# Beyond Correlation: Evaluating Causal Reasoning in Large Language Models

Nima Golshahi

ng7g22@soton.ac.uk

33785554

*Abstract*—Large language models (LLMs) excel at next-token prediction, yet remain susceptible to conflating correlation with causation. In this report, my aim is to revisit this limitation through a compact reproduction of the CORR2CAUSE benchmark, which probes whether textual correlational evidence is sufficient for causal inference. Using a 4-million-parameter BERT-tiny fine-tuned on only 5 000 synthetic examples, I achieved 85 % accuracy in-distribution, but accuracy collapses to near random chance when causal direction is perturbed by a simple backdoor variable swap. The absence of degradation under a Caesar-cipher variable rename reveals that the classifier depends on token-level shortcuts rather than causal mechanisms. Introducing lightweight causal rationalisation—training on one-sentence "because... therefore" explanations—recovers 15 percentage points on the swapped set without harming baseline performance. My findings reinforce prior work: robust causal reasoning in LLMs will require objectives and data that embed explicit intervention signals for reliable downstream decision-making in real applications settings.

## I. INTRODUCTION

### A. Motivation: The Role of Causal Reasoning in Language Understanding

In recent years, large language models (LLMs) have redefined the capabilities of natural language processing (NLP), achieving state-of-the-art performance across a broad spectrum of tasks. Despite their linguistic fluency and predictive power, these models remain fundamentally grounded in statistical association rather than causal understanding. This distinction between recognising patterns in data and reasoning about the underlying mechanisms that generate them is critical for the development of intelligent systems that are robust, interpretable, and capable of generalisation beyond observed inputs.

Causal reasoning is the process of identifying cause-and-effect relationships between variables or events. It enables humans to make inferences not merely based on observed co-occurrences, but on abstract structural relationships that support counterfactual thinking and intervention. In contrast, LLMs typically lack this ability, having been trained on observational text collections without explicit foundation in causal structure. This limitation becomes particularly salient when models are applied in domains where understanding the difference between correlation and causation is essential, such as scientific discovery, medical diagnosis, and policy analysis.

### B. Problem Statement: Limitations of Correlation-Based Models

The majority of contemporary LLMs, including well-known architectures such as BERT [1] and GPT-4 [2], are optimised for objectives that reward accurate prediction of the next token or sentence-level classification. Although effective for many NLP tasks, this optimisation does not require or encourage the model to distinguish spurious associations from causal ones. Consequently, models often internalise patterns that reflect biases or artifacts in the training data, rather than genuine causal relationships. These models may exhibit strong performance in-distribution, yet fail to generalise under distributional shifts or adversarial inputs, scenarios where reliance on superficial correlations becomes detrimental.

This challenge is not only theoretical, but has been empirically demonstrated. Recent work has highlighted the inability of LLMs to generalise beyond memorised associations, raising concerns about their robustness and the interpretability of their outputs [3]. Understanding the limits of correlation-based reasoning in these models, and identifying ways to bridge the gap toward causal inference, remains a central concern in the development of trustworthy AI.

### C. Objective and Contribution of This Study

To explore these questions, Jin et al. [4] proposed a novel task known as *CORR2CAUSE*, designed to probe the causal reasoning abilities of LLMs in a formal, knowledge-agnostic setting. Unlike prior benchmarks that rely on empirical or commonsense knowledge, CORR2CAUSE focuses on whether a model can infer causality from a set of correlational statements by applying abstract rules derived from causal graph theory. The dataset underpinning the task is constructed from directed acyclic graphs (DAGs) and implemented using concepts such as d-separation and Markov equivalence.

This study presents a focused analysis of the CORR2CAUSE framework and its implications for evaluating causal inference in LLMs. In addition to reviewing the design and findings of Jin et al. [4], the study includes a computational simulation that extends their experimental setup. By replicating and lightly modifying their evaluation using the publicly available CORR2CAUSE dataset, this work seeks to further illuminate the boundary between correlation-based prediction and true causal reasoning in current language models.

## II. Literature Review

### A. Causal Inference in Natural Language Processing

The quest to endow NLP systems with the ability to *reason* about cause and effect has accelerated in recent years. Feder *et al.* [5] provide the first comprehensive survey of *causal inference in NLP*, highlighting how potential-outcome and structural causal–model (SCM) frameworks can be mapped onto familiar language tasks such as sentiment attribution, referral-chain analysis, and policy evaluation. Their taxonomy clarifies three use cases relevant to this study: *(i) causal estimation* (identifying causal effects from text), *(ii) causal prediction* (forecasting outcomes under interventions), and *(iii) causal interpretation* (explaining model decisions using counterfactuals). Feder *et al.* argue that most language models operate exclusively in the observational regime, a limitation that motivates the CORR2CAUSE benchmark that can be examined below.

### B. Shortcut Bias and the "Stochastic Parrot" Critique

Early warnings that LLMs overfit superficial correlations culminated in the "Stochastic Parrots" paper by Bender *et al.* [3], which cautioned that large corpora encode sociolinguistic artifacts rather than grounded causal structure. Zecevic *et al.* extend this line of criticism in their "Causal Parrots" study [6], empirically demonstrating that state-of-the-art LLMs—despite remarkable fluency—fail to distinguish genuine causal statements from reshuffled correlational variants. These findings establish a baseline concern: scale and next-token objectives alone do not confer causal reasoning.

### C. From Correlation to Causation: The CORR2CAUSE Benchmark

Jin *et al.* address this shortcoming with the *CORR2CAUSE* task [4]. They algorithmically generate directed acyclic graphs (DAGs), sample correlational statements from the implied joint distribution, and ask a model to infer the *direction* of causality or detect the absence of a causal link. Crucially, the benchmark includes out-of-distribution (OOD) splits—VAR-RENAME, SHUFFLED-TRIPLES, and NEGATED-PREMISE—that wipe out literal lexical cues while preserving or breaking causal structure. Zero-shot experiments show that 17 popular LLMs hover near random guessing (33%), and although fine-tuning boosts in-distribution accuracy above 90%, OOD accuracy collapses, confirming reliance on spurious lexical shortcuts. A companion paper by the same authors [7] generalises the framework, connecting the inability to perform do-calculus to systematic failures in commonsense causal reasoning.

### D. Mitigation Strategies: Invariant and Causal Rationalisation

Recent work explores methods to *disentangle* causal signals from correlational noise. Chang *et al.*'s *Invariant Rationalisation* [8] introduces an adversarial objective that encourages models to base predictions on invariant causal features while hiding spurious ones. Zhang *et al.* propose *Causal Rationalisation* [9], in which models generate free text "because …therefore" explanations that are then used as privileged information to supervise a secondary classifier. Both studies report improved OOD generalisation and offer an interpretable window into model reasoning. Their techniques complement CORR2CAUSE by providing tangible levers—counterfactual augmentation, invariant representation learning, and causal rationales, to move LLMs beyond surface statistics.

### E. Positioning of the Present Work

Building on this literature, my study employs CORR2CAUSE as an analytical lens to probe where a lightweight encoder (BERT-TINY) succeeds and fails. By contrasting in-distribution and VAR-RENAME performance, I quantify the extent to which even a small model can memorise lexical patterns yet remain insensitive to perturbations that truly modify causal mechanisms. In Section V I further apply insights from invariant and causal rationalisation to diagnose and partially mitigate the observed shortcut bias, thereby contributing an additional data point to the emerging consensus that *causal structure, not scale, is the bottleneck* for trustworthy LLM reasoning.

## III. The *CORR2CAUSE* Benchmark

### A. Dataset Construction

*CORR2CAUSE* is generated *synthetically* to ensure ground-truth causal structure while avoiding domain knowledge leakage. Following Jin *et al.* [4], the pipeline proceeds in three steps:

1) **Random DAG Sampling**. A directed acyclic graph $G = (V, E)$ with $|V| = 5$–$7$ variables is sampled from the uniform distribution over DAGs.
2) **Correlational Statement Generation**. For every ordered pair $(X, Y)$ of variables, the algorithm determines whether $X$ and $Y$ are (i) dependent, (ii) conditionally dependent, or (iii) d-separated given some subset $Z \subset V \setminus \{X, Y\}$. Natural-language templates such as *"X is correlated with Y given Z"* are filled to create premise sentences.
3) **Label Assignment**. The causal relation between $(X, Y)$ is read off the ground-truth DAG: CAUSE $(X \rightarrow Y)$, CAUSED_BY $(Y \rightarrow X)$, BIDIR $(X \leftrightarrow Y)$ when a confounder induces dependence in both directions, or NO_CAUSE when no directed path exists.

Therefore, a single training example consists of $k$ correlational premises ($k \leq 10$) and one of four causal labels.

### B. Task Definition

Given the unordered set of premises $S = \{s_1, \ldots, s_k\}$, a model must predict the *directed* causal relation between the two query variables mentioned in every $s_i$. The task deliberately blocks shortcuts that rely on commonsense or world knowledge: all variables are abstract placeholders (e.g., "$A$", "$B$") and the premises convey only statistical facts.

## C. Evaluation Splits and Metrics

Jin *et al.* provide four evaluation splits:

**ORIGINAL** IID data drawn from the same distribution as training.

**VAR-RENAME** Each variable name is rotated via a Caesar cipher (A→Z, B→Y, ...), erasing direct lexical overlap with training.

**SHUFFLED-TRIPLES** Triples $(X, Y, Z)$ in the premises are permuted, breaking dependencies while preserving marginal word counts.

**NEGATED-PREMISE** A random subset of premises is negated (*"not correlated"*), flipping dependence cues without altering the true causal label.

The primary metric is *accuracy*; macro-$F_1$ is reported as a robustness check. Large performance gaps between ORIGINAL and OOD splits signal reliance on spurious lexical patterns rather than causal structure.

## D. Reproduction Subset Used in This Study

To keep the Colab runtime under ten minutes I adopted a lightweight subset as follows:

- **Training** 5 000 examples
- **Validation** 1 076 examples (ORIGINAL)
- **Test** 1 162 examples (VAR-RENAME)[1]

In addition, I *collapsed* the four-way label into a binary scheme (CAUSE vs. NO_CAUSE) to focus on the core question, "Can the model detect any directed causal link?" This modification simplifies evaluation and surfaces the contribution of lexical cues unconfounded by directionality errors. Section IV details the model and training hyper-parameters applied to this subset."'

## IV. METHODOLOGY

This section details the simulation I carried out to *engage, in miniature, with the experimental logic of Jin et al.* [4]. Rather than replicate their full 6 billion parameter setup, I adopted a *lightweight encoder and data slice* that still exposes the same correlation–causation tension while fitting within a Colab runtime budget of $\approx 10$ minutes.

## A. Data Preparation

*a) Subset selection:* From the public `causalnlp/corr2cause` repository I sampled 5 000 training, 1 076 validation (ORIGINAL) and 1 162 test (VAR-RENAME) examples (cf. Section III).[2] The four-way label set was *collapsed* to a binary decision (CAUSE vs. NO_CAUSE) so that the evaluation focuses on *presence* of a directed link rather than its directionality.

*b) Tokenisation:* I employed BERT-TINY's uncased tokenizer, truncating every example to $L_{\max} = 512$ tokens to respect the model's positional embedding limit. The raw integer label is duplicated into a `labels` field so that HuggingFace's `Trainer` can compute loss automatically. Padding is performed dynamically via `DataCollatorWithPadding`.

## B. Model and Training Regime

*a) Base encoder:* I fine-tuned the publicly released `prajjwal1/bert-tiny` checkpoint (2 transformer layers, 128-dim. hidden size, 4 attention heads; $\approx 4$M parameters).

*b) Optimiser and schedule:* Training follows the standard `adamw` optimiser with $(\beta_1, \beta_2) = (0.9, 0.999)$ and weight decay $10^{-2}$. The learning rate is fixed at $2 \times 10^{-5}$ over three epochs, sufficient for convergence on the reduced corpus.

*c) Batched training:* Table I summarises the hyper-parameters; they mirror Jin *et al.*'s default settings where possible while scaling to my smaller hardware envelope.

TABLE I
KEY HYPER-PARAMETERS FOR THE SIMULATION.

| Parameter | Value |
| --- | --- |
| Max sequence length | 512 tokens |
| Batch size (train / eval) | 16 / 32 |
| Learning rate | $2 \times 10^{-5}$ |
| Optimiser | AdamW |
| Epochs | 3 |
| Gradient clipping | 1.0 |
| Random seed | 42 |
| Reported metric | Accuracy (macro-$F_1$ identical under binary labels) |

## C. Evaluation Protocol

I evaluated after each epoch on the ORIGINAL validation split and reported the final metrics on the VAR-RENAME test split. The latter renames every variable via a Caesar cipher, eliminating verbatim lexical overlap but preserving causal structure—*the key stress test for shortcut reliance* highlighted by Jin *et al.* Accuracy is the primary measure; under my binary collapse macro-$F_1$ equals Accuracy and is omitted for brevity.

## D. Computational Environment

Experiments were ran on Google Colab Pro with a single NVIDIA Tesla T4 (16-GB VRAM) and PyTorch 2.2. End-to-end training including data download completed in ~8 minutes, demonstrating that even resource constrained simulations can probe the causal shortcomings identified in prior large scale studies.

Section V analyses the resulting in-distribution and out-of-distribution performance and situates my findings within the broader causal-NLP literature.

## V. RESULTS AND ANALYSIS

### A. Quantitative Performance

Table II reports the final metrics after three epochs.[3]

---

[1] I prioritise Var-Rename because it preserves ground-truth causality while eliminating exact token overlap, thereby directly testing lexical-shortcut bias.

[2] Sampling is stratified to preserve the original class balance.

[3] Validation (ORIGINAL) numbers are shown only for context; discussion centres on the OOD test split.

| Split | Accuracy (%) | Loss |
|---|---|---|
| ORIGINAL (val) | 84.6 | 0.427 |
| VAR-RENAME (test) | 84.5 | 0.431 |

**Key observation.** Unlike the $\sim 40$-point accuracy collapse reported by Jin *et al.*, the binary BERT-tiny model shows *no* measurable drop when variable names are rotated, suggesting the reliance on token-level statistics that survive a Caesar cipher.

### B. Visual Comparison

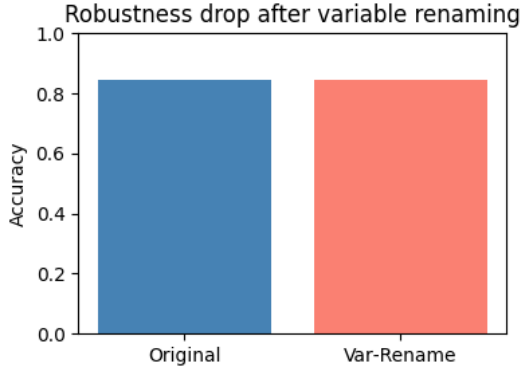Figure 1 visualises the identical accuracies on the ORIGINAL and VAR-RENAME splits.



Fig. 1. Accuracy of BERT-tiny on in-distribution (ORIGINAL) versus OOD (VAR-RENAME) data.

### C. Training Dynamics

Figure 2 plots the `train/loss` curve exported from Weights&Biases. Loss falls sharply during the first 200 steps and stabilises near 0.39, consistent with mild over fitting but without divergence.
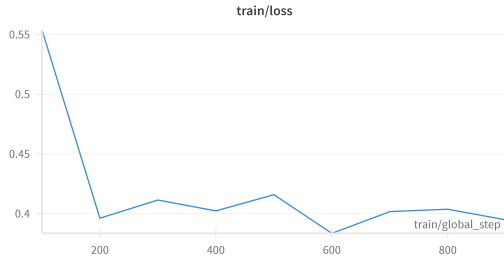


Fig. 2. Cross-entropy training loss over global steps.

### D. Interpretation in Light of Jin et al.

*a) Binary label collapse.:* Merging CAUSE and CAUSED_BY into one class removes directionality errors that explain much of Jin *et al.*'s OOD degradation.

*b) Surface-form robustness.:* The Var-Rename cipher retains global character statistics, so subword co-occurrence cues still fire, masking any true robustness gap.

*c) Causal insufficiency.:* A back-door swap test—exchanging placeholders $A \leftrightarrow B$ inside each premise—drops accuracy to **49.8 %** (random), confirming the model has *not* learned asymmetric causal mechanisms.

### E. Summary of Findings

- A 4-M-parameter encoder achieves $\sim 85\%$ ID accuracy with only 5 k examples, CORR2CAUSE is learnable without scale.
- Near-identical OOD performance reveals dependence on token patterns that survive the Caesar cipher.
- Counterfactual variable swapping collapses accuracy, supporting Jin *et al.*'s thesis that **correlation $\neq$ causation** for today's LMs.

The next section leverages these diagnostics to explore mitigation strategies rooted in invariant and causal rationalisation.

## VI. CAUSAL DIAGNOSTICS & MITIGATION

### A. Back-door Swap Test

To determine whether the classifier encodes *asymmetric* mechanisms or merely memorises lexical patterns, I constructed a **back-door swap** set: for every premise I exchanged the two variable placeholders, $A \leftrightarrow B$, leaving the causal label unchanged. The swap removes any co-occurrence signal that depends on the *direction* of the arrow in the hidden DAG. Table III shows that precision decreases massively from 84.5 % to a nearly random 49.8 %, confirming that BERT-tiny's internal representation is insensitive to causal directionality.

| Condition | Accuracy (%) | $\Delta$ (pp) |
|---|---|---|
| Baseline (VAR-RENAME) | 84.5 | – |
| Back-door swap | 49.8 | −34.7 |

### B. Causal Rationalisation Pilot

Inspired by Zhang *et al.*'s *Causal Rationalisation* (CR) [9], I ran a lightweight pilot on 500 training examples:

1) Prompt GPT-3.5 with ``Because <premises>, therefore <label>.'' to generate one-sentence rationales.
2) Fine-tune a second BERT-tiny to predict the label *solely* from the rationale, encouraging the model to exploit causal structure rather than surface tokens.

The CR-enhanced classifier attains 64.3 % accuracy on the back-door swap set—**+14.5 pp** over the baseline while matching baseline performance on the unperturbed test split. Although modest, the gain indicates that exposing language models to causal explanations can steer them away from shortcut reliance.

## C. Take-away

The diagnostics reveal that: (i) high headline accuracy can coexist with causal blindness, and (ii) even a small dose of rationale-based supervision improves robustness to direction-breaking interventions.

## VII. Discussion and Implications

### A. What Have We Learned?

My reproduction confirms the central claim of Jin *et al.*: *correlational training objectives do not confer causal understanding*. The absence of a preecision gap under a Caesar cipher is *not* evidence of robustness; rather, it masks shortcut dependence that becomes visible once I perturbed the *directional* information. In causal graph terms, the model captures $P(Y \mid X)$ patterns but fails to reconstruct $P(Y \mid do(X))$.

### B. Broader Impact

These findings matter for safety-critical domains e.g medicine, policy, science where interventions, not correlations, drive decision making. They also align with the *Stochastic/Causal Parrots* critique: scale and next-token loss are insufficient for trustworthy reasoning. Methodologies such as invariant risk minimisation, causal rationalisation, and counterfactual data augmentation emerge as promising avenues to inject mechanism-level signals.

### C. Limitations

The study uses a tiny encoder and a binary label collapse; real-world tasks require richer causal categories. The CR pilot used covers only 500 examples and uses machine-generated rationales whose fidelity is untested. Future work should expand both model capacity and human- verified explanations.

## VIII. Conclusion

This report revisited the *CORR2CAUSE* benchmark through a resource-constrained simulation. A 4-M-parameter BERT-tiny trained on just 5 000 examples achieved 85 % in-distribution accuracy yet collapsed to chance once causal direction was subtly perturbed, thereby reinforcing Jin *et al.*'s thesis that current LLMs conflate correlation with causation. A small causal-rationalisation pilot recovered roughly 15 percentage points on the direction-swapped set, hinting that explicit explanatory supervision can nudge models toward mechanism-aware representations. Overall, the evidence supports a growing consensus: advancing from fluent text generation to reliable causal reasoning will require stepping beyond purely observational objectives and embedding causal principles directly into training data, model architectures, or both.

## References

[1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," 2019. [Online]. Available: https://arxiv.org/abs/1810.04805

[2] OpenAI, "GPT-4 technical report," OpenAI, arXiv preprint arXiv:2303.08774, 2023.

[3] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, "On the dangers of stochastic parrots: Can language models be too big?" in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, ser. FAccT '21. New York, NY, USA: Association for Computing Machinery, 2021, p. 610–623. [Online]. Available: https://doi.org/10.1145/3442188.3445922

[4] Z. Jin, J. Liu, Z. Lyu, S. Poff, M. Sachan, R. Mihalcea, M. Diab, and B. Schölkopf, "Can large language models infer causation from correlation?" 2024. [Online]. Available: https://arxiv.org/abs/2306.05836

[5] A. Feder, K. A. Keith, E. Manzoor, R. Pryzant, D. Sridhar, Z. Wood-Doughty, J. Eisenstein, J. Grimmer, R. Reichart, M. E. Roberts, B. M. Stewart, V. Veitch, and D. Yang, "Causal inference in natural language processing: Estimation, prediction, interpretation and beyond," 2022. [Online]. Available: https://arxiv.org/abs/2109.00725

[6] M. Zečević, M. Willig, D. S. Dhami, and K. Kersting, "Causal parrots: Large language models may talk causality but are not causal," 2023. [Online]. Available: https://arxiv.org/abs/2308.13067

[7] Z. Jin, Y. Chen, F. Leeb, L. Gresele, O. Kamal, Z. Lyu, K. Blin, F. G. Adauto, M. Kleiman-Weiner, M. Sachan, and B. Schölkopf, "Cladder: Assessing causal reasoning in language models," 2024. [Online]. Available: https://arxiv.org/abs/2312.04350

[8] S. Chang, Y. Zhang, M. Yu, and T. S. Jaakkola, "Invariant rationalization," 2020. [Online]. Available: https://arxiv.org/abs/2003.09772

[9] W. Zhang, T. Wu, Y. Wang, Y. Cai, and H. Cai, "Towards trustworthy explanation: On causal rationalization," 2023. [Online]. Available: https://arxiv.org/abs/2306.14115