

How does varying the complexity of Structural Causal Models in causal fine-tuning affect the OOD F1 score stab

Assignee Research

June 10, 2026

Abstract

The causal capabilities of large language models (LLMs) are a matter of significant debate, with critical implications for the use of LLMs in societally impactful domains such as medicine, science, law, and policy. We conduct a "behaviorial" study of LLMs to benchmark their capability in generating causal arguments. Across a wide range of tasks, we find that LLMs can generate text corresponding to correct causal arguments with high probability, surpassing the best-performing existing methods. Algorithms based on GPT-3.5 and 4 outperform existing algorithms on a pairwise causal discovery task (9

1 Introduction

This paper examines: Causal Reasoning and Large Language Models: Opening a New Frontier for Causality. Research question: How does varying the complexity of Structural Causal Models in causal fine-tuning affect the OOD F1 score stability of tabular foundation models on benchmarks like TabFact?.

2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.1/10.

3 Results

12 papers retrieved. 7 claims extracted; 7 independently verified. Quality review score: 8.1/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
Algorithms based on GPT-3.5 and 4 outperform existing algorithms on a pairwise causal discovery task (97%, 13 points gain)	✓	0.31
Algorithms based on GPT-3.5 and 4 outperform existing algorithms on a counterfactual reasoning task (92%, 20 points gain)	✓	0.28
Algorithms based on GPT-3.5 and 4 achieve 86% accuracy in determining necessary and sufficient causes in vignettes.	✓	0.21
LLMs generalize to novel datasets that were created after the training cutoff date.	✓	0.21
LLMs exhibit unpredictable failure modes.	✓	0.17
LLMs bring capabilities so far understood to be restricted to humans, such as using collected knowledge to generate caus	✓	0.36
LLMs may be used by human domain experts to save effort in setting up a causal analysis.	✓	0.25

References

- <https://doi.org/10.1016/j.inffus.2023.101805>
- <https://doi.org/10.48550/arxiv.2305.00050>
- <https://doi.org/10.1038/s41586-024-08328-6>