 5.0 | 4.8 | 5.6 | 6.7 | 8.5 | 10.0 | 12.5 | 5.0 | 7.1 | 6.3 | 8.1 | 5.5 | 7.4 |
| CLOMO | 4.4 | 5.2 | 5.0 | 5.9 | 7.2 | 9.0 | 10.6 | 13.0 | 5.4 | 7.5 | 6.7 | 8.6 | 6.0 | 7.9 |
| RNN-Topo | 4.6 | 5.4 | 5.1 | 6.0 | 7.8 | 9.5 | 11.0 | 13.5 | 5.6 | 7.7 | 7.2 | 9.1 | 6.5 | 8.4 |
| CVQA-Bool | 4.8 | 5.6 | 5.2 | 6.1 | 11.2 | 13.3 | 14.5 | 16.8 | 9.0 | 11.1 | 12.1 | 14.0 | 10.4 | 12.5 |
| CVQA-Count | 5.0 | 5.8 | 5.3 | 6.3 | 12.0 | 14.1 | 15.2 | 17.6 | 9.8 | 11.9 | 12.9 | 14.9 | 11.2 | 13.3 |
| COCO | 5.1 | 6.0 | 5.4 | 6.4 | 12.7 | 14.8 | 15.9 | 18.3 | 10.5 | 12.7 | 13.6 | 15.6 | 11.9 | 14.0 |
| Arithmetic | 4.5 | 5.2 | 5.1 | 6.0 | 9.5 | 11.6 | 12.8 | 15.3 | 7.4 | 9.5 | 9.0 | 11.1 | 8.3 | 10.4 |
| MalAlgoQA | 4.7 | 5.3 | 5.2 | 6.1 | 10.1 | 12.2 | 13.4 | 15.9 | 7.9 | 10.0 | 9.7 | 11.8 | 8.9 | 11.0 |
| HumanEval-Exe | 5.2 | 6.2 | 5.5 | 6.5 | 13.3 | 15.4 | 16.5 | 18.9 | 11.0 | 13.2 | 14.2 | 16.2 | 12.5 | 14.7 |
| Open-Critic | 5.3 | 6.3 | 5.6 | 6.6 | 13.9 | 16.0 | 17.0 | 19.6 | 11.6 | 13.8 | 14.8 | 16.9 | 13.1 | 15.3 |
| Code-Preference | 5.1 | 6.1 | 5.5 | 6.5 | 13.6 | 15.7 | 16.7 | 19.2 | 11.3 | 13.5 | 14.5 | 16.5 | 12.8 | 15.0 |

models struggle to propagate visual changes into numerical or behavioral predictions. For instance, recognizing that "removing a ball" should reduce the count of possible goals involves chaining visual detection, counting logic, and causal propagation—steps that current LLMs rarely integrate coherently.

Symbolic and Mathematical Reasoning. Datasets built on arithmetic or algorithmic transformations highlight another bottleneck: reliance on memorized patterns instead of causal mechanisms. In variable identification (Task I), explicit quantities are correctly recognized. In causal graph construction (Task II), rules linking operations (e.g., "base conversion influences the final number") are applied consistently. However, in outcome reasoning (Task IV), models frequently fail to simulate the correct causal pathway, often defaulting to template-based responses rather than computing the actual counterfactual result. This suggests that while symbolic data supports high precision in explicit structure, it exposes the weakness of LLMs in mechanistic simulation of causal processes.

Code-Based Reasoning. Programming-oriented datasets such as code execution, code generation, or preference tasks are particularly difficult across all stages. In Task I, identifying exposures and outcomes is hindered by the abstractness of programming constructs (e.g., "function signature" as exposure, "program output" as outcome). Mediators (such as intermediate execution states) are even harder to capture, as they are not explicitly represented in the code text but must be inferred from semantics. In Task II, while models can generate plausible causal graphs describing dependencies among variables or functions, these graphs often overgeneralize or miss critical execution details. Task IV is especially challenging: even when interventions like "changing a loop to recursion" are recognized, LLMs often fail to simulate downstream program behavior, producing outcomes that are logically plausible but incorrect. This reflects a persistent gap between syntactic recognition and semantic reasoning.

Cross-Cutting Observations. Across modalities, two key impediments recur:

- (1) Complex modalities impede variable identification. Images and code introduce higher error rates in Task I, since grounding or semantic parsing must precede causal reasoning.
- (2) Implicit mediators bottleneck outcome reasoning. Regardless of modality, when mediators are abstract or not explicitly present, performance in Task IV drops substantially. LLMs can identify interventions reliably, but they fail to propagate their effects along causal chains to yield consistent counterfactual outcomes.

These findings suggest that LLMs are "eligible" for decompositional reasoning in structured settings with explicit causal variables (e.g., clean text or arithmetic). However, when confronted with modality complexity or implicit causal pathways, their reasoning capacity is significantly impeded.

B.3 WORKING MEMORY PERSPECTIVE TO INTEPRET COUNTERFACTUAL REASONING

Prior research (Zhang et al., 2024a) has demonstrated that language models exhibit notable difficulty in temporally storing and manipulating information even in n-back tasks that are cognitively simpler than explicit reasoning. This underlying limitation in working memory capacity poses constraints

on long-term and multi-step reasoning. To explore the connection between memory bottlenecks and mediator identification challenges, we conducted additional experiments in more depth.

Experiments on Working Memory. To examine how working memory affects mediator reasoning, we designed a controlled n-back Mediator Recall task. While most benchmarks involve only single-step mediation, the Open-Critic dataset (code modality) includes examples of multi-step causal mediation. For instance, adjusting code inputs requires reasoning over prior inputs and transformations. In this task, the mediator must be inferred from causal variables presented n steps earlier in the input. We vary n from 1 to 3 and report F1 scores consistent with Table 9.

Table 9: LLM performance (F1) in n-back mediator recall

| Model | 1-hop | 2-hop | 3-hop |
|--------------|-------|-------|-------|
| GPT-o4 | 72.2% | 63.5% | 9.7% |
| Qwen | 58.3% | 39.6% | 12.1% |
| Gemini | 66.4% | 26.1% | 7.5% |
| LLaMA4-Scout | 45.5% | 47.2% | 3.6% |

Findings and insights. These results reveal a sharp performance drop as the number of intermediate steps increases. From a working memory perspective, this suggests that current LLMs struggle to retain or reconstruct causal paths to mediators when they are separated by multiple reasoning hops. This degradation highlights a key constraint in long-horizon causal reasoning.

These findings align with our earlier observation that mediator reasoning is a consistent bottleneck in decompositional analysis. By framing this in terms of working memory capacity, we offer a mechanistic explanation for why LLMs falter on such tasks and why enhanced memory mechanisms (e.g., intermediate supervision or tool-assisted retrieval) may be necessary for progress.

C LARGE LANGUAGE MODEL (LLM) USAGE DISCLOSURE

Large language models were used only for minor grammar revision and sentence-level polishing during manuscript preparation. They were not employed in ideation, methodological design, experimental execution, or result analysis. The scientific contributions, benchmarks, and evaluations presented in this work were entirely conceived and developed by the authors. LLM involvement was minimal in the research process.