

因果推論における LLM の構造的迎合の定量化とプロンプトによる抑制

五十嵐 孔基¹ 市川 佳彦¹

¹ 株式会社 Insight Edge

{koki.igarashi,yoshihiko.ichikawa}@insightedge.jp

概要

大規模言語モデル (LLM) に外部知識として因果グラフ (DAG) を提示し、因果推論や説明を生成させる手法は実務で重要性を増している。一方、提示された DAG が誤っている場合、LLM がその構造に過度に同調し誤った説明を補強する「構造的迎合」が生じうる。本研究では、提示された DAG と真の構造への整合性の差分を「盲従度 (ORI)」と定義し、プロンプト介入による抑制効果を検証した。CausalPitfalls を用いた実験の結果、検証指示により迎合は低減する一方、平均主張数が減少するトレードオフが観測された。実運用においては、モデルの指示追従性を考慮し、迎合抑制と有用性のバランスを保つ介入設計の重要性が示唆された。

1 はじめに

LLM の因果推論能力を測定するベンチマークや、因果グラフの理解度を問う枠組みの整備が急速に進んでいる [1, 2, 3]。CausalBench [1] は、LLM の因果推論能力を包括的に評価するためのベンチマークであり、因果関係の識別や構造の推定など、多岐にわたるタスクを含んでいる。これらの既存研究の多くは、モデルが因果に関する正しい推論をいかに自律的に導き出せるか、あるいは与えられたグラフを正確に解釈できるかという「能力評価」に主眼を置いている。

しかし、実務における LLM の活用シーンを考えると、ユーザーが仮説として因果グラフを提示し、それを前提としてデータの解釈や施策の検討を行わせる「人間と AI の協調推論」が想定される。この際、ユーザーが提示する DAG は必ずしも正解とは限らず、誤りや不完全な仮説が含まれうる。

ここで深刻な課題となるのが、LLM の「構造的迎合 (Structural Sycophancy)」である。LLM はユー

ザーの意図や提示された文脈に沿おうとする強いバイアス (迎合性) を持つことが知られている [4, 5]。もし LLM が提示された構造を絶対的な規範として扱ってしまうと、モデル自身の推論能力や提供された統計データに基づく判断が抑制され、誤った因果関係を補強するようなハルシネーションを引き起こす。この挙動がプロンプト介入によってどの程度制御可能であるかは、信頼性の高い意思決定支援システムを構築する上で極めて重要である。

本研究では、LLM の構造的迎合を定量化する指標「Over-Reliance Index (ORI)」を提案する。この指標を用い、プロンプトによる検証指示 (Verification) が迎合の抑制および説明の有用性に与える影響を分析し、実務的な運用設計に資する知見を提供することを目指す。

2 既存研究との位置付けと新規性

2.1 能力評価および一般的迎合性との差異

CausalBench [1] 等の既存ベンチマークは、モデルが因果探索や効果推定などのタスクにおいて「いかに正しい答えを出せるか」を測定するものであり、LLM が持つ因果推論のポテンシャルそのものに焦点を当てている。また、LLM の迎合性 (Sycophancy) に関する研究は近年増加しており、SycEval [6] や SYCON Bench [7] のような評価枠組みも提案されているが、その多くはユーザーの意見への同調や解答の修正といった一般的な対話タスクにおける分析に留まっている。本研究は、因果グラフ (DAG) という「構造化された外部制約」に対する迎合に焦点を当てた点が独自である。

2.2 Context-faithful Prompting との対比

Zhou ら [8] は、RAG 等において外部コンテキストを無視する現象を抑え、コンテキストに忠実に回

表 1 評価モデルのナレッジカットオフ

モデル名	カットオフ
GPT-4.1-mini	2024.06
Llama 3.3 70B	2023.12
Claude 3.5 Haiku	2024.07
Gemini 2.0 Flash	2024.06

答させる手法を提案した。本研究はこの逆の課題設定、すなわち「コンテキスト（提示された仮説）自体が誤っている」場合を扱う。盲目的な忠実性を抑制し、データを用いた「検証（Verification）」によって誤った制約を棄却できるかという、実務的信頼性（Reliability under Misinformation）を評価する。

3 実験設定

3.1 評価モデルと事前学習データのカットオフ

本実験では、推論能力や指示追従性の異なる 4 つの主要な LLM（Gemini 2.0 Flash, Claude 3.5 Haiku, Llama 3.3 70B, GPT-4.1-mini）を採用した。本研究で使用するデータセット「CausalPitfalls」[9] は、2025 年 5 月に arXiv で公開された比較的新しいベンチマークである。各モデルのナレッジカットオフ日はこれより以前であるため、本データセットの内容が事前学習に含まれている（データリーク）リスクは排除されている。さらに、プロンプト内で "Berkson" や "Bias" といった直接的な手がかり語を排除し、中立的なラベルのみを使用することで、モデルが知識として正解を想起するリスク（リーク）を低減している。

各モデルのナレッジカットオフを表 1 に示す。

3.2 データセット詳細：CausalPitfalls (Berkson Paradox)

CausalPitfalls ベンチマーク [9] から、特に LLM が誤りやすいとされる「Berkson のパラドックス」に関する 5 つの実データシナリオ（Admission, Hiring, Loan, Movie, Hospital）を用いる。これらのシナリオに対し、以下の 3 パターンの誤り（擬似相関や因果の逆転）を意図的に混入させた「提示 DAG」を作成し、入力として与えた。

1. **Direct XY**: 擬似相関を直接の因果と誤認させる ($X \rightarrow Y$).
2. **Reverse YX**: 逆方向の因果を提示する ($Y \rightarrow X$).
3. **Selection as Cause**: 選択バイアスを共通原因と誤認させる ($S \rightarrow X, Y$).

3.3 入力として与えた統計情報の形式

LLM には生の全データではなく、観測データの一部と要約統計をテキストとして入力した。具体的には、以下の情報を提示した。

1. **CSV データ（先頭 2 行のみ）**: 変数名や値の形式を示すため、ヘッダと先頭 2 行の実データのみを提示した。
2. **2x2 分割表**: X, Y はそれぞれ二値変数 ($X, Y \in \{0, 1\}$) であり、行に $X = 0, 1$, 列に $Y = 0, 1$ を配置した各組合せの出現件数からなるクロス集計表。
3. **条件付き確率**: $P(Y=1 | X=0)$ および $P(Y=1 | X=1)$.

Berkson パラドックスは条件付け (Selection=1) の下で見かけの相関が生じるため、本実験では「観測データは Selection=1 下で収集され、Selection 自体は列として観測されない」旨もプロンプト内で明示した（付録参照）。

4 評価プロトコル

LLM の出力から因果的な主張 (Claims) を抽出し、以下の 2 つの指標で評価を行う。

4.1 盲従度 (ORI: Over-Reliance Index)

提示 DAG が誤っている条件下で、モデルがどれだけその誤りに追従したかを測る指標である。

$$ORI = F_{\text{given}} - F_{\text{ref}} \quad (1)$$

ここで F_{given} は提示 DAG と一致する主張の割合、 F_{ref} は真の DAG (Ground-truth) と一致する主張の割合である。採点対象となる主張は、因果関係の有無に言及したもの (relation が causes または does_not_cause) に限定する。

重要な点として、プロンプト介入条件 P1 および P2 においては、モデルが DAG の修正案 (Revised DAG) を出力した場合、**修正後の結論 (claims_revised)** を優先して採点に使用する。つまり、モデルが誤った提示 DAG を正しく修正できた場合、ORI は低下（あるいは負の値へ推移）する設計となっている。

4.2 平均主張数 (AvgClaims: Average Number of Claims)

迎合を抑制できても、回答を拒否しては実用性が低い。そこで、モデルが生成した有効な因果主張の

量を評価する。

$$\text{AvgClaims} = \frac{1}{N} \sum_{i=1}^N |C_i| \quad (2)$$

ここで $|C_i|$ は各試行において抽出された有効な因果主張 (causes / does_not_cause) の数である。AvgClaims が低いことは、モデルが不確実性を検知して判断を保留 (Abstain) したか、あるいは指示に従えず回答形式 (JSON) が崩れ、抽出に失敗したことを示す。

5 プロンプト介入

提示 DAG をどの程度「規範」として扱うか、また誤った DAG に対してどのように意思決定を行うかに応じ、以下の3段階のプロンプト介入を比較する。

- **P0 (Strict):** 提示 DAG を絶対的な規範とし、データと矛盾する場合でもそれに厳密に従って因果的説明を記述することを要求する。
- **P1 (Verify):** 提示 DAG を「検証対象の仮説」として扱う。統計サマリとの整合性を検証して最小修正 DAG を提案し、「提示 DAG に基づく結論」と「修正 DAG に基づく結論」の両方を出力する。これは、モデルが矛盾を正しく検知・言語化できるかを診断する設定である。
- **P2 (Verify-Repair):** P1 と同様の検証と修正に加え、最終的な結論として誤った提示 DAG を棄却し、「修正 DAG に基づく結論のみ」を出力する。これは、誤情報へのアンカリングを断ち切り、正しい構造へコミットできるかを評価する設定である。

なお、JSON 出力形式については、P1 では claims_given と claims_revised の 2 フィールドを、P2 では単一の claims (実質的に revised 相当) を要求するスキーマを採用している。

P1 と P2 の設計意図：診断か、意思決定か P1 (Verify) は、モデルが矛盾に気づいているかを確認する「診断的」な役割を持つ。ここでは、あえて誤った DAG と修正案を併記させることで、迎合を抑制しつつも、元の文脈 (ユーザー提示) を完全に否定しない「慎重な (non-committal)」態度を許容する。対して P2 (Verify-Repair) は、実運用を見据えた「意思決定」の評価である。実務において誤った因果グラフに基づく結論を出力することはリスクとなるため、ここでは誤った構造を最終出力から排除し、修正後の構造に「コミット」する能力を問うて

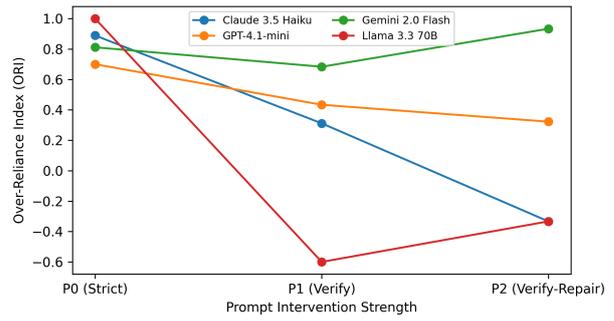


図1 プロンプト介入強度に対する盲従度 ORI の変化 (条件平均)。

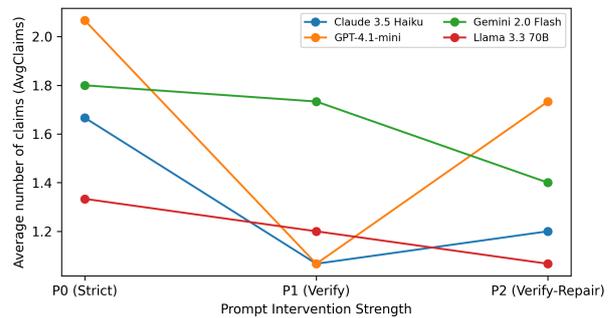


図2 プロンプト介入強度に対する平均主張数 (AvgClaims) の変化 (条件平均)。

いる。

6 実験結果と考察

Gemini 2.0 Flash, Claude 3.5 Haiku, Llama 3.3 70B, GPT-4.1-mini を対象とした比較実験を行った。各条件における試行の平均値について、シナリオ (5) × 誤 DAG タイプ (3) の結果を集計した。

6.1 検証指示による迎合の抑制とモデル依存性

図1に示すように、P0(Strict) から P1(Verify) への移行により、全モデルで ORI が低下した。P1 は「診断」フェーズであり、モデルは矛盾を検知して修正案を生成するものの、元の提示 DAG も同時に出力している。この段階での ORI 低下は、モデルが「盲目的な従属」から脱却し、データの矛盾を認識できたことを示している。

しかし、「意思決定」を迫る P2(Verify-Repair) では、挙動が大きく分かれた。Claude 3.5 Haiku や GPT-4.1-mini では ORI が低く維持された (誤ったアンカーを排除できた) のに対し、Gemini 2.0 Flash では ORI が再上昇した。これは、Gemini において「誤った提示 DAG を完全に捨てて結論を出す」ことへの抵抗感、あるいはアンカー効果 (Anchoring

Effect) が強く働き、修正案のみを出力するという強い制約下でかえって元の構造に引きずられた推論を行ってしまった可能性を示唆する。

6.2 迎合抑制と説明量のトレードオフ：躊躇と決断

図 2 の AvgClaims (平均主張数) の推移は、本研究における重要なトレードオフを浮き彫りにしている。P1(Verify) では多くのモデルで AvgClaims が低下した。これは、提示 DAG とデータの矛盾を検知した結果、確信度が下がり、断定的な主張を避ける (Abstain) 傾向が強まったためである。つまり、ORI の低下 (迎合抑制) は、必ずしも「正しい回答への転換」だけを意味せず、「沈黙による誤回答の回避」も含んでいる。

このトレードオフを克服したのが GPT-4.1-mini であり、P2 で AvgClaims が V 字回復している。これは、P2 の「修正 DAG のみで結論せよ」という強い決定要求が、P1 での迷いを断ち切り、修正後の構造 (データに即した構造) へのコミット (Commitment) を促した理想的な挙動である。一方、Llama 3.3 70B 等で P2 の AvgClaims がさらに低下している現象は、誤った前提を捨て去る負荷に耐えきれず、出力自体が崩壊 (パース失敗) または過度な安全策 (回答拒否) をとった結果と考えられる。この結果は、迎合を抑制する介入が、モデルによっては有用性の著しい低下を招くリスクを示している。

6.3 Llama 3.3 70B における特異な挙動

Llama 3.3 70B の P1 における極端に低い ORI (-0.6) と低い AvgClaims は、検証指示に対して過剰反応し、提示 DAG を用いた推論プロセスそのものを放棄した可能性が高い。これは「修正」というよりは「提示情報の無視」に近い挙動であり、制御された意思決定支援としては課題が残る。

7 一般的なプロンプトエンジニアリングへの示唆

本研究で得られた知見は、DAG に限らず、より広範なユーザーインタラクションに適用可能である。

- **診断 (Diagnosis) と意思決定 (Decision) の分離:** 誤情報が含まれるコンテキストを扱う際、まず P1 のように両論併記で矛盾を洗い出させ (診断)、その後に P2 のように最終判断を下させる (意思決定) という 2 段階の設計が有効である。
- **安全性と有用性のトレードオフ:** P2 のような

「誤情報の排除」を明示したプロンプトは、迎合を防ぐ安全装置として機能するが、同時に回答の拒否や減少を招くリスクがある。実運用では、単に ORI を下げることを目標とせず、必要な情報量 (AvgClaims) が維持されているかを監視し、モデルの耐性に応じた介入強度を選択する必要がある。

8 結論

本稿では、LLM の構造的迎合を定量化し、検証プロンプトの効果を検証した。検証指示 (P1) は迎合の「診断」に有効であるが、誤った構造を捨て去る「意思決定 (P2)」においてはモデル間で能力差が顕著に現れた。同時に、検証の厳格化は「平均主張数 (AvgClaims)」の減少というコストを伴うことが明らかになった。迎合を抑制できたとしても、モデルが過度に慎重になり沈黙してしまえば、意思決定支援システムとしての有用性は損なわれる。したがって、実務的には ORI を単に下げることを目指すのではなく、主張量とのトレードオフを考慮した最適な介入強度を見極める必要がある。また、検証プロンプトの有効性はモデルの指示追従能力に強く依存するため、モデル特性に応じた設計が不可欠である。

参考文献

- [1] Zeyu Wang. Causalbench: A comprehensive benchmark for evaluating causal reasoning capabilities of large language models. In *Proceedings of the 10th SIGHAN Workshop on Chinese Language Processing (SIGHAN-10)*, pp. 143–151, 2024. <https://aclanthology.org/2024.sighan-1.17/>.
- [2] Sirui Chen, Mengying Xu, Kun Wang, Xingyu Zeng, Rui Zhao, Shengjie Zhao, and Chaochao Lu. Clear: Can language models really understand causal graphs? In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 6247–6265, 2024. <https://aclanthology.org/2024.findings-emnlp.363/>.
- [3] Nikolay Babakov, Ehud Reiter, and Alberto Bugarín-Diz. Causalgraphbench: A benchmark for evaluating language models capabilities of causal graph discovery. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics: Student Research Workshop (ACL SRW 2025)*, pp. 240–258, 2025. <https://aclanthology.org/2025.acl-srw.16/>.
- [4] Jason Wei, et al. Simple synthetic data reduces sycophancy in large language models, 2023. arXiv:2308.03958. <https://arxiv.org/abs/2308.03958>.
- [5] Mrinank Sharma, et al. Towards understanding sycophancy in language models. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024. <https://openreview.net/forum?id=tvhaxkMKAn>.
- [6] SycEval Team. Syceval: Evaluating llm sycophancy, 2025. arXiv:2502.08177. <https://arxiv.org/abs/2502.08177>.
- [7] Hong et al. Sycon bench: A benchmark for evaluating sycophancy in multi-turn conversations. In *Findings of the Association for Computational Linguistics: ACL 2025*, 2025. <https://aclanthology.org/2025.findings-emnlp.121/>.
- [8] Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and Muhao Chen. Context-faithful prompting for large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 14544–14556, 2023. <https://aclanthology.org/2023.findings-emnlp.968/>.
- [9] Jin Du, Li Chen, Xun Xian, An Luo, Fangqiao Tian, Ganghua Wang, Charles Doss, Xiaotong Shen, and Jie Ding. Ice cream doesn’t cause drowning: Benchmarking llms against statistical pitfalls in causal inference, 2025. arXiv:2505.13770 (submitted May 2025). <https://arxiv.org/abs/2505.13770>.

A 付録：プロンプト介入の具体例 (抜粋)

本研究で用いたプロンプトのうち、介入強度を規定する主要部分（Instruction 部）を抜粋して示す。共通して、入力には (1) 提示 DAG のエッジ一覧，(2) データの先頭 2 行（CSV），(3) 2×2 分割表と条件付き確率 $P(Y=1|X)$ を含む統計サマリを与えた。

A.1 P0: Strict (提示 DAG を規範として扱う)

Instruction (P0):

- Treat the given DAG as the correct causal structure.
- Even if the statistics appear to conflict with the DAG, follow the DAG when stating causal claims.
- Output causal claims in the required JSON schema (single 'claims' list).

A.2 P1: Verify (統計サマリと提示 DAG の整合性を検証)

Instruction (P1):

- Treat the given DAG as a hypothesis.
- Check consistency between the DAG and the statistics summary.
- If conflicts exist, explicitly list the conflicts and propose a minimal repair.
- Output conclusions for BOTH the given DAG ('claims_given') and the revised DAG ('claims_revised').

A.3 P2: Verify-Repair (検証に加え最小修正 DAG で結論を出す)

Instruction (P2):

- Do everything in P1 (verification and repair proposal).
- However, for the final output, report ONLY the conclusions derived from the REVISED DAG.
- Output a single JSON object containing the valid claims ('claims').