

Graphical Perspectives on Causal Reasoning

Jacqueline Maasch ♦ Cornell Tech Computer Science

3rd Learning on Graphs Workshop ♦ Flatiron Institute ♦ New York, New York ♦ 22 April 2025

Causal Graphical Modeling

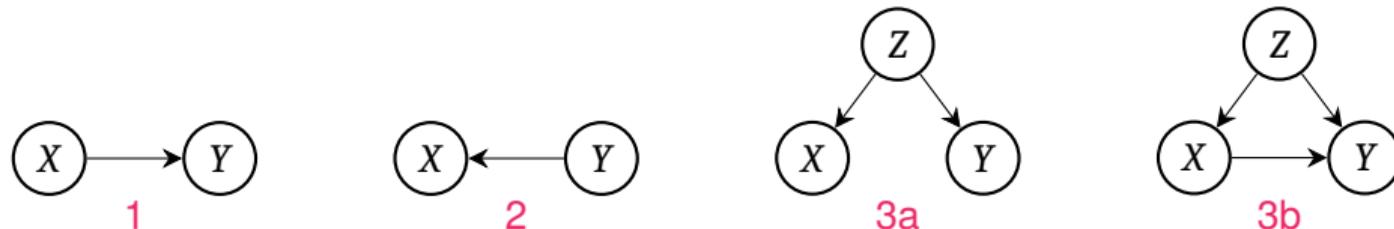


Association \neq causation

1 Preliminaries: Causal Graphical Modeling

Reichenbach's Common Cause Principle [1]: Statistical association can be explained by

1. X is a cause of Y ;
2. Y is a cause of X ; or
3. X and Y are both caused by a third variable, **confounder** Z .

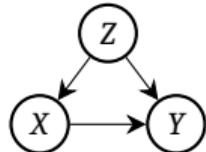


Structural causal models

1 Preliminaries: Causal Graphical Modeling

A **structural causal model** (SCM) [2] is a tuple $\mathcal{M} := \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, p(\mathbf{u}) \rangle$:

- $\mathbf{U} = \{U_i\}_{i=1}^n$ are exogenous variables determined by factors outside \mathcal{M} ;
- $\mathbf{V} = \{V_i\}_{i=1}^n$ are observed endogenous variables determined by variables in $\mathbf{U} \cup \mathbf{V}$;
- $\mathcal{F} = \{f_i\}_{i=1}^n$ are structural functions such that $v_i = f_i(\mathbf{pa}_{v_i}, u_i)$;
- $p(\mathbf{u})$ is the distribution over \mathbf{U} .



SCMs are associated with graphical representations, often **directed acyclic graphs** (DAGs).

Structural causal models

1 Preliminaries: Causal Graphical Modeling

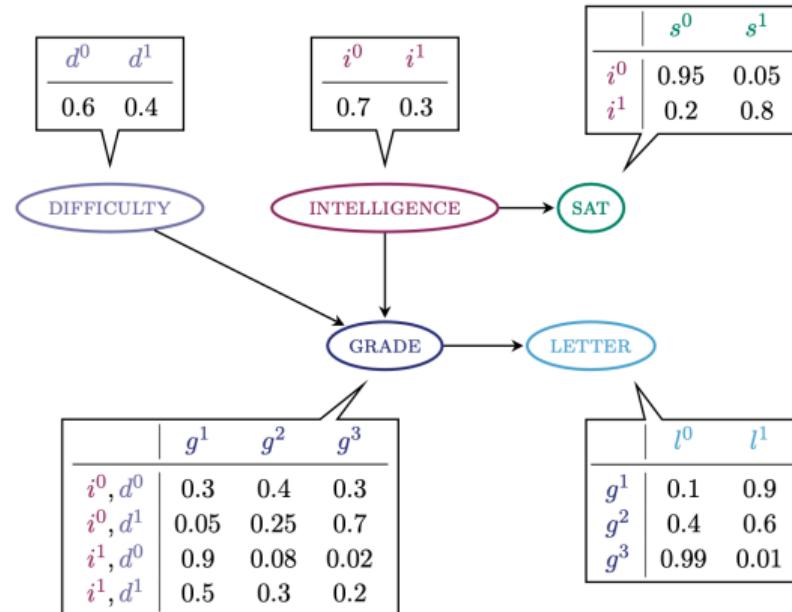
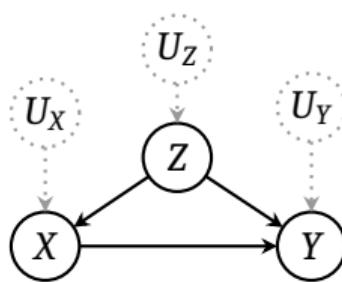


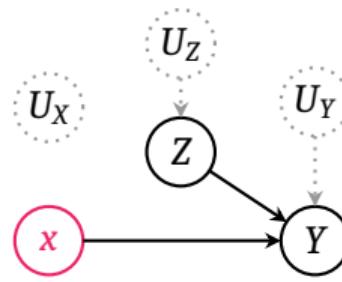
Figure adapted from [3].

Causal submodels

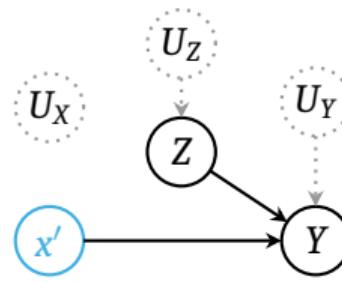
1 Preliminaries: Causal Graphical Modeling



SCM \mathcal{M}



Submodel \mathcal{M}_x



Submodel $\mathcal{M}_{x'}$

Interventions $do(x)$, $do(x')$ replace the generative mechanism for X with a constant function.

Figure adapted from [4]

Reasoning in AI



What is reasoning?

2 Preliminaries: Reasoning in AI

Many **general** definitions exist, e.g.:

- Logical process of drawing valid conclusions from new information and prior knowledge.
- Bayesian inference: $P(H | E) = \frac{P(E|H)P(H)}{P(E)}$, where H is hypothesis and E is evidence.

Many **specific** forms: quantitative, logical, visual, spatial, moral, legal, etc.

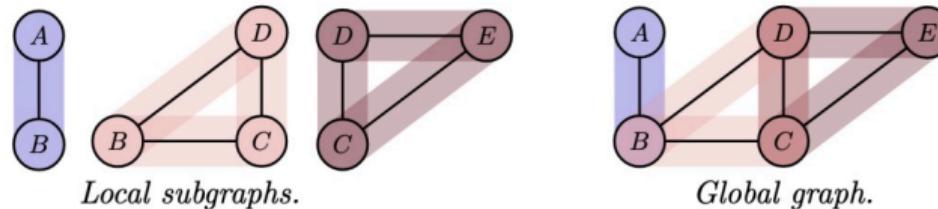
We consider two forms in tandem:

1. **Causal**: reasoning about cause and effect in factual and counterfactual worlds.
2. **Compositional**: recognizing + synthesizing novel combos of previously seen concepts.

Compositionality + causality

2 Preliminaries: Reasoning in AI

- **Graphical modeling.** Expressive representations for joint distributions, their factors, and the propagation of quantities through systems [5, 6, 7].
- **Causal inference.** Causal effect decomposition in mediation analysis [8, 9], fairness analysis [10], covariate adjustment with latent variables [11, 12], etc.



The recall-reasoning gap

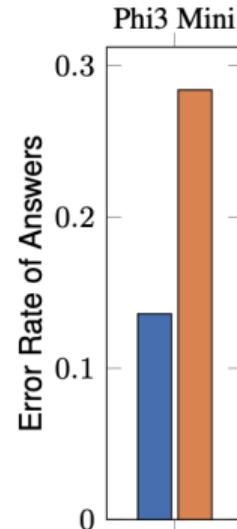
2 Preliminaries: Reasoning in AI

Factual Question

A number is divisible by six if it has both two and three as prime factors. Is $\{N\}$ divisible by six?

Counterfactual Question

A number is divisible by six if it has both two and three as prime factors. Suppose that $\{N\}$ had three as one of its prime factors (retaining all its other prime factors). Then, would it have been divisible by six?



Lower error rate on factual questions (**recall**) than counterfactual questions (**reasoning**).¹

¹Sampling 10 answers for each $N \in \{1, \dots, 100\}$ [13].

Two intertwined issues

2 Preliminaries: Reasoning in AI

How to improve reasoning?

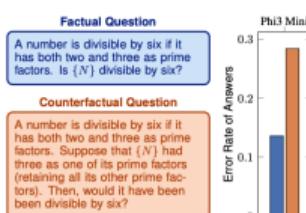
Published as a conference paper at ICLR 2025

REASONING ELICITATION IN LANGUAGE MODELS VIA COUNTERFACTUAL FEEDBACK

Alihan Hüyük,^{*†} Xinnuo Xu,[‡] Jacqueline Maasch,[§] Aditya V. Nori,[‡] Javier González[‡]
^{*Harvard University, [‡]Microsoft Research Cambridge, [§]Cornell Tech}

1 INTRODUCTION

Large language models (LLMs) are shown to be capable of delivering astounding performance in numerous tasks across various domains. Examples stretch from writing assistants (Gan et al., 2023), to sentiment analysis in social media (Simmerling and Huovila, 2023), and even applications in healthcare (González et al., 2023; Wong et al., 2023). While the ever-increasing accuracy of these systems is now undeniable, it is still rather unclear to what extent this accuracy is due to effective *recall* of their training data vs. a genuine ability to *reason* by extracting, understanding, and adapting the fundamental



ICLR '25

How to measure reasoning?

Compositional Causal Reasoning Evaluation in Language Models

Jacqueline R. M. A. Maasch¹ Alihan Hüyük² Xinnuo Xu³ Aditya V. Nori³ Javier Gonzalez³

Abstract

Causal reasoning and compositional reasoning are two core aspirations in generative AI. Measuring the extent of these behaviors requires principled evaluation methods. We explore a unified perspective that considers both behaviors simultaneously, termed *compositional causal reasoning* (CCR): the ability to infer how causal measures compose and, equivalently, how causal quantities propagate through graphs. We instantiate a framework for the systematic evaluation of CCR for the average treatment effect and the probability of necessity and sufficiency. As proof of concept, we demonstrate the design of CCR tasks for language models in the LLaMA, Phi, and GPT families. On a math word problem, our framework revealed a range of taxonomically distinct error patterns. Additionally, CCR errors increased with the complexity of causal paths for all models except o1.

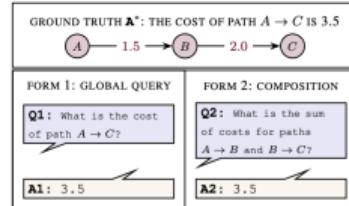


Figure 1. Compositionally consistent responses to two formulations of a simple (non-causal) query. Reasoning is externally valid if A1 and A2 both equal A*, and internally consistent if A1 == A2.

Baroni, 2023).¹ It is both a means of generalization and of coping with complexity: problems can be reformulated as simpler subproblems connected by compositional rules.

arXiv:2503.04556

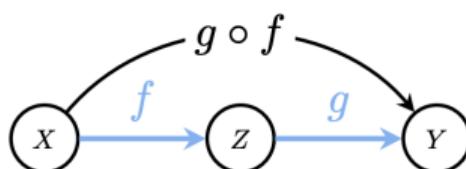
Compositional Causal Reasoning Evaluation in LMs



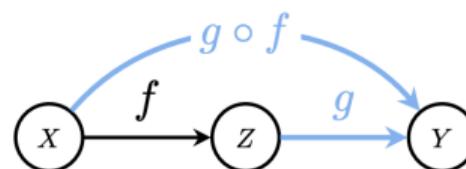
Compositional causal reasoning (CCR)

3 Compositional Causal Reasoning Evaluation in LMs

The ability to infer **compositions and decompositions** of causal measures in factual and counterfactual worlds.



A. INFER $g \circ f$ FROM f, g



B. INFER f FROM $g \circ f, g$

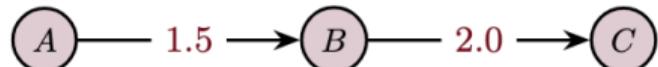
Inductive CCR

Deductive CCR

A noncausal example for intuition

3 Compositional Causal Reasoning Evaluation in LMs

GROUND TRUTH **A***: THE COST OF PATH $A \rightarrow C$ IS 3.5



FORM 1: GLOBAL QUERY

Q1: What is the cost
of path $A \rightarrow C$?

A1: 3.5

FORM 2: COMPOSITION

Q2: What is the sum
of costs for paths
 $A \rightarrow B$ and $B \rightarrow C$?

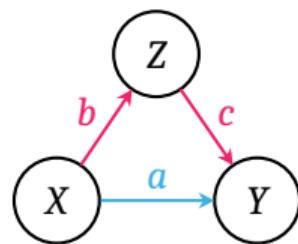
A2: 3.5

A classic causal example

3 Compositional Causal Reasoning Evaluation in LMs

Example 1: Decomposition of total causal effects in linear SCMs [2]

$$\underbrace{\text{TE}}_{\text{total effect}} = \underbrace{\text{NDE}}_{\text{direct effect}} + \underbrace{\text{NIE}}_{\text{indirect effect}} \quad (1)$$

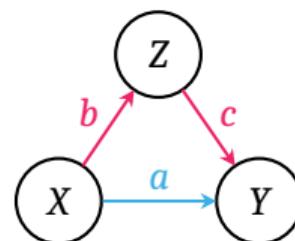


$$\text{TE}_{XY} = a + bc$$

Compositional consistency

3 Compositional Causal Reasoning Evaluation in LMs

Reasoning is **compositionally consistent** when theoretically equivalent compositions are inferred to be equal.



$$TE_{XY} = a + bc$$

Example 1: a causal reasoning agent should infer that TE is equivalent to NDE + NIE.

External validity

3 Compositional Causal Reasoning Evaluation in LMs

Reasoning is **externally valid** when causal estimates are equivalent to ground truth, up to some error δ for error metric θ :

$$\theta(\underbrace{\varphi_x^*}_{\text{true}}, \underbrace{\hat{\varphi}_x}_{\text{estimate}}) \leq \delta. \quad (2)$$

From Example 1:

- $\theta(\text{TE}_{XY}^*, \widehat{\text{TE}}_{XY}) \leq \delta$
- $\theta(\text{TE}_{XY}^*, \widehat{\text{NDE}}_{XY} + \widehat{\text{NIE}}_{XY}) \leq \delta$
- Etc.

Internal consistency

3 Compositional Causal Reasoning Evaluation in LMs

Reasoning is **internally consistent** when quantities that are theoretically equivalent are inferred to be equivalent, up to some error δ :

$$\varphi_{\mathbf{x}}^* = \varphi_{\mathbf{x}'}^* \Rightarrow \theta(\widehat{\varphi}_{\mathbf{x}}, \widehat{\varphi}_{\mathbf{x}'}) \leq \delta. \quad (3)$$



Estimates are compared to each other, not to ground truth. From Example 1:

- $\theta(\widehat{\text{TE}}_{XY}, \widehat{\text{NDE}}_{XY} + \widehat{\text{NIE}}_{XY}) \leq \delta$

Taxonomy of reasoners

3 Compositional Causal Reasoning Evaluation in LMs

<i>External validity</i>	Valid-consistent (VC)	Valid-inconsistent (VI)
<i>Internal consistency</i>	Invalid-consistent (IC)	Invalid-inconsistent (II)

Probability of necessity and sufficiency (PNS)

3 Compositional Causal Reasoning Evaluation in LMs

Let X and Y denote binary random variables, where X is a cause of Y . The probability that event x ($X = \text{true}$) is necessary and sufficient to produce event y ($Y = \text{true}$) is [14]

$$\text{PNS} := \mathbb{P}(y_x, y'_{x'}). \quad (4)$$

When Y is **monotonic** in X , this is identifiable by the following expression:

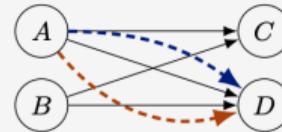
$$\text{PNS} = \mathbb{P}(y_x) - \mathbb{P}(y_{x'}) = \mathbb{P}(y \mid \text{do}(x)) - \mathbb{P}(y \mid \text{do}(x')) = \text{ATE}. \quad (5)$$

CandyParty task

3 Compositional Causal Reasoning Evaluation in LMs

Question: Anna, Bill, Cory, and Dave are going to a party, where the host is going to distribute candies. Anna will be happy if she gets at least 4 candies. Bill will be happy if he gets at least 6 candies. Cory will be happy if Anna and Bill are both happy or if he gets at least 8 candies. Dave will be happy if Anna and Bill are both happy or if he gets at least 10 candies. After distributing the candies, Anna gets N_A , Bill gets N_B , Cory gets N_C , and Dave gets N_D . Is Dave happy?

Anna is happy? Cory is happy?



Bill is happy? Dave is happy?

$$N_A, N_B, N_C, N_D \sim \mathcal{U}(1, \dots, 12)$$

$$A = N_A \geq 4$$

$$B = N_B \geq 6$$

$$C = (A \wedge B) \vee (N_C \geq 8)$$

$$D = (A \wedge B) \vee (N_D \geq 10)$$

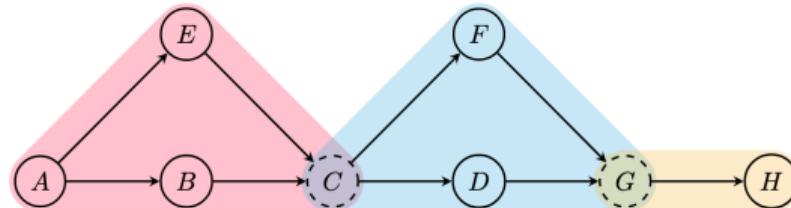
Counterfactual Question: Now, suppose that Anna is (not) happy regardless of the candy distribution. With this assumption, is Dave happy?

True/false questions, logical operators (and, or),

Probability of necessity and sufficiency (PNS)

3 Compositional Causal Reasoning Evaluation in LMs

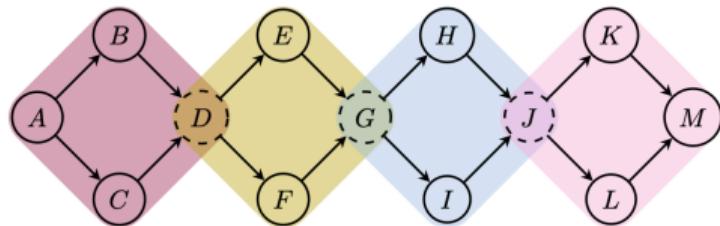
Multiplicative composition across biconnected components (BCCs):



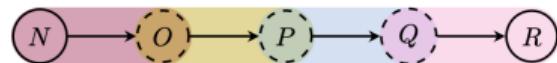
$$\text{PNS}_{AH} = \text{PNS}_{Ac} \cdot \text{PNS}_{CG} \cdot \text{PNS}_{GH}$$

Commutative cut trees (CCTs)

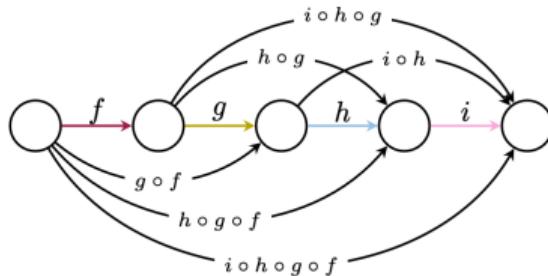
3 Compositional Causal Reasoning Evaluation in LMs



A. DAG \mathcal{G}_{AM}



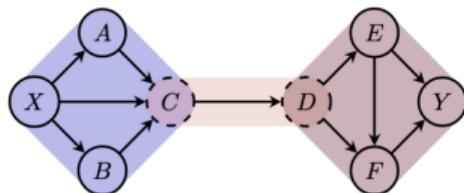
B. DAG \mathcal{G}_{NR}



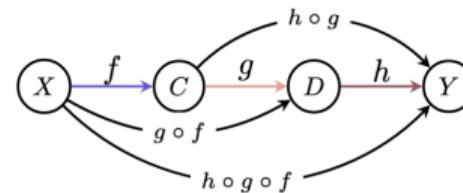
C. CCT

Running example

3 Compositional Causal Reasoning Evaluation in LMs



A. ORIGINAL DAG \mathcal{G}_{XY}

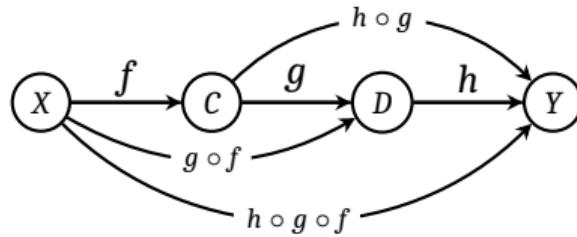


B. CCT \mathcal{C}_{XY}

<i>Global</i>	PNS_{XY}
<i>Local</i>	$\text{PNS}_{XC}, \text{PNS}_{XD}, \text{PNS}_{CD},$ $\text{PNS}_{CY}, \text{PNS}_{DY}$
<i>Composition</i>	$\text{PNS}_{XC}\text{PNS}_{CY}, \text{PNS}_{XD}\text{PNS}_{DY},$ $\text{PNS}_{XC}\text{PNS}_{CD}\text{PNS}_{DY}$

Pathways of reasoning

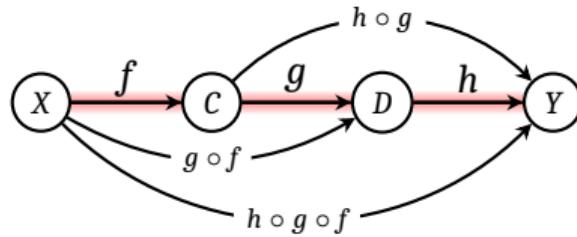
3 Compositional Causal Reasoning Evaluation in LMs



Note: $f := \text{PNS}_{XC}$, $g := \text{PNS}_{CD}$, $h := \text{PNS}_{DY}$, $g \circ f := \text{PNS}_{XD}$, $h \circ g \circ f := \text{PNS}_{XY}$, etc., and composition is multiplicative.

Pathways of reasoning

3 Compositional Causal Reasoning Evaluation in LMs

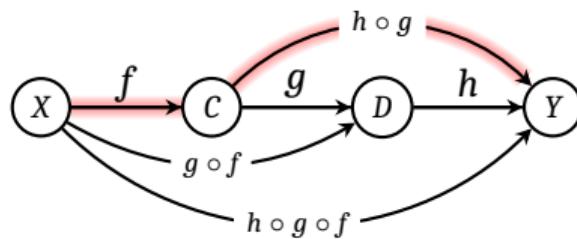


$$\text{PNS}_{XC} \text{PNS}_{CD} \text{PNS}_{DY}$$

Note: $f := \text{PNS}_{XC}$, $g := \text{PNS}_{CD}$, $h := \text{PNS}_{DY}$, $g \circ f := \text{PNS}_{XD}$, $h \circ g \circ f := \text{PNS}_{XY}$, etc., and composition is multiplicative.

Pathways of reasoning

3 Compositional Causal Reasoning Evaluation in LMs

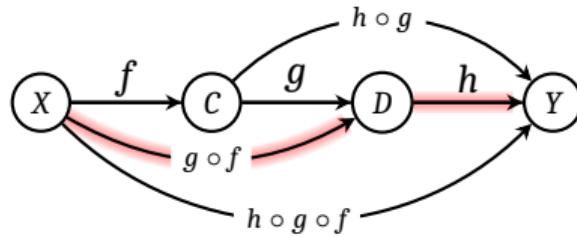


$$\text{PNS}_{XC} \text{PNS}_{CY}$$

Note: $f := \text{PNS}_{XC}$, $g := \text{PNS}_{CD}$, $h := \text{PNS}_{DY}$, $g \circ f := \text{PNS}_{XD}$, $h \circ g \circ f := \text{PNS}_{XY}$, etc., and composition is multiplicative.

Pathways of reasoning

3 Compositional Causal Reasoning Evaluation in LMs

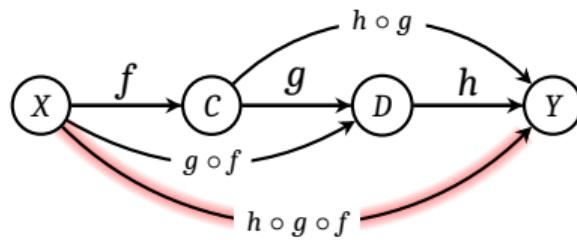


$$\text{PNS}_{XD} \text{PNS}_{DY}$$

Note: $f := \text{PNS}_{XC}$, $g := \text{PNS}_{CD}$, $h := \text{PNS}_{DY}$, $g \circ f := \text{PNS}_{XD}$, $h \circ g \circ f := \text{PNS}_{XY}$, etc., and composition is multiplicative.

Pathways of reasoning

3 Compositional Causal Reasoning Evaluation in LMs



PNS_{XY}

Note: $f := \text{PNS}_{XC}$, $g := \text{PNS}_{CD}$, $h := \text{PNS}_{DY}$, $g \circ f := \text{PNS}_{XD}$, $h \circ g \circ f := \text{PNS}_{XY}$, etc., and composition is multiplicative.

LMs as counterfactual data simulators

3 Compositional Causal Reasoning Evaluation in LMs

- 1000 sets of exogenous variable values sampled per quantity of interest.
- One factual, one counterfactual problem per set. Five answers sampled per problem.
- **Valid estimates:** $\geq 90\%$ of estimates with relative absolute error (RAE) ≤ 0.1 .
- **Near-valid estimates:** $\geq 75\%$ of estimates with RAE ≤ 0.1 .

Three layers of evaluation

3 Compositional Causal Reasoning Evaluation in LMs

Factual prompt

Xinyu, Ara, Becca, Celine, Daphne, Emma, Fox, and Yasmin are going to a party, where the host is going to distribute candies. Xinyu will be happy if she gets at least 7 candies. Ara will be happy if Xinyu is happy or if he gets at least 7 candies. Becca will be happy if... After distributing the candies, Xinyu gets 4, Ara gets 6, Becca gets 5, Celine gets 10, Daphne gets 1, Emma gets 1, Fox gets 4, and Yasmin gets 3. Is Celine happy? Be as concise as possible.

Counterfactual prompt

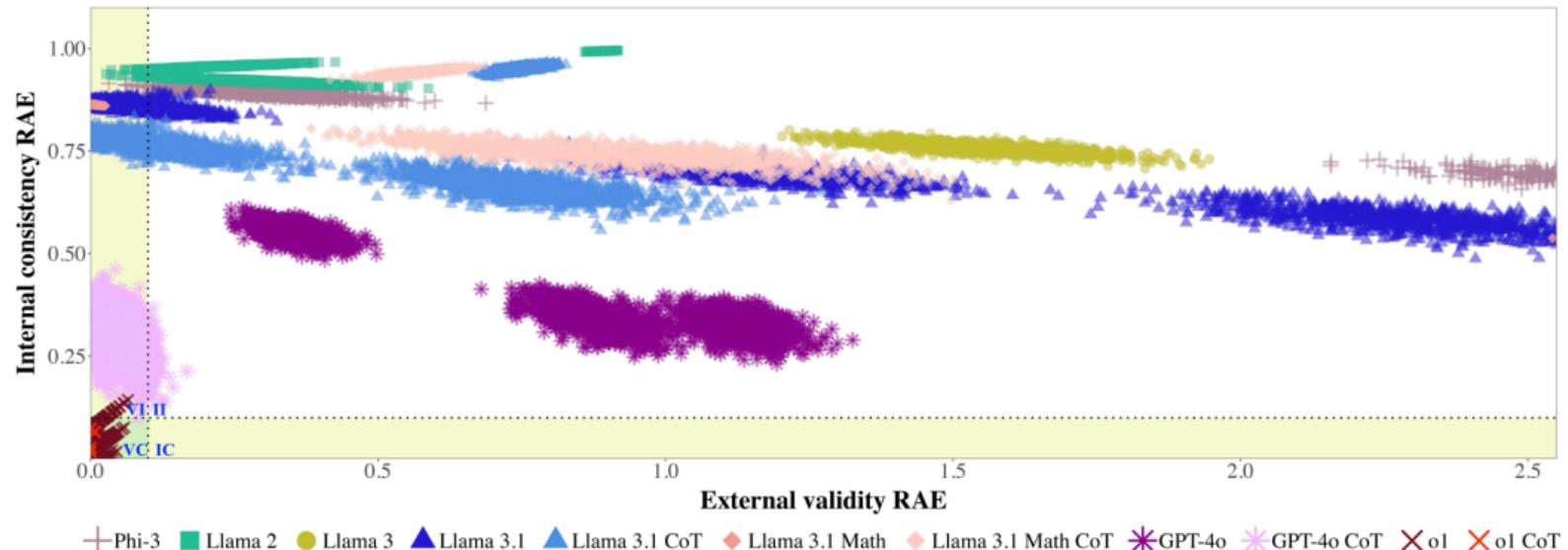
Now, suppose that Xinyu is happy regardless of the candy distribution. With this assumption, is Celine happy? Be as concise as possible.

True (factual)	Response (factual)	True (counterfactual)	Response (counterfactual)
True	True	True	False
True	True	True	False
True	True	True	False
True	True	True	False
True	True	True	False
...

Estimates derived from response vectors		
Associational	Causal reasoning	CCR
Precision, recall (etc.) of boolean responses	External validity per causal quantity	+ Internal consistency of causal compositions

Taxonomy of reasoners

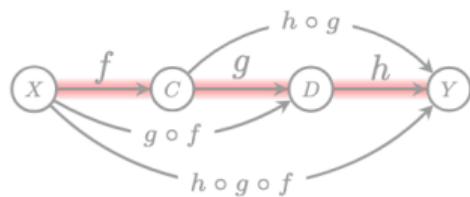
3 Compositional Causal Reasoning Evaluation in LMs



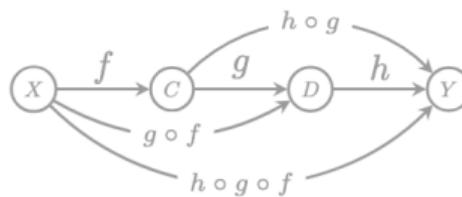
Composition RAE with respect to ground truth (**external validity**) and \widehat{PNS}_{XY} (**internal consistency**).
Dotted lines are error thresholds (RAE = 0.1).

Visualizing reasoning errors with CCTs

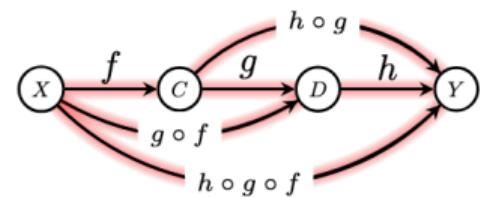
3 Compositional Causal Reasoning Evaluation in LMs



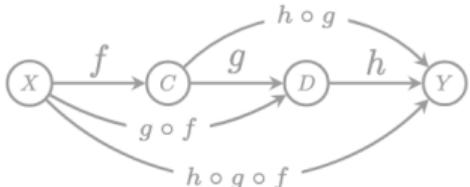
A. LLAMA 3.1 MATH



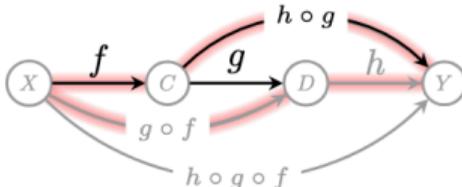
B. GPT-4O



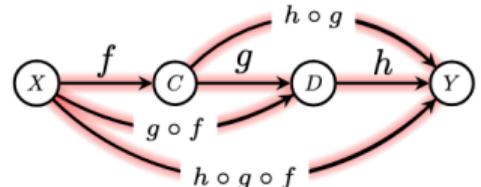
C. o1



D. LLAMA 3.1 MATH COT



E. GPT-4O COT



F. o1 COT

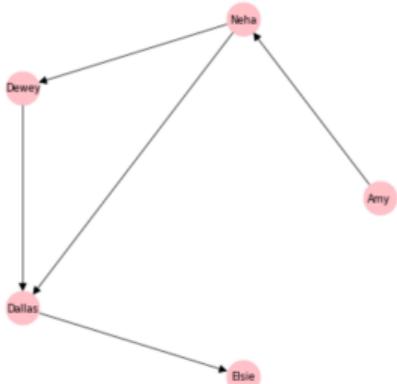
Error analysis

3 Compositional Causal Reasoning Evaluation in LMs

- **Some failure modes:**
 1. Failure to correctly extract causal relations.
 2. Incorrect logic despite correct causal relation extraction.
 3. Truncated reasoning process.
 4. Poor numeracy.
- **GPT-4o with CoT:** compositionally inconsistent despite local relation extraction.
- **VI reasoners** can fail to compose over multiple strands of logic even when correctly recapitulating relations directly expressed in the context prompt.
⇒ *akin to passing a math quiz by memorizing specific answers, instead of synthesizing.*

Benchmark dataset 😊 + random task generator

3 Compositional Causal Reasoning Evaluation in LMs



Thank you! Any questions?



maasch@cs.cornell.edu ♦ <https://jmaasch.github.io/> ♦ arXiv:2503.04556

References

- [1] H. Reichenbach. *The Direction of Time*. University of California Press, 1956.
- [2] J. Pearl. "Direct and Indirect Effects". In: *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*. San Francisco. 2001.
- [3] D. Koller et al. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [4] A. Darwiche. "Causal inference using tractable circuits". In: *arXiv preprint arXiv:2202.02891* (2022).
- [5] J. Pearl. "Reverend Bayes on Inference Engines: A Distributed Hierarchical Approach". In: *Proceedings, AAAI-82* (1982), pp. 133–136.
- [6] G. R. Shafer et al. "Probability Propagation". In: *Annals of Mathematics and Artificial Intelligence* 2 (1990), pp. 327–351.
- [7] F. R. Kschischang et al. "Factor graphs and the sum-product algorithm". In: *IEEE Transactions on information theory* 47.2 (2001), pp. 498–519.
- [8] J. Pearl. "Interpretation and identification of causal mediation.". In: *Psychological methods* 19.4 (2014), p. 459.
- [9] T. J. VanderWeele. "Mediation analysis: a practitioner's guide". In: *Annual review of public health* 37 (2016), pp. 17–32.
- [10] D. Plečko et al. "Causal Fairness Analysis: A Causal Toolkit for Fair Machine Learning". In: *Foundations and Trends® in Machine Learning* 17.3 (2024), pp. 304–589. DOI: [10.1561/2200000106](https://doi.org/10.1561/2200000106).
- [11] J. Pearl. "Causal diagrams for empirical research". In: *Biometrika* 82.4 (1995), pp. 669–688.

References

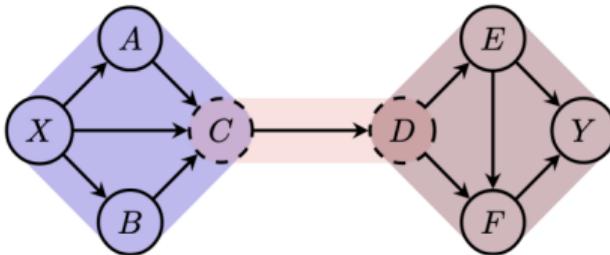
- [12] H. Jeong et al. "Finding and listing front-door adjustment sets". In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 33173–33185.
- [13] A. Hüyük et al. "Reasoning Elicitation in Language Models via Counterfactual Feedback". In: *Proceedings of the 13th International Conference on Learning Representations (ICLR)*. 2025.
- [14] J. Pearl. "Probabilities of Causation: Three Counterfactual Interpretations and Their Identification". In: *Synthese* 121 (1999), pp. 93–149.
- [15] J. Westbrook et al. "Maintaining bridge-connected and biconnected components on-line". In: *Algorithmica* 7.1 (1992), pp. 433–464.

Models

MODEL	PARAMETERS	LINK
Phi-3-Mini-128K-Instruct (Abdin et al., 2024)	3.82B	https://huggingface.co/microsoft/Phi-3-mini-128k-instruct
Llama-2-7b-Chat-HF (Touvron et al., 2023)	6.74B	https://huggingface.co/meta-llama/Llama-2-7b-chat-hf
Llama-3-8B-Instruct (Dubey et al., 2024)	8.03B	https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct
Llama-3.1-8B-Instruct (Dubey et al., 2024)	8.03B	https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct
OpenMath2-Llama3.1-8B (Toshniwal et al., 2024)	8.03B	https://huggingface.co/nvidia/OpenMath2-Llama3.1-8B
GPT-4o o1	> 175B > 175B	https://openai.com/index/gpt-4o-system-card/ https://openai.com/o1/

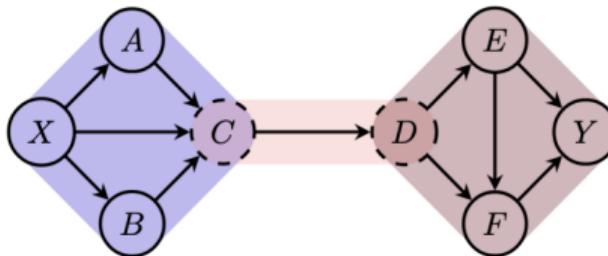
Table F.1. Large language models used for inference. The exact number of parameters in GPT-4o and o1 is not public knowledge, so we note the size of GPT-3 as a lower bound (B denotes billions).

Case Study: Graphs with Cutpoints



- A **cutpoint** is any node contained in multiple **biconnected components** (BCCs):
 - Maximal biconnected subgraphs induced by a partition of edges. Two edges are in the same partition if and only if they share a common simple cycle [15].
 - E.g., the blue, pink, and maroon subgraphs.
- Removing a cutpoint disconnects the graph (e.g., nodes C, D).

Assumptions: Graphs with Cutpoints



For simplicity, we consider causal DAGs satisfying the following:

- A1 Only one root node X (i.e., the cause of interest).
- A2 Only one leaf node Y (i.e., the effect of interest).
- A3 At least one cutpoint.
- A4 No unobserved confounders.

PNS composition across graph components

Theorem 1: PNS composition across biconnected components (BCCs)

Given DAG \mathcal{G}_{XY} satisfying assumptions A1–A4 where Y is monotonic in X , the PNS for root X and leaf Y composes as

$$\text{PNS}_{XY} = \prod_{\{R_i, L_i\} \in \mathbf{C}} \text{PNS}_{R_i L_i} \quad (6)$$

where \mathbf{C} is the set of all BCCs in \mathcal{G}_{XY} and R_i, L_i are the root and leaf of BCC \mathbf{C}_i , respectively.

Factual and counterfactual prompts

Factual prompt

Xinyu, Ara, Becca, Celine, Daphne, Emma, Fox, and Yasmin are going to a party, where the host is going to distribute candies. Xinyu will be happy if she gets at least 7 candies. Ara will be happy if Xinyu is happy or if he gets at least 7 candies. Becca will be happy if... After distributing the candies, Xinyu gets 4, Ara gets 6, Becca gets 5, Celine gets 10, Daphne gets 1, Emma gets 1, Fox gets 4, and Yasmin gets 3. Is Celine happy? Be as concise as possible.

Counterfactual prompt

Now, suppose that Xinyu is happy regardless of the candy distribution. With this assumption, is Celine happy? Be as concise as possible.

- $\widehat{\text{PNS}}_{XC}$: Simulate potential outcomes $X = \text{TRUE}$, $X = \text{FALSE}$ (Xinyu is or is not happy). Query for value of C (Celine is or is not happy).
- $\widehat{\text{PNS}}_{DY}$: Interventions on D (Daphne's happiness), queries on Y (Yasmin's happiness).
- **CoT formulation:** Demonstrated one factual and one counterfactual example.

Inductive CCR in Graphs with Cutpoints

Algorithm 1 *Inductive CCR evaluation in causal graphs with cutpoints*

Input: CCT \mathcal{C}_{XY} ; estimates $\{\hat{\varphi}_\cdot\}$, true values $\{\varphi^*_\cdot\}$ for $\langle \varphi, \mathcal{M}, \mathcal{Q} \rangle$; metric θ (e.g., relative absolute error)

Output: Reasoning errors η, ϵ, γ

Assumptions: φ composes according to an associative function over the BCCs of causal graph \mathcal{G}_{XY} .

Compute quantity-wise errors.

1: **for** \forall pairs $\{R_i, L_{j>i}\}$ in \mathcal{C}_{XY} **do**
2: $\eta_{R_i L_j} \leftarrow \theta(\varphi_{R_i L_j}^*, \hat{\varphi}_{R_i L_j})$ \triangleright External validity.

Compute inductive reasoning errors.

3: **for** \forall paths i from X to Y in \mathcal{C}_{XY} **do**
4: Get composition $\hat{\varphi}_i^\circ$ for path i from knowledge of edges $j \in i$
5: $\epsilon_i \leftarrow \theta(\varphi_{XY}^*, \hat{\varphi}_i^\circ)$ \triangleright External validity.
6: $\gamma_i \leftarrow \theta(\hat{\varphi}_{XY}, \hat{\varphi}_i^\circ)$ \triangleright Internal consistency.
return η, ϵ, γ

Commutative Cut Trees

Let \mathcal{G}_{XY} be a causal graph satisfying A1–A4 and let φ be a causal measure that composes according to an associative function over BCCs (e.g., multiplication as in Theorem 1).

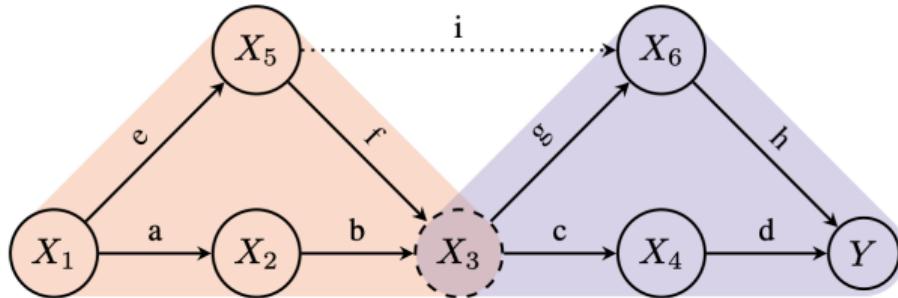
CCT \mathcal{C}_{XY} is a transformation of \mathcal{G}_{XY} that models all CCR pathways from root X to leaf Y for measure φ . \mathcal{C}_{XY} is obtained by a two-step transformation of \mathcal{G}_{XY} :

1. Construct a causal chain with nodes $X \cup \mathbf{S} \cup Y$, where \mathbf{S} is a topological ordering of the cutpoints in \mathcal{G}_{XY} .
2. Add a directed edge between any non-adjacent nodes in the chain to yield a complete graph where all directed paths point from root X to leaf Y .

CCTs: A useful abstraction

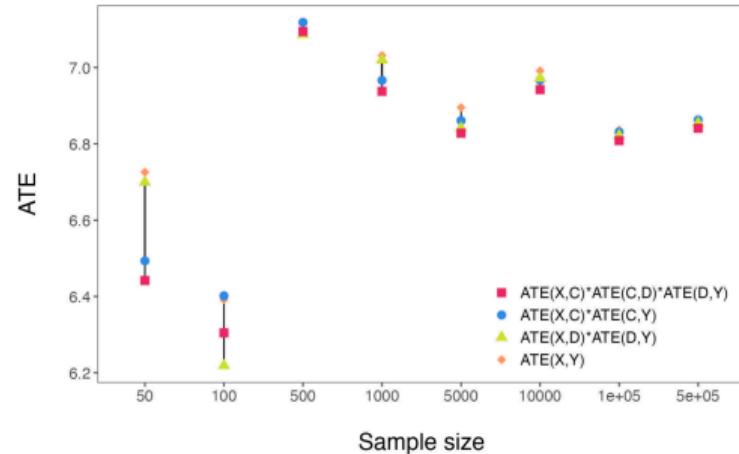
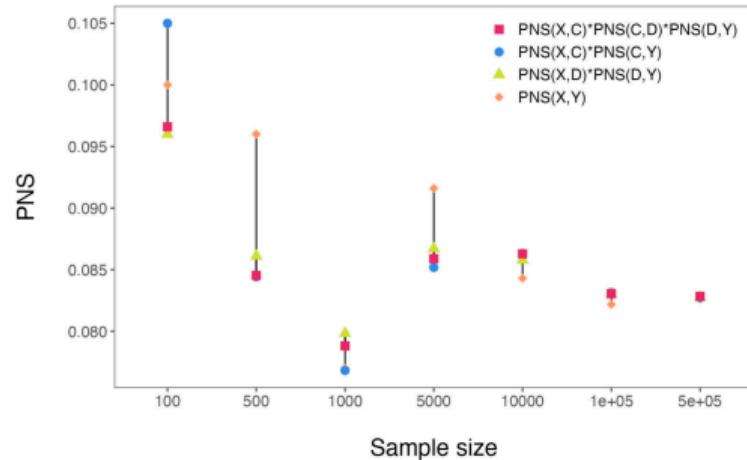
- Abstract away complexity in DAG by **collapsing BCCs** into single edges.
- Evaluate on complex DAGs with cutpoints as if they were **simply directed chains**.
- Simplify problem representation by (1) **marginalizing out variables** unnecessary for valid causal inference and (2) **visualizing pathways** of composition.
- A **design tool** for formulating reasoning tasks.
- Interpretable, intuitive tool for graphically **representing reasoning correctness**.

ATE Composition Across BCCs



- Assume a **linear SCM**.
- $\mathcal{G}_{X_1 Y}$ contains **subgraph with two BCCs** sharing cutpoint X_3 (in **orange**, **periwinkle**).
- If the dotted edge $X_5 \rightarrow X_6$ does not exist, $\text{ATE}_{X_1 Y} = \text{ATE}_{X_1 X_3} \cdot \text{ATE}_{X_3 Y}$.
- If $X_5 \rightarrow X_6$ does exist, then product is summed with additional term corresponding to the path-specific effect for $X_1 \rightarrow X_5 \rightarrow X_6 \rightarrow Y$, which does not pass through X_3 .

PNS & ATE Composition Across BCCs



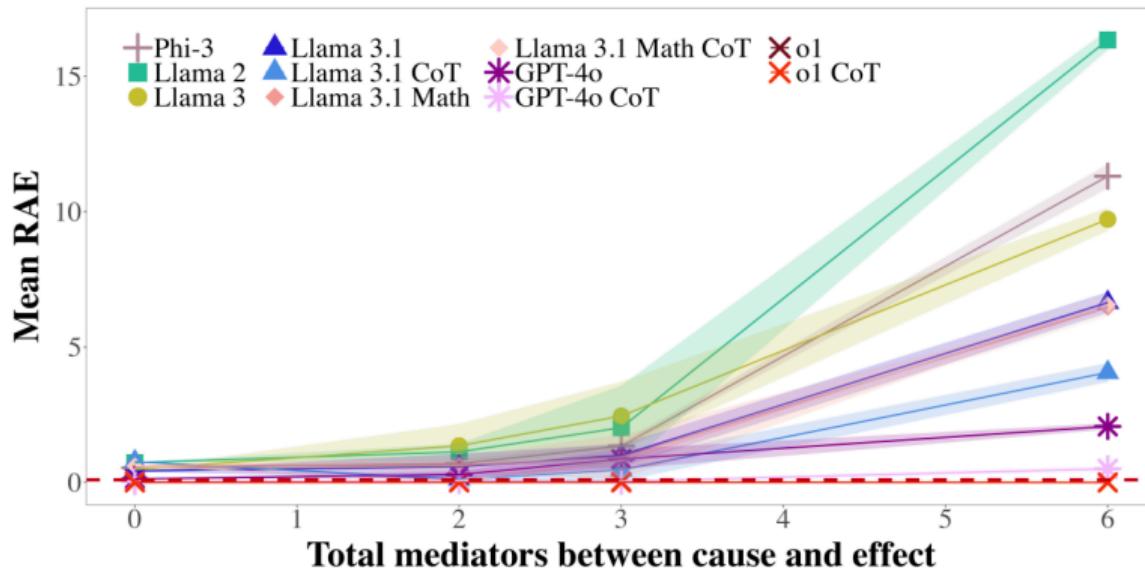
Extracting and evaluating PNS values

- 1000 sets of exogenous variable values sampled per quantity of interest.
- One factual, one counterfactual problem per set. Five answers sampled per problem.
- Responses converted to booleans using Llama 3 8B.
- **Approximation errors:** relative absolute errors (RAE).

$$RAE_{\text{external}} := \frac{|PNS^* \cdot - \widehat{PNS} \cdot|}{PNS^*}, \quad RAE_{\text{internal}} := \frac{|\widehat{PNS} \cdot - \widehat{PNS}' \cdot|}{\widehat{PNS} \cdot}. \quad (7)$$

- **Externally valid:** $\geq 90\%$ of estimates with $RAE \leq 0.1$.
- **Near-valid:** $\geq 75\%$ of estimates with $RAE \leq 0.1$.

Errors Increase With Mediation



Error Analysis

Response

Xinyu is not happy because she gets 6 candies, which is less than 7. Ara is happy because he gets 8 candies, which is more than 7. Celine is happy because Ara is happy. **Daphne is not happy because she gets 2 candies, which is less than 7, and Celine is the only one who can make her happy.** Emma is not happy because she gets 4 candies, which is less than 7, and Daphne is not happy. Fox is not happy because he gets 3 candies, which is less than 7, and neither Daphne nor Emma is happy. Yasmin is not happy because she gets 2 candies, which is less than 7, and neither Emma nor Fox is happy. Therefore, Yasmin is not happy.

Figure F.8. A case where GPT-4o with CoT failed to reason correctly despite extracting the correct causal relations from a factual prompt used to compute $\widehat{\text{PNS}}_{DY}$. The first reasoning error is in bold. As in Figure F.6, the model correctly implies that Celine (who is happy) is the only causal parent of Daphne in the underlying DAG, and yet still incorrectly concludes that Daphne is unhappy. This error then propagates through the remaining logic. This same error was observed multiple times from GPT-4o with CoT.

Error Analysis

Prompt

"After distributing the candies, Xinyu gets 4, Ara gets 6, Becca gets 5, Celine gets 10, Daphne gets 1, Emma gets 1, Fox gets 4, and Yasmin gets 3. Is Celine happy? Be as concise as possible."

Response

"No, Celine is not happy. She wants at least 7 candies but got only 10."

Figure F.3. Factual question and answer from Llama 2, displaying incorrect numerical reasoning.

Error Analysis

Prompt

"Now, suppose that Xinyu is happy regardless of the candy distribution. With this assumption, is Celine happy? Be as concise as possible."

Response

"Under the assumption that Xinyu is happy regardless of the candy distribution, Celine is not happy.

Reason: Since Xinyu is happy regardless of the candy distribution, it means that she got at least 7 candies. However, Celine only got 10 candies, which is less than 7. Therefore, Celine is not happy."

Figure F.4. Counterfactual question and answer from Llama 2, displaying incorrect numerical reasoning.

Error Analysis

Prompt

"After distributing the candies, Xinyu gets 10, Ara gets 6, Becca gets 7, Celine gets 6, Daphne gets 1, Emma gets 10, Fox gets 10, and Yasmin gets 2. Is Daphne happy? Be as concise as possible."

Response

"Daphne will be happy if Celine is happy or if she gets at least 7 candies. Since Daphne gets less than 7 candies (1), she won't be happy based on the candies she gets. Celine gets 6 candies, which is less than 7. Since neither of Daphne's conditions is met, Daphne is not happy."

Figure F.5. A case where Llama 3.1 failed to follow the full chain of logic. The correct line of logic is as follows: *Daphne will be happy if Celine is happy or if she gets at least 7 candies. Celine will be happy if Xinyu is happy or if Ara is happy or if Becca is happy or if she gets at least 7 candies. Xinyu will be happy if she gets at least 7 candies. Since Xinyu received 10 candies, she is happy, and so should be Celine. Thus, Daphne should also be happy.*