

Causality for Natural Language Processing

Zhijing Jin

2024

*To my parents, Xuefang Jin and Jian Wu,
for their unconditional love and endless inspirations to me.*

Abstract

Causal reasoning is a cornerstone of human intelligence and a critical capability for artificial systems aiming to achieve advanced understanding and decision-making. This thesis delves into various dimensions of causal reasoning and understanding in large language models (LLMs). It encompasses a series of studies that explore the causal inference skills of LLMs, the mechanisms behind their performance, and the implications of causal and anticausal learning for natural language processing (NLP) tasks. Additionally, it investigates the application of causal reasoning in text-based computational social science, specifically focusing on political decision-making and the evaluation of scientific impact through citations. Through novel datasets, benchmark tasks, and methodological frameworks, this work identifies key challenges and opportunities to improve the causal capabilities of LLMs, providing a comprehensive foundation for future research in this evolving field.

Zusammenfassung

The "Zusammenfassung" is a machine-translated version of the abstract via <https://deepl.com> and corrected with the help of Jan Schneider.

Kausales Denken ist ein Grundpfeiler menschlicher Intelligenz und eine entscheidende Fähigkeit für künstliche Systeme, die ein fortgeschrittenes Verständnis und fundierte Entscheidungsfindung anstreben. Diese Arbeit untersucht verschiedene Aspekte des kausalen Denkens und Verstehens in großen Sprachmodellen, nämlich LLMs. Sie umfasst eine Reihe von Studien, die die Fähigkeiten von LLMs zur kausalen Inferenz, die Mechanismen hinter ihrer Leistungsfähigkeit und die Auswirkungen von kausalem und antikausalem Lernen auf Aufgaben der natürlichen Sprachverarbeitung, oder NLP beleuchten. Darüber hinaus wird die Anwendung kausalen Denkens in der textbasierten computergestützten Sozialwissenschaft untersucht, wobei der Fokus auf politischer Entscheidungsfindung und der Bewertung des wissenschaftlichen Einflusses durch Zitationen liegt. Durch neuartige Datensätze, Benchmark-Aufgaben und methodische Rahmenbedingungen identifiziert diese Arbeit zentrale Herausforderungen und Chancen zur Verbesserung der kausalen Fähigkeiten von LLMs und bietet eine umfassende Grundlage für zukünftige Forschung in diesem sich entwickelnden Bereich.

Acknowledgements

Throughout my PhD, I am deeply grateful to Professor Bernhard Schölkopf for being a constant source of inspiration and a remarkable role model as a scientist. His foundational work on causal inference for machine learning inspired me to start my own journey of causality for NLP. His broad knowledge across AI, statistics, physics, and philosophy often bring interesting insights in the way we do science. Moreover, as an incredible supervisor, Bernhard creates one of the best nurturing environments to do curiosity-driven research, gives me extensive support, and keeps inspiring with his sharp insights.

I really appreciate my co-supervisor, Professor Mrinmaya Sachan, for his detailed mentorship, unwavering support, and guidance throughout my PhD. His close supervision during the design of my research agenda and teaching of how to write well have been crucial to my success. Thanks to him, I have enjoyed a resourceful and fulfilling PhD experience across Germany and Switzerland.

I owe endless thanks to Professor Rada Mihalcea, whose loving and inspiring presence has been invaluable to me. She is one of the kindest people I have ever met, and being her mentee is one of the best experiences a person could ever imagine. I appreciate our shared vision and passion, which drive us to establish various lines of research on NLP for Social Good, and the ACL Year-Round Mentorship program for our community.

I have also been fortunate to collaborate with and receive mentorship from several people. I really admire Mona Diab for her caring nature and strong sense of social responsibility, which made me really enjoy our responsible AI research together. Collaborating with Kun Zhang has shown me the brilliance of technical causal inference research and the importance of responsible community service. Ryan Cotterell has taught me how to connect formal methods with NLP, and how to write precise, beautiful academic papers. Working with Asli Celikyilmaz at Meta has taught me novel idea generation and efficient project management to bring those ideas to realization.

My PhD journey has been supported by Max Tegmark's recognition of my research agenda and his belief in my potential to contribute greatly to AI safety. The fellowship and funding support from his Future of Life Institute have significantly expedited my research. Collaborating with and learning from Yejin Choi has shown me what it means to be a great NLP researcher with a sharp mind and how to clearly position papers within the community. A final big moment at the end of my PhD was my encounter with Geoffrey Hinton, who identified my potential and encouraged me to apply to the University of Toronto, where I am fortunate to become an assistant professor, in the hometown of deep learning.

Rome wasn't built in a day, and I owe a lot of thanks to the people who have guided me into research. During my undergrad years, I was fortunate to receive smart and patient side-by-side guidance from Di Jin, who was a PhD student in Peter Szolovits's lab at MIT CSAIL at the time. His ability to identify research problems and work through technical methods led to my first burst of papers and, more importantly, helped me understand what research truly is. My second burst of papers was inspired by Qipeng Guo during our time working together at Zheng Zhang's Amazon AI lab. Qipeng, a senior PhD student back then, was excellent in teaching, philosophical thinking, and research skills. It was during this period that I began to develop my research tastes. As I began my PhD

in causal inference, I am deeply grateful to Julius von Kügelgen and Luigi Gresele, who have strong technical backgrounds and kindly spent tremendous one-on-one time exploring the multifaceted beauty of statistical causal inference with me. Another turning point in the middle of my PhD was my collaboration with Sydney Levine and Max Kleiman-Weiner, who sparked my strong curiosity about moral reasoning. This curiosity evolved into one of the three pillars of my knowledge enterprise, alongside logic and causality. To deepen my exploration of these three pillars, the many inspiring discussions with Felix Leeb and Ari Holtzman, fueled by our shared interest in the philosophy of science and understanding how things around us work, have been invaluable. At the end of my PhD in 2024, I had great fortune to deliver several talk tours which gave the opportunity for me to talk to many profound thinkers, and I was intensively inspired by the chats with Peter Spirtes, Judea Pearl, Geoffrey Hinton, Yoshua Bengio, Caroline Uhler, Dominik Janzing, Elias Bareinboim, and Eric Lander, for many aspects of causality, from its philosophical foundation, mathematical formalization, to interdisciplinary applications. As an important source of knowledge, I feel really grateful to many life-changing books I encountered during my PhD, especially Bertrand Russell's *The History of Western Philosophy*, and Walter Isaacson's biography of *Leonardo da Vinci*, which opens my awareness of the fundamental methodology that people use for thinking.

I want to acknowledge all my talented mentees that I am very fortunate to have advised and worked with: Abhinav Lalwani, Ahmad Khan, Amélie Reymond, Andreas Opedal, András Strausz, Ayush Kaushal, David Guzman, David Jenny, Dmitrii Kharlapenko, Fernando Gonzalez, Flavio Schneider, Francesco Ortu, Giorgio Piatti, Harshvardhan Srivastava, Hongyuan Liu, Ishan Kumar Agrawal, Jad Beydoun, Jiarui Liu, Jiayi Zhang, Jingwei Ni, João Afonso Miguel, Justus Mattern, Kevin Blin, Lea Künstler, Luise Woehlke, Navreet Kaur, Neemesh Yadav, Nils Heil, Ojasv Kamal, Roberto Ceraolo, Rongwu Xu, Samuel Simko, Sawal Acharya, Steven Wang, Sumon Kanti Dey, Tejas Vaidhya, Vetha Raghuram, Xiaoyu Xing, Yahang Qi, Yihuai Hong, Yiwen Ding, Yongyi Yang, Yuchun Dai, Yuen Chen, Yuxin Ren, and Zhiheng Lyu.

Science could not be made possible without the people behind. At Max Planck Institute for Intelligent Systems and ETH, I received lots of help from directors/professors, colleagues, our incredible administrative assistants, Sabrina, Ann-Sophie, Lidia, Linda, Patrizia, Paulina, and our research engineer Vincent, along with many others. I am also grateful to Professors Georg Martius and Seong Joon Oh for kindly being on my thesis committee, and administrative support at the University of Tuebingen.

Finally, my PhD years would not have been as inspiring without the 24/7 support of my dad. His intelligence and love provided not only emotional support but also practical help in solving unforeseen out-of-distribution problems beyond both our life experiences. As my PhD coincided with the COVID-19 period challenging for physical and mental health, the optimistic and resilient character you see in me today owes to my parents' unconditional love and the cherished friendship with my 90+ year-old German neighbor, Reinhard Beyer.

Publications

The following peer-reviewed publications are at the core of my PhD research and covered in this dissertation:

Equal Contribution (*); Co-Supervision (†); Mentees I Advised (*Name*); Oral (🗣️).

Part I. Causal Reasoning in LLMs

- **Can Large Language Models Infer Causation from Correlation?**
Zhijing Jin*, *Jiarui Liu**, *Zhiheng Lyu*, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona Diab†, Bernhard Schölkopf†. [ICLR 2024](#)
- **CLadder: Assessing Causal Reasoning in Language Models**
Zhijing Jin*, *Yuen Chen**, Felix Leeb*, Luigi Gresele*, *Ojasv Kamal*, *Zhiheng Lyu*, *Kevin Blin*, *Fernando Gonzalez*, Max Kleiman-Weiner, Mrinmaya Sachan, Bernhard Schölkopf. [NeurIPS 2023](#)

Part II. Causal Understanding of How LLMs Work

- **Competition of Mechanisms: Tracing How Language Models Handle Facts and Counterfactuals**
*Francesco Ortu**, **Zhijing Jin***, Diego Doimo, Mrinmaya Sachan, Alberto Cazzaniga†, Bernhard Schölkopf†. [ACL 2024](#)
- **A Causal Framework to Quantify the Robustness of Mathematical Reasoning in Language Models**
Alessandro Stolfo*, **Zhijing Jin***, Kumar Shridhar, Bernhard Schölkopf, Mrinmaya Sachan. [ACL 2023](#)

Part III. Causality among the Learning Variables

- **Causal Direction in Data Collection Matters: Implications of Causal and Anticausal Learning in NLP**
Zhijing Jin*, Julius von Kuegelgen*, *Jingwei Ni*, *Tejas Vaidhya*, *Ayush Kaushal*, Mrinmaya Sachan, Bernhard Schölkopf. [EMNLP 2021](#). 🗣️ **Oral Presentation**
- **On the Causal Nature of Sentiment Analysis**
*Zhiheng Lyu**, **Zhijing Jin***, *Fernando Gonzalez*, Rada Mihalcea, Bernhard Schölkopf, Mrinmaya Sachan. [Findings of EMNLP 2024](#)

Part IV. Causality for Text-Based Computational Social Science

- **Mining the Cause of Political Decision-Making from Social Media: A Case Study of COVID-19 Policies across the US States**
Zhijing Jin, Zeyu Peng, *Tejas Vaidhya*, Bernhard Schölkopf, Rada Mihalcea. [Findings of EMNLP 2021](#)
- **CausalCite: A Causal Formulation of Paper Citations**
*Ishan Kumar**, **Zhijing Jin***, Ehsan Mokhtarian, Siyuan Guo, *Yuen Chen*, Mrinmaya Sachan, Bernhard Schölkopf. [Findings of ACL 2024](#)

My entire list of publications is available at

Google Scholar: scholar.google.com/citations?user=Mdr6wjUAAAAJ

Citations: 3,266 **h-index:** 23 (by December 2024)

Contents

1	Introduction	1
1.1	Overview	1
1.2	Structure of the Thesis	4
1.3	Contributions and Impact	4
I	Causal Reasoning in LLMs	5
2	Can LLMs Infer Causation from Correlation?	7
2.1	Introduction	7
2.2	Preliminaries: Causal Inference	9
2.3	Dataset Construction	10
2.4	Experiments	14
2.5	Related Work	18
2.6	Conclusion	18
3	CLadder: Assessing Causal Reasoning in LLMs	20
3.1	Introduction	20
3.2	Preliminaries on Causal Inference	22
3.3	Composing the CLADDER Dataset	24
3.4	Our CAUSALCoT Model	28
3.5	Testing LLMs with CLADDER	30
3.6	Related Work	32
3.7	Discussion of Limitations and Future Work	32
3.8	Conclusion	33
II	Causal Understanding of How LLMs Work	34
4	Competition of Mechanisms: Tracing How LLMs Handle Facts and Counter-factuals	36
4.1	Introduction	36
4.2	Related Work on Interpretability	38
4.3	Problem Setup	39
4.4	Method and Background	39
4.5	Experimental Setup	40
4.6	Results and Findings	41
4.7	Discussion and Future Work	47
4.8	Conclusion	48

5	A Causal Framework to Quantify the Robustness of Mathematical Reasoning in LLMs	50
5.1	Introduction	50
5.2	Problem Setup	52
5.3	A Causal Framework	53
5.4	Experimental Setup	57
5.5	Results	58
5.6	Related Work	61
5.7	Conclusion	62
III	Causality among the Learning Variables	67
6	Causal Direction in Data Matters: Implications of Causal and Anticausal Learning in NLP	69
6.1	Introduction	69
6.2	Categorization of Common NLP Tasks	71
6.3	Implications of ICM for Causal and Anticausal Learning	72
6.4	Validity of ICM for NLP Data Using MDL	74
6.5	SSL for Causal vs. Anticausal Models	77
6.6	DA for Causal vs. Anticausal Models	79
6.7	How to Use the Findings in this Study	80
6.8	Limitations and Future Work	80
6.9	Related Work	81
6.10	Conclusion	82
7	On the Causal Nature of Sentiment Analysis	83
7.1	Introduction	83
7.2	Problem Formulation of SA	85
7.3	Causal Discovery of Sentiment and Review	86
7.4	Improving Sentiment Classifiers with Causal Alignment	91
7.5	Related Work	93
7.6	Conclusion	94
IV	Causality for Text-Based Computational Social Science	96
8	Mining the Cause of Political Decision-Making from Social Media: A Case Study of COVID-19 Policies across the US States	98
8.1	Introduction	98
8.2	Related Work	100
8.3	Governor-Targeted Public Opinion	101
8.4	Collection of Policies and Confounders	103
8.5	Mining Decisive Factors of COVID Policies	105
8.6	Fine-Grained Analyses	109
8.7	Additional Discussions	110
8.8	Conclusion	111
9	CausalCite: A Causal Formulation of Paper Citations	113
9.1	Introduction	113
9.2	Problem Formulation	115

9.3	The TEXTMATCH Method	117
9.4	Performance Evaluation	122
9.5	Findings	124
9.6	Related Work	126
9.7	Conclusion	127
10	Conclusion & Future Work	129
A	Appendix	132
A.1	Additional Materials for Chapter 2	132
A.2	Additional Materials for Chapter 3	135
A.3	Additional Materials for Chapter 4	151
A.4	Additional Materials for Chapter 5	153
A.5	Additional Materials for Chapter 6	160
A.6	Additional Materials for Chapter 7	161
A.7	Additional Materials for Chapter 8	170
A.8	Additional Materials for Chapter 9	174

Introduction

1.1 Overview

Causality is a fundamental aspect of human cognition and intelligence, underlying our understanding of the world and our ability to make decisions. In the field of natural language processing (NLP), the capability to infer and reason about causality is increasingly recognized as a critical component of intelligent systems.

Despite the recent advancement of large language models (LLMs) (Radford et al., 2019; Devlin et al., 2019; Brown et al., 2020; Zhang et al., 2022; OpenAI, 2023; Ignat et al., 2024, *inter alia*), a key question still remains: Can these models understand and reason about causality? This is a critical skill before we can trust AI agents to be integrated into decision-making systems. Moreover, even if LLMs succeed at some extent of reasoning, they still lack transparency of how their decisions are made, forming a strong need for interpretability (Luo and Specia, 2024; Räuker et al., 2023; Zou et al., 2023).

To bridge the gap, this thesis explores various facets of causal reasoning in LLMs. We present a series of studies that collectively advance the knowledge of how well these models perform causal reasoning (Part I), how their decisions are made (Part II), how causality among learning variables influences NLP tasks (Part III), and how causality and NLP can together analyze social problems (Part IV). Below we introduce an overview of the four parts and their corresponding chapters.

1.1.1 Causal Reasoning in LLMs (Part I)

In Part I, we investigate two formal causal reasoning skills in LLMs: causal discovery (Chapter 2) and causal effect reasoning (Chapter 3), which existing models struggle to address. We propose these pure reasoning tasks independent of empirical knowledge, and contribute symbolically grounded datasets that are carefully constructed. Based on our proposed data and test pipeline, we first look into the initial performance of LLMs off the shelf, explore how fine-tuning improves their performance, and propose chain-of-thought reasoning to ground the inference skills of LLMs in formal steps.

Causal Discovery in LLMs (Chapter 2) The ability to distinguish between correlation and causation is crucial for intelligent AI systems. In Chapter 2, we propose a benchmark dataset, CORR2CAUSE, specifically designed to test the causal discovery skills of LLMs. This novel task requires models to determine causal relationships from a set of correlational statements. Our large-scale dataset consists of over 200K samples, on which we evaluate seventeen existing LLMs. The findings reveal a significant shortcoming: these models perform close to random on this task, indicating a fundamental gap in their causal inference abilities. Although fine-tuning improves performance somewhat, the models still struggle with generalization, performing well only on in-distribution settings but failing on out-of-distribution queries. This study underscores the challenges in guiding

future research to enhance the pure reasoning skills and generalizability of LLMs.

Causal Effect Reasoning in LLMs (Chapter 3) To investigate the other skill, causal effect reasoning, we propose a new task in Chapter 3 inspired by the “causal inference engine” concept postulated by Pearl and Mackenzie (2018). Our dataset contains 10K samples derived from causal graphs and queries, which are translated into natural language. We evaluate LLMs using this dataset and introduce a chain-of-thought prompting strategy, CAUSALCoT. The results highlight the challenges LLMs face in causal reasoning, providing deep insights into their limitations and suggesting directions for future improvements.

1.1.2 Causal Understanding of How LLMs Work (Part II)

Beyond understanding the causal reasoning skills in LLMs, a natural following question is to interpret how models make their decisions, which leads to our exploration in Part II. We look into two types of interpretation, one to inspect and intervene their internal states, called intrinsic interpretability (Chapter 4), and the other to perturb the input-output space to capture behavioral tendencies, called behavioral interpretability (Chapter 5). We will introduce each in the following.

Intrinsic Interpretability (Chapter 4) We present a novel formulation to understand the inner mechanisms of LLMs. Most existing research focuses on analyzing a single mechanism, such as how models copy or recall factual knowledge. In this work, we propose a formulation of *competition of mechanisms*, which focuses on the interplay of multiple mechanisms instead of individual mechanisms and traces how one of them becomes dominant in the final prediction. We uncover how and where mechanisms compete within LLMs using two interpretability methods: logit inspection and attention modification. Our findings show traces of the mechanisms and their competition across various model components and reveal attention positions that effectively control the strength of certain mechanisms.

Behavioral Interpretability (Chapter 5) Different from previous work that uses heuristics to perform behavioral tests, we systematically propose a pipeline to probe the difference between the desired causal mechanisms, and model-learned causal mechanisms. To implement our framework, we focus on mathematical reasoning problems. By grounding our analysis in a causal graph, we assess the causal effect of different input factors such as problem text, operands, and operators on the model’s output. This framework enables a detailed analysis of the models’ robustness and highlights the critical role of causal factors in shaping their responses. The study reveals that robustness does not continuously improve with model size; however, the GPT-3 Davinci model (175B) shows significant advancements in both robustness and sensitivity compared to other variants. Despite the relative improvement, our results indicate that there is still substantial room for enhancing the robustness and generalization capabilities of LLMs through a better understanding of their causal mechanisms.

1.1.3 Causality among the Learning Variables (Part III)

Apart from causal inference to improve performance and to interpret the models, we also explore the causal relationships between the input and output variables in NLP tasks in

Part III, which provides valuable insights into the design and evaluation of models. This part includes studies that examine the implications of causal and anticausal learning and their impact on NLP tasks (Chapter 6), and also provide an example discovering the causal relationship between the input-output learning variables (Chapter 7).

Implications of Causal and Anticausal Learning in NLP (Chapter 6) We introduce the principle of independent causal mechanisms (ICM) to the NLP community. By inspecting the causal directions between the input and output variables, we categorize common NLP tasks into clearly causal (where the input variable serves as a cause of the output variable, or is collected before the annotation of the output), anticausal (where the output variable serves as a cause of the input variable, or is collected before the annotation of the input), or mixed. After an extensive meta-analysis of over 130 published studies, the paper demonstrates that the causal direction of data collection significantly affects learning outcomes. This work is the first to apply the ICM principle to NLP and offers constructive suggestions for future modeling choices based on causal insights.

Discovering Causal and Anticausal Sentiment Analysis (Chapter 7) We further extend this exploration to cases where the variable causal relations are not evident *a priori*, but discovered using interdisciplinary insights. In Chapter 7, we reformulate sentiment analysis (SA) a combination of causal discovery and prediction tasks. For the first task, we discover the causal relation in SA tasks, namely between a review and its sentiment, by using the peak-end rule from psychology. In this way, we classify reviews based on whether the review primes the sentiment or vice versa. Knowing the causal direction, we inspect how it affects the prediction performance of LLMs when the prediction task aligns with the same causal direction or is opposite. To improve model performance, we construct causal prompts that reflect the underlying causal graph, which lead to substantial improvements, demonstrating the importance of considering causal mechanisms in SA tasks.

1.1.4 Causality for Text-Based Computational Social Science (Part IV)

Finally, we explore in Part IV the application of causal inference to text-based computational social science to address real-world problems. We conduct two studies that examine the causes behind political decision-making (Chapter 8) and paper citations (Chapter 9).

Causal Policy Analysis (Chapter 8) As the first application study, we investigate how public opinion on social media influences policy decisions during the COVID-19 pandemic. By analyzing textual Twitter data and controlling for confounders such as case increases and unemployment rates, the study conducts causal inference to identify trends in political decision-making across different states. This research highlights the dynamic interaction between public sentiment and political actions in a rapidly evolving context.

Causal Analysis of Paper Citations (Chapter 9) For the second case study of paper citations, we propose a new method to evaluate the significance of scientific papers through causal impact on subsequent research. Using TEXTMATCH, a novel causal inference method adapted for high-dimensional text embeddings, we assess the causal influence of papers on their follow-ups. The effectiveness of our framework is demonstrated through various

criteria, providing a more accurate measure of a paper's true impact and suggesting ways for future researchers to utilize this metric.

1.2 Structure of the Thesis

This thesis is organized into four main parts, corresponding to the thematic areas outlined above. Each part includes a detailed examination of the relevant studies, methodologies, and findings, offering a comprehensive view of causal methods for NLP:

1. **Causal Reasoning in LLMs** (Part I): Assessing the ability of LLMs on causal inference by introducing formal benchmarks, and improving the model performance by chain-of-thought reasoning and finetuning.
2. **Causal Understanding of How LLMs Work** (Part II): Investigating the internal mechanisms and robustness of LLMs by causal interventions and causal mediation analysis.
3. **Causality among Learning Variables** (Part III): Exploring the causal or anticausal relation among the learning variables, and drawing insights to understand and improve model performance.
4. **Causality for Text-Based Computational Social Science** (Part IV): Applying causal inference on text data to uncover social and political insights.

1.3 Contributions and Impact

This body of work contributes to the growing field of causal inference in machine learning, providing critical insights into the capabilities and limitations of LLMs in reasoning about causality. By using causal methods to intervene on LLMs, we pave the way for developing more robust and interpretable AI systems. Moreover, the applications in computational social science highlight the broader societal implications of causal reasoning with text data, offering new tools and perspectives for analyzing complex social phenomena.

Part I

Causal Reasoning in LLMs

Can LLMs Infer Causation from Correlation?

Causal inference is the study that discovers cause-effect relationships from interventional and/or observational data (Pearl, 2009b; Spirtes et al., 2000; Zhang and Hyvärinen, 2009). Causality started as a philosophical subject (Beebe et al., 2009; Russell, 2004; Kant, 1781), and in the recent centuries integrated into statistics with established concepts and tools (Fisher and Ford, 1927; Rubin, 1980; Spirtes et al., 1993; Pearl, 2009b). In the era of LLMs, we formulate such formally-grounded pure causal inference as a target reasoning capability for language models.

In this work, we propose the first benchmark dataset to test the pure causal inference skills of large language models (LLMs). Specifically, we formulate a novel task `CORR2CAUSE`, which takes a set of correlational statements and determines the causal relationship between the variables. We curate a large-scale dataset of more than 200K samples, on which we evaluate seventeen existing LLMs. Through our experiments, we identify a key shortcoming of LLMs in terms of their causal inference skills, and show that these models achieve almost close to random performance on the task. This shortcoming is somewhat mitigated when we try to re-purpose LLMs for this skill via finetuning, but we find that these models still fail to generalize – they can only perform causal inference in in-distribution settings when variable names and textual expressions used in the queries are similar to those in the training set, but fail in out-of-distribution settings generated by perturbing these queries. `CORR2CAUSE` is a challenging task for LLMs, and can be helpful in guiding future research on improving LLMs’ pure reasoning skills and generalizability. Our data is available at <https://huggingface.co/datasets/causalnlp/corr2cause>, and our code is at <https://github.com/causalNLP/corr2cause>.

2.1 Introduction

Causal inference, i.e., the ability to establish the correct causal relationships between variables or events, is fundamental to human intelligence. There are two distinct ways this causal inference capability can be acquired: one through empirical knowledge, e.g., we know from common sense that touching a hot stove will get us burned; the other through *pure causal reasoning*, as causality can be formally argued and reasoned about using known procedures and rules from causal inference (Spirtes et al., 2001; Pearl, 2009b; Peters et al., 2017). One example is that we have the a priori knowledge that the correlation between A and B does not necessarily imply causality. This is a formal rule that holds true regardless of the realizations of the variables A and B.

With the rise of large language models (LLMs) (Radford et al., 2019; Devlin et al., 2019; Ouyang et al., 2022; Zhang et al., 2022; OpenAI, 2023, *inter alia*), a crucial research question is whether they can do causal reasoning well. Recent studies have pointed out that LLMs are “causal parrots,” which recite the causal knowledge in the training data (Zeče-

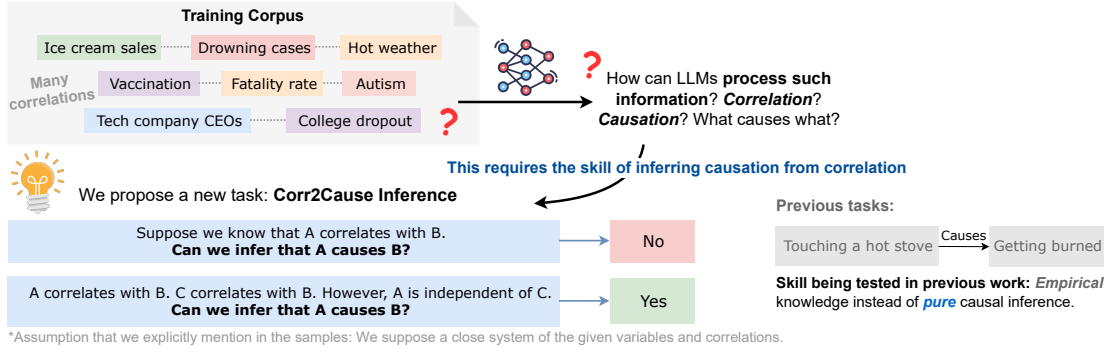


Figure 2.1: Illustration of the motivation behind our task and dataset.

vić et al., 2023). Moreover, the vast majority of studies frame causal reasoning as a skill to navigate around empirical knowledge (Gordon et al., 2012; Sap et al., 2019a,b; Qin et al., 2019; Bhagavatula et al., 2020), and also treat LLMs as a knowledge base when evaluating its causal skills (Kıcıman et al., 2023; Tu et al., 2023; Xie et al., 2023). However, all the above lines of research frame causality as empirical knowledge, thus relying heavily on the quality and the coverage of the training data, overlooking the great potential of the formal causal reasoning skills to process correlational information to causal conclusions.

Drawing inspirations from technical studies on causal discovery (Spirtes et al., 2001; Spirtes and Zhang, 2016; Glymour et al., 2019), we formulate a novel task for NLP, *correlation-to-causation inference* (CORR2CAUSE), which is an important skill for LLMs. Imagine the scenario in Figure 2.1, where the training corpus does not tediously cover every causal relation, but more pervasively talk about correlations, such as which events tend to co-occur. Learning a good CORR2CAUSE skill can enable LLMs to draw causal relations behind the mere correlational information on the surface. For example, several decades ago, there might be an observation that female university students tend to perform better, but behind the correlational statistics is the causal graph that female students have to achieve extra good performance to get into universities as the first place.

To this end, we collect the CORR2CAUSE dataset, the first dataset to test the pure causal reasoning abilities of LLMs. All the questions in this dataset are centered around testing when it is valid or invalid to infer causation from correlation. To systematically compose this dataset, we ground our generalization process in the formal framework of causal discovery (Spirtes et al., 2001; Glymour et al., 2016; Spirtes and Zhang, 2016), which provides rules about how to deduce causal relations among variables given their statistical correlation in the observational data. We generate more than 200K data points, and label a correlation-causation statement pair as valid if and only if there is a bijective mapping between the statistical correlation and the underlying causality.

Based on our CORR2CAUSE dataset with 200K samples, we investigate two main research questions: (1) How well do existing LLMs perform on this task? (2) Can existing LLMs be re-trained or re-purposed on this task and obtain robust causal inference skills? Through extensive experiments, we show empirically that none of the 17 existing LLMs we investigate perform well on this pure causal inference task. We also show that although LLMs can demonstrate better performance after being finetuned on the data, the causal inference skills attained by them are not robust. In summary, our contributions are as follows:

1. We propose the novel task of CORR2CAUSE, to probe an aspect of LLM’s reasoning ability, *pure causal inference*;
2. We compose a dataset of over 200K samples, using insights from causal discovery;
3. We evaluate the performance of 17 LLMs on our dataset, finding that all of them perform poorly, close to the random baseline;
4. We further explored whether LLMs can learn the skill through finetuning, and find that LLMs fail to robustly acquire this skill in out-of-distribution settings. Finally, we suggest future work to explore more ways to enhance the pure causal inference skill in LLMs.

2.2 Preliminaries: Causal Inference

2.2.1 Directed Graphical Causal Models (DGCMs)

A directed graphical causal model (DGCM) is a commonly used representation to express the causal relations among a set of variables. Given a set of N variables $\mathbf{X} = \{X_1, \dots, X_N\}$, we can encode the causal relations among them using a directed graph $\mathcal{G} := (\mathbf{X}, \mathbf{E})$, where \mathbf{E} is the set of directed edges. Each edge $e_{i,j} \in \mathbf{E}$ represents a causal link $X_i \rightarrow X_j$, meaning that X_i is a direct cause of X_j . In the context of this work, we take the common assumption of directed acyclic graphs (DAGs), which most causal discovery methods use (Glymour et al., 2019), as graphs with cycles can make the causal discovery process arbitrarily hard.

Following the graph-theoretic terminology, we use an analogy of the ancestry tree to denote the relations between two variables. For example, we call X_i as a *parent* of X_j if there is a directed edge $X_i \rightarrow X_j$ in the graph, and, thus, X_j is a *child* of X_i . Similarly, we denote X_i as an *ancestor* of X_j if there exists a directed path from X_i to X_j , and, thus, X_j is a *descendent* of X_i . Note that a parent is a special case of an ancestor where the directed path has a length of 1.

For convenience, we also introduce the notions for some special three-variable relations. Given two variables X_i and X_j , we call a third variable X_k a *confounder* (i.e., *common cause*) if X_k is a parent of both X_i and X_j ; a *collider* (i.e., *common effect*) if X_k is a child of both X_i and X_j ; and a *mediator* if X_k is both a child of X_i , and a parent of X_j .

2.2.2 D-Separation and Markov Property

D-Separation D-separation (Pearl, 1988) is a fundamental concept in graphical models used to determine whether two sets of nodes \mathbf{X} and \mathbf{Y} in a DAG \mathcal{G} are conditionally independent given a third set of nodes \mathbf{Z} , where the three sets are disjoint. We say that \mathbf{X} and \mathbf{Y} are d-separated by \mathbf{Z} if all paths between any node in \mathbf{X} and any node in \mathbf{Y} are *blocked* by the conditioning set \mathbf{Z} . A path between \mathbf{X} and \mathbf{Y} is blocked by \mathbf{Z} if there exists a node $A \in \mathbf{Z}$ which satisfies one of the following conditions: A is the parent node in a fork structure on the path (i.e., $\cdot \leftarrow A \rightarrow \cdot$); A is the mediator node in a chain structure on the path (i.e., $\cdot \rightarrow A \rightarrow \cdot$); or in any collider structure on the path (i.e., $\cdot \rightarrow A \leftarrow \cdot$), \mathbf{Z} does not contain A or its descendants.

Markov Property The Markov property in a DAG \mathcal{G} states that each node X_i is conditionally independent of its non-descendants given its parents, namely $X_i \perp\!\!\!\perp \mathbf{NonDe}(X_i) \mid \mathbf{PA}(X_i)$, where $\mathbf{NonDe}(X_i)$ denotes the non-descendants of X_i excluding itself, and $\mathbf{PA}(X_i)$ denotes

the parents of X_i . Using the Markov property, we can factorize the joint distribution of all the nodes in the graph into $P(X_1, \dots, X_N) = \prod_{i=1}^N P(X_i | \mathbf{PA}(X_i))$. To infer the causal graph from probability distributions, a common assumption is faithfulness, namely the validity to infer all the d-separation sets in the graph from the independence relations in the probability distribution. In our work, we also take this broadly taken assumption which holds for most real-world scenarios.

Markov Equivalence of Graphs We denote two DAGs as Markov equivalent if they induce the same joint distribution $P(\mathbf{X})$. The set of DAGs that are Markov equivalent to each other is called a Markov equivalence class (MEC). Causal graphs in the same MEC can be easily identified since they have the same skeleton (i.e., undirected edges) and V-structures (i.e., structures in the form of $A \rightarrow B \leftarrow C$ where A and C are not connected).

Obviously, there is a one-to-many mapping (i.e., surjection) between the causal graph and statistical distribution. Namely, each causal graph sufficiently determines a statistical distribution, but from a statistical distribution, we cannot necessarily induce a unique causal graph. This is why we say “correlation does not necessarily mean causation”.

2.2.3 Causal Discovery

Causal discovery aims to learn the causal relations by analyzing statistical properties in the observational data (Spirtes et al., 2001; Glymour et al., 2016; Spirtes and Zhang, 2016; Glymour et al., 2019). It can be achieved through constraint-based methods (Spirtes et al., 2001), score-based methods (Chickering, 2002), or other methods taking advantage of the functional causal models (Shimizu et al., 2006; Hoyer et al., 2008; Zhang and Hyvärinen, 2009).

To fit for the spirit of this paper to infer from correlation (expressed in natural language) to causation, we base our dataset design on the widely-used Peter-Clark (PC) algorithm (Spirtes et al., 2001). The PC algorithm is based on the principles of conditional independence and the causal Markov assumption, which allows it to efficiently identify causal relationships among variables in a given dataset. The algorithm first starts with a fully connected undirected graph among all the variables. Then it removes the edge between two variables if there is an unconditional or conditional independence relationship between them. Afterwards, it orients the directed edges whenever there is a V-structure. And finally, it iteratively checks the direction of the other edges until the entire causal graph is consistent with all the statistical correlations.

2.3 Dataset Construction

We introduce the construction of our dataset in this section. We start with our task formulation for CORR2CAUSE, and then briefly give an overview of the data generation process, followed by detailed descriptions of each step. We conclude the section with the overall statistics of the dataset.

2.3.1 Task Formulation

Given a set of N variables $\mathbf{X} = \{X_1, \dots, X_N\}$, we have a statement s about all the correlations among the variables, and a hypothesis h describing the causal relation r between

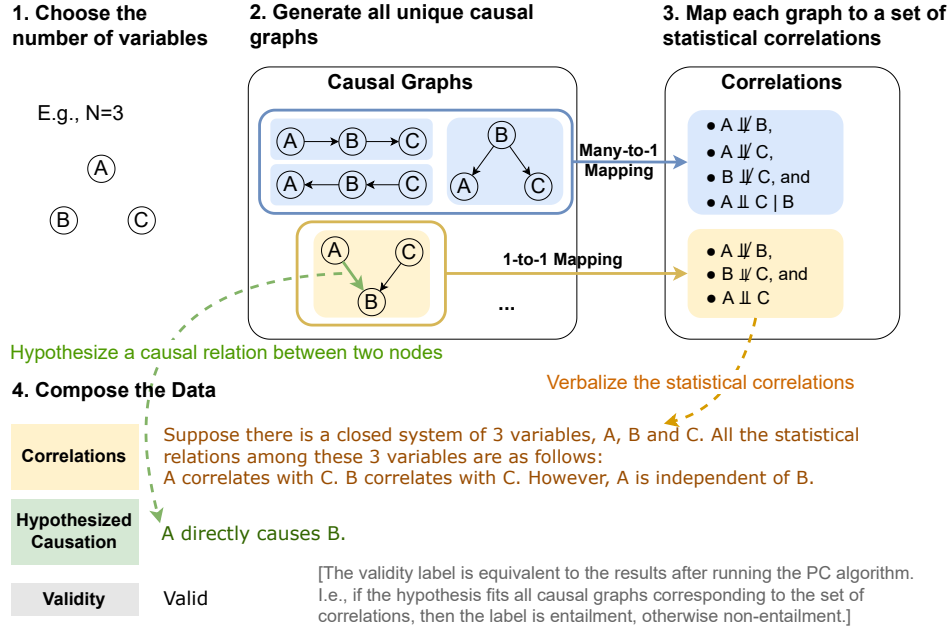


Figure 2.2: Pipeline of the data construction process.

the pair of variables X_i and X_j . The task is to learn a function $f : (s, h) \mapsto v$ which maps the correlation statement s and the causal relation hypothesis h to their validity $v \in \{0, 1\}$, which takes the value 0 if this inference is invalid, and the value 1 if this inference is valid.

2.3.2 Overview of the Data Generation Process

We base the construction our dataset on several concepts of causal inference, including the DGCM, d-separation, and MECs, as introduced in Section 2.2.

As in the overview of our data generation process in Figure 2.2, we first choose the number N of variables (Step 1) and generate all the unique DGCMs with N nodes (Step 2), which we will introduce in the Section 2.3.3. Then we collect all the d-separation sets from these graphs to identify MECs (Step 3) in Section 2.3.4. Then, in Step 4, we create the formal form of data in Section 2.3.5. For each correspondence of the MEC to causal graphs, we compose the correlation statement based on the statistical relations in the MEC, and hypothesize a causal relation between two variables, and produce the validity $v = 1$ if the hypothesis is a shared property of all causal graphs in the MEC, and $v = 0$ if the hypothesis is not necessarily true for all the MEC graphs. Finally, we introduce the verbalization process in Section 2.3.6.

2.3.3 Constructing the Graphs with Isomorphism Checks

The first step of the data generation is to compose the causal graphs, as in Step 1 and 2 of Figure 2.2. For a set of N variables $X = \{X_1, \dots, X_N\}$, there are $N(N - 1)$ possible directed edges, since each node can link to any node other than itself. To remove cycles in the graph, we make the nodes in topological order, which only allows edges $X_i \rightarrow X_j$, where $i < j$. We achieve this by limiting the adjacency matrix of the graph to only having non-zero values above the diagonal, resulting in $N(N - 1)/2$ possible directed edges for the DAGs.

# Nodes	# Unique DAGs	# Edges/DAG	# MECs	# DAGs/MEC
2	2 out of 2	0.50	2	1.0
3	6 out of 2^3	1.67	5	1.2
4	31 out of 2^6	3.48	20	1.55
5	302 out of 2^{10}	5.89	142	2.13
6	5,984 out of 2^{15}	8.77	2,207	2.71
Total	6,325	8.60	2,376	2.66

Table 2.1: Statistics about the source causal graphs in our dataset. Given the number of nodes, we report the number of unique DAGs, average number of edges per DAG, number of MECs, and average number of DAGs per MEC.

At the first glance, for N nodes, there should be $2^{N(N-1)/2}$ possible DAGs (i.e., the power set of all edges). However, there could be isomorphic graphs in this set. To avoid this, we perform a graph isomorphism check (McKay and Piperno, 2014), and reduce the set so that only unique DAGs are retained, and we show their statistics in Table 2.1. Although we can handle large graphs, we mostly focus on smaller graphs that can still lead to a reasonably sized dataset, so we empirically set $N = 6$, but future work can use our open-sourced codes to extend to more nodes.

2.3.4 Programmatically Generating the D-Separation Sets

Based on the set of unique DAGs, we then programmatically generate the d-separation sets by graph theoretical conditions, as in Step 3 of Figure 2.2. To realize this step, we code an efficient graph-theoretic algorithm to check for all the chain, fork, and collider structures to automatically identify the set of nodes that d-separate each pair of nodes. Using the d-separation sets and the faithfulness assumption, we form the statistical correlations as follows. For each pair of nodes, they are conditionally independent given the variables in the d-separation set. If the d-separation set is empty, then the two nodes are unconditionally independent. If no d-separation set can be found for the two nodes, then they are directly correlated.

Moreover, using the d-separation sets, we are able to cluster causal graphs to MECs. We achieve it by tracing the mapping between the causal graphs and the set of statistical correlations, and backtracking the graphs with the same d-separation sets to group them in the same MEC. We show in Table 2.1 that each MEC contains on average 2.66 DAGs.

2.3.5 Composing the Hypotheses and Label

After generating the set of correlations based on the d-separation sets, we now generate the causal hypotheses. For the causal relation r , we focus on six common causal relations between two nodes introduced in Section 2.2.1: Is-Parent, Is-Child, Is-Ancestor (excluding the parents), Is-Descendant (excluding the children), Has-Confounder (i.e., there exists a confounder, or common cause, of the two nodes), and Has-Collider (i.e., there exists a collider, or common effect, of the two nodes). In this way, the set of hypotheses contains all six meaningful causal relations between every pair of variables, resulting in a total size of $6 \cdot N(N-1)/2 = 3N(N-1)$ hypotheses for a graph with N variables.

To generate the ground-truth validity label, we start from the correlation sets in Step 3, then look up all the causal graphs in the same MEC corresponding to the given set of

Causal Relation	Hypothesis Template
Is-Parent	{Var i} directly causes {Var j}.
Is-Ancestor	{Var i} causes something else which causes {Var j}.
Is-Child	{Var j} directly causes {Var i}.
Is-Descendant	{Var j} is a cause for {Var i}, but not a direct one.
Has-Collider	There exists at least one collider (i.e., common effect) of {Var i} and {Var j}.
Has-Confounder	There exists at least one confounder (i.e., common cause) of {Var i} and {Var j}.

Table 2.2: Templates for each causal relation in the hypothesis. We use {Var i} and {Var j} as placeholders for the two variables.

correlations, and check the necessity of the hypothesized causal relation. If the causal relationship proposed in the hypothesis is valid for all causal graphs within the MEC, then we generate the validity $v = 1$; otherwise, we generate $v = 0$. A special case of valid samples is that when the size of the MEC is 1, then there is a bijective mapping between the causal graph and the d-separation sets, so any hypothesis stating the causal properties of that unique causal graph is valid.

2.3.6 Verbalizing into Language

Finally, as in the last step of Figure 2.2, we convert all the information above to text data for our CORR2CAUSE task. For the correlation statement, we verbalize the set of correlations in Step 3 into a natural language statement s . When two variables cannot be d-separated, i.e., $A \not\perp\!\!\!\perp B$, then we describe them as “A correlates with B” since they are directly correlated and cannot be independent by any condition. And if two variables have a valid d-separation set C , then we describe them as “A is independent of B given C.” In the special case when the d-separation set is empty, we directly say “A is independent of B.” In addition, we disambiguate the setting by starting the correlation statement with the setup of a closed system of the given variables, and no hidden variables: “Suppose there is a closed system of N variables, A, B, ... All the statistical relations among these N variables are as follows:”. Finally, to verbalize the hypothesis, we feed the causal relation triplet (X_i, r, X_j) into their hypothesis templates in Table 2.2. For example, we turn the triplet $(A, \text{Is-Parent}, B)$ into “A directly causes B”, as in the example of Figure 2.2.

2.3.7 Statistics of the Resulting Data

We show the statistics of our CORR2CAUSE dataset in Table 2.3. Overall, our dataset contains 207,972 samples, where 18.57% of the samples have the positive label (i.e., with validity=1). The average length of the premise is 424.11 tokens, and hypothesis 10.83 tokens. We split the data into 205,734 training samples, 1,076 development and 1,162 test samples.¹ Since the main purpose of the dataset is to benchmark the performance of LLMs, we prioritize the test and development sets to have a comprehensive coverage over all sizes of graphs. Specifically, we iterate through the subset of our data for each N, and split it entirely for only the test and development sets if the data is less than 1K, which is

¹Note for our dataset v2.0: We notice that our original data (v1.0) has duplication due to symmetric relations and verbalizations of the hypothesis. E.g., Is-Parent(A, B) has the exact hypothesis verbalization as Is-Child(B, A). Hence, for our data v2.0, we perform a careful de-duplication, and update the data statistics in Table 2.3. See more version comparison details in Appendix A.1.4. Note that, due to the symmetry, the current version is a random sample half of the size of the original version, so the modeling results in the experiment section roughly hold.

the case for $N = 2$ and 3. For the other subsets that are larger, we randomly sample up to 1K or 10% of the data, whichever is smaller, to the test and development sets. We set the cap to be 1K in order to form a reasonable computation budget, since many LLMs are expensive to query in the inference mode. Aside from the test and valid sets, all the rest of the data goes into the training set.

	Overall	Statistics by the Number of Nodes N				
		$N = 2$	$N = 3$	$N = 4$	$N = 5$	$N = 6$
# Samples	207,972	12	90	720	8,520	198,630
# Test	1,162	6	48	72	514	522
# Dev	1,076	6	42	72	482	474
# Train	205,734	0	0	576	7,524	197,634
# Tokens/Premise	424.11	31.5	52.0	104.0	212.61	434.54
# Tokens/Hypothesis	10.83	10.83	10.83	10.83	10.83	10.83
% Positive Labels	18.57	0.00	3.33	7.50	13.01	18.85
Vocab Size	65	49	53	55	57	61

Table 2.3: Statistics of our CORR2CAUSE dataset, and by subsets. We report the total number of samples (# Samples); splits of the test (# Test), development (# Dev) and training sets (# Train); number of tokens per premise (# Tokens/Premise) and hypothesis (# Tokens/Hypothesis); percentage of the positive labels (% Positive Labels), and vocabulary size by the number of unique tokens (Vocab Size). Note that the number of unique graphs and MECs are in Table 2.1.

2.4 Experiments

2.4.1 Experimental Setup

We set up a diverse list of LLMs for the experiments on our CORR2CAUSE dataset. To *test existing LLMs*, we first include six commonly used BERT-based NLI models in the transformers library (Wolf et al., 2020): BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), BART (Lewis et al., 2020), DeBERTa (He et al., 2021), DistilBERT (Sanh et al., 2019), and DistilBART (Shleifer and Rush, 2020). Apart from these BERT-based NLI models, we also evaluate the general-purpose autoregressive LLMs based on GPT (Radford et al., 2019): GPT-3 Ada, Babbage, Curie, Davinci (Brown et al., 2020); its instruction-tuned versions (Ouyang et al., 2022), text-davinci-001, text-davinci-002, and text-davinci-003; and GPT-3.5 (i.e., ChatGPT), and the latest GPT-4 (OpenAI, 2023) by April 2023, using the OpenAI API (<https://openai.com/api/>) with temperature 0. We also evaluate the recent, more efficient models, LLaMa (Touvron et al., 2023) and Alpaca (Taori et al., 2023a).

When inspecting the behavior of *finetuned models*, we adopt a large set of models, including GPT-based models (GPT-3 Ada, Babbage, Curie, and Davinci) using the OpenAI fine-tuning API for classification at <https://platform.openai.com/docs/guides/fine-tuning>, open-sourced decoder-only models (GPT2, GPT2-Large, GPT2-XL, LLaMA-7B, and LLaMA2-7B), BERT-based models from scratch (BERT-Base, BERT-Large, RoBERTa-Base, and RoBERTa-Large), and BERT-Based NLI models (BERT-Base MNLI, BERT-Large MNLI, RoBERTa-Base MNLI, and RoBERTa-Large MNLI) using the transformers library (Wolf et al., 2020). See training details in Appendix A.1.1.

For the *random baselines*, we provide “always majority” to predict the majority class 100%

	F1	Precision	Recall	Accuracy
Random Baselines				
Always Majority	0.0	0.0	0.0	84.77
Random (Proportional)	13.5	12.53	14.62	71.46
Random (Uniform)	<u>20.38</u>	15.11	31.29	62.78
BERT-Based Models				
BERT MNLI	2.82	7.23	1.75	81.61
RoBERTa MNLI	22.79	34.73	16.96	82.50
DeBERTa MNLI	14.52	14.71	14.33	74.31
DistilBERT MNLI	20.70	24.12	18.13	78.85
DistilBART MNLI	26.74	15.92	83.63	30.23
BART MNLI	33.38	31.59	35.38	78.50
LLaMa-Based Models				
LLaMa-7B	26.81	15.50	99.42	17.36
Alpaca-7B	<u>27.37</u>	15.93	97.37	21.33
GPT-Based Models				
GPT-3 Ada	0.00	0.00	0.00	84.77
GPT-3 Babbage	27.45	15.96	97.95	21.15
GPT-3 Curie	26.43	15.23	100.00	15.23
GPT-3 Davinci	27.82	16.57	86.55	31.61
GPT-3 Instruct (text-davinci-001)	17.99	11.84	37.43	48.04
GPT-3 Instruct (text-davinci-002)	21.87	13.46	58.19	36.69
GPT-3 Instruct (text-davinci-003)	15.72	13.4	19.01	68.97
GPT-3.5	21.69	17.79	27.78	69.46
GPT-4	<u>29.08</u>	20.92	47.66	64.60

Table 2.4: Overall performance. We report F1 (main metric), precision, recall and accuracy. For the main metric, F1 score, we use the **bold** font to highlight the overall best performance, and underline to highlight the best performance within each category of models.

of the time, “random (uniform)” to uniformly sample a label (i.e., 50% for each), and “random (proportional)” to sample a label from a Bernouli distribution proportional to the development set label distribution.

2.4.2 The Corr2Cause Skill in Existing LLMs

We show the performance of seventeen LLMs in Table 2.4. We can see that pure causal inference is a very challenging task across all existing LLMs. Among all the LLMs, the best performance is 33.38% F1 by BART MNLI, which is even higher than the latest GPT-based model, GPT-4. Notably, many models are worse than random guess, which means that they totally fail at this pure causal inference task. The observation still holds for few-shot chain-of-thought prompts tested in Appendix A.1.7.

2.4.3 Finetuned Performance

Next, we address the question: *Can we re-purpose LLMs to learn this task?* The experimental results in Table 2.5a of 17 models finetuned on our CORR2CAUSE seem very strong at first sight. Most models see a substantial increase, among which the finetuned BERT-based NLI models demonstrate the strongest performance. The best-performing one, RoBERTa-Large MNLI, achieves 94.74% F1 score on this task, as well as very high precision, recall and accuracy scores.

	F1	Precision	Recall	Accuracy	F1 (Paraph.)	F1 (Var. Ref.)
<i>Finetuned GPT-Based Models Using OpenAI API</i>						
GPT-3 Ada	79.85	70.47	92.11	92.92	61.73	41.57
GPT-3 Babbage	78.19	69.98	88.60	92.48	62.34	43.28
GPT-3 Curie	81.23	75.00	88.60	93.77	64.93	45.32
GPT-3 Davinci	<u>85.52</u>	80.26	91.52	95.28	65.01	46.96
<i>Finetuned Open-Sourced Decoder-Only Models</i>						
GPT2	89.18	88.03	90.35	96.66	56.76	31.70
GPT2-Large	94.29	92.18	96.49	98.22	55.95	31.99
GPT2-XL	<u>94.30</u>	91.94	96.78	98.22	60.32	43.95
LLaMA-7B	91.98	88.62	95.61	97.46	56.41	53.92
LLaMA2-7B	92.92	90.11	95.91	97.77	52.24	49.47
<i>Finetuned BERT-Based Models</i>						
BERT-Base	69.29	54.42	95.32	87.13	61.13	35.20
BERT-Large	85.26	77.51	94.74	95.01	63.64	38.54
RoBERTa-Base	87.60	78.47	99.12	95.73	65.58	53.12
RoBERTa-Large	<u>89.10</u>	82.54	96.78	96.39	65.05	60.20
<i>Finetuned BERT-Based NLI Models</i>						
BERT-Base MNLI	89.88	85.49	94.74	86.51	65.56	31.50
BERT-Large MNLI	90.19	84.44	96.78	96.79	67.24	52.04
RoBERTa-Base MNLI	94.27	90.35	98.54	98.17	57.42	62.83
RoBERTa-Large MNLI	<u>94.74</u>	92.24	97.37	98.35	55.45	67.87

(a) Performance of finetuned models on the original test set.

(b) F1 scores of finetuned models on the perturbed test sets by paraphrasing (Paraph.) and variable refactorization (Var. Ref.).

Table 2.5: Performance of finetuned models on the original test set and perturbed test sets.

2.4.4 Fine-Grained Performance by Causal Relation

In addition to the overall results mentioned above, we conduct a fine-grained analysis to check the performance of the strongest finetuned model, RoBERTa-Large MNLI, by our six causal relation types. As in Table 2.6a, the model is very good at judging relations such as Is-Parent, Is-Descendant and Has-Confounder, all with more than 96% F1 scores, whereas it is several points weaker on the Has-Collider relations. This could be due to that the collider relation is the most special type, requiring identification of the V-structure based on both the unconditional independence based on the two variables only and correlations whenever conditioned on a common descendant. We also conduct error analysis for non-finetuned models in Appendix A.1.6.

2.4.5 Robustness Analysis

Looking at the very high performance of the finetuned models, we raise the next question: *Did the models really robustly learn the causal inference skills?*

Two Robustness Tests We design two simple robustness tests: (1) paraphrasing, and (2) variable refactorization. For (1) paraphrasing, we simply paraphrase the hypothesis by changing the text template for each causal relation to some semantically-equivalent alternatives in Appendix A.1.3. For (2) variable refactorization, we reverse the alphabet

Relation Type	F1	Precision	Recall	Accuracy	F1	Precision	Recall	Accuracy
Is-Parent	96.18	95.45	96.92	98.67	74.80	79.31	70.77	91.73
Is-Ancestor	93.94	93.94	93.94	98.93	45.45	90.91	30.30	93.60
Is-Child	95.73	94.92	96.56	98.67	73.39	78.43	68.97	92.27
Is-Descendant	96.55	93.33	100	99.47	29.41	83.33	17.86	93.60
Has-Collider	92.19	87.41	97.52	94.64	70.70	75.00	66.90	82.04
Has-Confounder	98.67	97.37	100	99.73	70.42	73.53	67.57	94.37

(a) Fine-grained performance of RoBERTa-Large by causal relation type on the original test set.

(b) Its fine-grained performance by relation type after variable refactorization.

Table 2.6: Fine-grained analysis of the best-performing model, RoBERTa-Large MNLI.

of the variable names, namely flipping A, B, C, to Z, Y, X and so on. The inspiration behind the two robustness tests comes from the spurious correlation analysis described in Appendix A.1.5.

Specifically, we adopt the common setup of text adversarial attack (Morris et al., 2020; Jin et al., 2020) to preserve the training set and keep the same saved models, but run the inference on the perturbed test set. In this way, we separate the possibilities of the models only overfitting on the training data vs. mastering the reasoning skills.

Results after Perturbation We can see from Table 2.5b that all the models drop drastically, by up to 39.29 on the paraphrased test set, and up to 62.30 after variable refactorization. The best-performing model, RoBERTa-Large MNLI, is especially sensitive towards paraphrasing, demonstrating the most drop among all models; however, it is the most robust against the variable refactorization, maintaining a high F1 score of 67.87. We conduct fine-grained analysis for RoBERTa-Large MNLI under perturbation in Table 2.6b. We can see the the main source of the performance drop of the model comes from the two classes, Is-Ancestor (decreasing to 45.45%) and Is-Descendant (decreasing to 29.41%), while the other classes stay relatively robust, keeping their F1 scores over 70%.

From this analysis, we make the following suggestions to future studies testing this CORR2CAUSE skill of LLMs. First, it is safe to use it as a test set to benchmark existing LLMs’ performance, since the data we generate is out-of-distribution from the training data of the current LLMs. Then, when testing finetuned models, it is very important to accompany adversarial attack together with the i.i.d. test set. We open-source our perturbed test sets for future work to test the generalizability skill.

2.4.6 Extension to Natural Stories

We envision our CORR2CAUSE dataset to be a foundation for future extensions to various settings, such as instantiating the variables with actual phenomena and situating the story in a more natural setting. For example, the *correlation does not imply causation* rule can be instantiated with the ice cream sales and swimming pool attendance as the two variables, and argue that ice cream sales does not necessarily affect swimming pool attendance, because their correlation could be due to a third variable, such as hot weather. We provide a case study for how to instantiate the symbolic expressions in our dataset to more natural stories, and find that LLMs such as GPT-4 can generate realistic, daily life stories that has

foreseeably broad applications. See more details in Appendix A.1.2.

2.5 Related Work

Existing Causal Reasoning Tasks A large body of existing research of causal reasoning in NLP focuses on leveraging empirical knowledge to do tasks such as inferring the cause and effect of why an agent perform certain tasks (Sap et al., 2019a), the motivation and emotional reaction in a social context (Sap et al., 2019b), how people achieve a given goal with a set of concrete steps (Zhang et al., 2020a), the development of a story given a different beginning (Qin et al., 2019), and how in general LLMs serve as a knowledge base of cause and effect (Willig et al., 2023; Kıcıman et al., 2023). In contrast, our CORR2CAUSE task focuses on the pure causal inference skill of models, which is a knowledge-dependent reasoning skill based on formally correct rules from causal inference.

Existing Logical and Inference Tasks Another related area of literature is logical and inference tasks, of which a well-established one is natural language inference (NLI), to identify the semantic relationship between a pair of sentences (MacCartney and Manning, 2008; Bowman et al., 2015). NLI datasets mainly focus on the set and paraphrase relations. For example, “a group of boys are playing football” can entail “some guys are playing football,” where “boys” are a sub-concept of “guys,” and “a group of” and “some” are paraphrases. Recently, there have been increasing efforts to extend the inference task to various logical inference skills such as deductive logic and propaganda techniques (Jin et al., 2022b; Alhindi et al., 2022). Our CORR2CAUSE dataset is the first dataset testing the correlation-to-causation inference skill, which is unique of its type.

2.6 Conclusion

In this work, we introduced a novel task, CORR2CAUSE, to infer causation from correlation, and collected a large-scale dataset of over 200K samples. We evaluated an extensive list of LLMs on this new task, and showed that off-the-shelf LLMs perform poorly on this task. We also show that it is possible to re-purpose LLMs on this task by finetuning, but future work needs to be aware of the out-of-distribution generalization problem. To avoid the Goodhart’s law, we recommend using this dataset to benchmark the pure causal inference skills for LLMs that have not seen this dataset. Given the limited reasoning abilities of current LLMs, and the difficulty of separating actual reasoning from training-corpus-derived knowledge, it is imperative that our community focus on work aiming to accurately disentangle and measure both abilities. We believe the present work is a first such step.

Limitations and Future Work

We identify several limitations of this work and open future directions: First, in the context of this work, we limit the causal graphs to two to six nodes, but future work can feel free to explore larger graphs. Another aspect is that we do not assume hidden confounders in this inference problem, so we welcome future work to generate an even more challenging dataset to infer the existence of hidden confounders, analogous to the causal

discovery algorithm of fast causal inference (FCI) (Spirtes et al., 2001). And also in general, explorations of other causal discovery algorithms are welcomed too. Finally, a lot of motivation behind proposing this task is inspired by the problem of invalid reasoning patterns in our daily reasoning (Jin et al., 2022b), which could fertilize the ground for more pervasive spread of fake news. We believe false causal inference is a prevalent type of fallacious beliefs, and welcome future work to connect the idea of this benchmark to more real-world false beliefs based on confusing correlation with causation.

CLadder: Assessing Causal Reasoning in LLMs

Apart from causal discovery covered in the previous chapter, we further explore causal effect reasoning in natural language, inspired by the “*causal inference engine*” postulated by Pearl and Mackenzie (2018). To this end, we compose a large dataset, CLADDER, with 10K samples: based on a collection of causal graphs and queries (associational, interventional, and counterfactual), we obtain symbolic questions and ground-truth answers, through an oracle causal inference engine. These are then translated into natural language. We evaluate multiple LLMs on our dataset, and we introduce and evaluate a bespoke chain-of-thought prompting strategy, CAUSALCoT. We show that our task is highly challenging for LLMs, and we conduct an in-depth analysis to gain deeper insights into the causal reasoning abilities of LLMs. Our data is open-sourced at <https://huggingface.co/datasets/causalNLP/cladder>, and our code can be found at <https://github.com/causalNLP/cladder>.

3.1 Introduction

Once we really understand the logic behind causal thinking, we could emulate it on modern computers and create an “artificial scientist”.

— Pearl and Mackenzie (2018)

Causal reasoning is believed to be one of the hallmarks of human intelligence (Penn and Povinelli, 2007; Harari, 2014). The ability to draw causal inferences from available information is crucial for scientific understanding and rational decision-making: for example, knowing whether smoking causes cancer might enable consumers to make a more informed decision (Doll and Hill, 1950, 1954); assessing the causal effect of a vaccine is essential for effective policy-making during a pandemic (Plotkin, 2005; De Serres et al., 2013; Whitney et al., 2014; Kekić et al., 2023); and understanding the interplay behind family background, education and income helps devise effective education policies (Card, 1999; Chetty et al., 2011; Heckman et al., 2006; Psacharopoulos and Patrinos, 2004).

Our opening quote therefore mirrors the aspirations of many scientists in artificial intelligence and causal inference: to construct a machine capable of performing sound causal reasoning, and able to answer causal questions at scale and with ease. Recent advances in large language models (LLMs) have brought about a paradigm shift in natural language processing (NLP) and artificial intelligence (Radford et al., 2019; Devlin et al., 2019; Brown et al., 2020; Zhang et al., 2022; OpenAI, 2023; Ignat et al., 2024, *inter alia*). These transformative developments raise the question of whether these machines are already capable of causal reasoning: *Do LLMs understand causality?*

Many previous works addressed the above question by focusing on *commonsense* causality (Zečević et al., 2023; Zhang et al., 2023; Ho et al., 2022), inspired by the literature that

Question: Imagine a self-contained, hypothetical world with only the following conditions, and without any unmentioned factors or causal relationships:

Physical vulnerability has a direct effect on the likelihood of **fatality** and **vaccination decision**. **Vaccination** has a direct effect on the **fatality rate**.

In the entire population, 50% of the people are vulnerable to a certain disease.
 For vulnerable and vaccinated people, the fatality rate is 4%. For vulnerable and unvaccinated people, the fatality rate is 7%.
 For strong and vaccinated people, the fatality rate is 1%. For strong and unvaccinated people, the fatality rate is 5.8%.
 Overall, the fatality rate for vaccinated people is 5%, while the fatality rate for unvaccinated people is 4.5%.

Does getting vaccinated increase the likelihood of death?

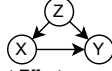
Ground-Truth Answer: No

CLadder

Correct steps to lead to the ground-truth answer:

1) Parse the **causal graph**: Confounding

Subskill: Causal Relation Extraction



2) Classify the **query type**: Average Treatment Effect

Subskill: Causal Question Classification

3) Formulate the query to its **symbolic form**:

$E[Y | do(X=1)] - E[Y | do(X=0)]$

Subskill: Formalization

4) Collect the **available data**:

$P(Z=1)=0.5$

$P(Y=1|Z=1,X=1)=0.04$, $P(Y=1|Z=1,X=0)=0.07$

$P(Y=1|Z=0,X=1)=0.01$, $P(Y=1|Z=0,X=0)=0.058$

$P(Y=1|X=1)=0.05$, $P(Y=1|X=0)=0.045$

Subskill: Semantic Parsing

5) Derive the **estimand** using **causal inference**:

$E[Y | do(X=1)] - E[Y | do(X=0)]$

Subskill: Formal Causal Inference

$= \sum_{Z=v} P(Z=v) [P(Y=1|Z=v,X=1) - P(Y=1|Z=v,X=0)]$ # remove "do" using do-calculus
 $= P(Z=0) [P(Y=1|Z=0,X=1) - P(Y=1|Z=0,X=0)]$
 $+ P(Z=1) [P(Y=1|Z=1,X=1) - P(Y=1|Z=1,X=0)]$ # turn the expression into terms in the available data

6) Solve for the estimand by plugging in the relevant data in Step 4:

$= 0.5 \cdot (0.01 - 0.058) + 0.5 \cdot (0.04 - 0.07)$ # plug in the numbers in the available data

$= -0.039$

Subskill: Arithmetics

< 0 # the effect size is negative, so the final answer is "No"

Figure 3.1: Example question in our CLADDER dataset featuring an instance of *Simpson’s paradox* (Pearl, 2022). We generate the following (symbolic) triple: (i) the causal query; (ii) the ground-truth answer, derived through a *causal inference engine* (Pearl and Mackenzie, 2018); and (iii) a step-by-step explanation. We then *verbalize* these questions by turning them into stories, inspired by examples from the causality literature, which can be expressed in natural language.

explores LLMs as *knowledge bases* (Petroni et al., 2019; Shin et al., 2020; Jiang et al., 2020) (we refer to this line of work as *causality as knowledge*). This involves assessing the alignment between commonsense knowledge about causal relationships in humans and LLMs. This line of work generally does not focus on evaluating how well models are capable of *causal reasoning*. For example, it may be difficult to rule out the possibility that LLMs perform potentially unreliable *amortized causal inference*, answering causal questions by a simple repetition of verbal patterns present in the texts composing their training data:¹² in other words, LLMs may just be “*causal parrots*” (Zečević et al., 2023).

In this work, we introduce a way to test the *formal causal reasoning in LLMs*. To this end, we introduce the CLADDER dataset. The specificity of CLADDER is that causal questions posed in natural language are *grounded in symbolic questions and ground truth answers*: the latter are derived through an oracle *causal inference engine* (CI engine) (Pearl and Mackenzie, 2018), which abides by the rules of the causal inference approach described by Pearl (2009b), based on graphical models and structural causal models (SCMs) (Pearl, 1995; Spirtes et al., 2001; Pearl, 2009b; Glymour et al., 2016; Peters et al., 2017). We compose more than 10,000 causal questions that cover a variety of causal queries across the three rungs of the *Ladder of Causation* (Pearl and Mackenzie, 2018; Bareinboim et al., 2022)—i.e., *associational* (Rung 1), *interventional* (Rung 2), and *counterfactual* (Rung 3). We consider

¹which may itself contain instances of fallacious causal reasoning.

²The extent to which this would imply an inaptitude of LLMs for causal reasoning has been questioned (Huszár, 2023).

several causal graphs, giving rise to scenarios which require different causal inference abilities. Additionally, we generate ground-truth explanations with step-by-step reasoning for more in-depth analysis of LLM behavior. Our symbolic questions and answers are then *verbalized*, by turning them into stories which can be expressed in natural language. To probe whether LLMs employ amortized causal inference, we construct stories with commonsensical, as well as anti-commonsensical and with nonsensical causal relations: in these latter cases, amortized causal inference is expected to fail, whereas formal causal reasoning would still yield the correct answer. An example question from CLADDER is shown in Figure 3.1.

Exploiting CLADDER, we also introduce a method to elicit sound causal reasoning in LLMs and help them solve challenging causality questions. Specifically, we develop **Causal-CoT**, a chain-of-thought prompting strategy (Wei et al., 2022b) inspired by the CI engine, which prompts the LLM to extract the causal graph, causal query, and available “data” (e.g., conditional or interventional *do*-probabilities (Goldschmidt and Pearl, 1992)) from the question, formalize them precisely, and perform correct causal inferences. Our experiments indicate that CAUSALCoT achieves an accuracy of 70.40%, which substantially improves the performance of vanilla GPT-4 by 8.37 points on CLADDER.

We summarize the *main contributions* of our work:

1. In contrast to most other works on causality in LLMs, focusing on *commonsense causal knowledge*, our goal is to assess the LLMs’ ability to perform *formal causal reasoning* (briefly reviewed in Section 3.2).
2. We introduce CLADDER (Section 3.3), a dataset containing more than 10K causal questions, spanning all three rungs of the ladder of causation, several causal graphs, and various stories for verbalization.
3. We develop CAUSALCoT (Section 3.4), a chain-of-thought prompting strategy to elicit formal causal reasoning in LLMs, inspired by the *causal inference engine*.
4. We perform extensive experiments on eight LLMs (Section 3.5), analyze fine-grained errors to showcase the limitations of LLMs in formal causal reasoning, and suggest directions for future research.

3.2 Preliminaries on Causal Inference

Our dataset design takes inspiration from the *Causal Inference Engine* as postulated by Pearl and Mackenzie (2018), see also (Pearl, 1995). We begin with a brief overview of the causality framework by Pearl (2009b).³ This framework was largely developed within the field of artificial intelligence, and therefore puts particular emphasis on *algorithmic* aspects of causal reasoning (e.g., (Pearl, 2011))—which makes it particularly suited for our work, where we want to algorithmically generate ground truth answers to causal queries, without having to appeal to common sense to assess the correctness of an answer.

³We refer to (Pearl et al., 2016; Bareinboim et al., 2022) for a comprehensive introduction. See also Appendix A.2.3 for further details.

3.2.1 The Ladder of Causation

The *Ladder of Causation*, introduced by Pearl and Mackenzie (2018), is a proposed taxonomy, and hierarchy, of causal inference tasks (Bareinboim et al., 2022). It consists of three distinct rungs.

Rung 1 (“seeing”). This describes statistical associations (“How often do I take an aspirin when I have a headache?”). Rung 1 deals with statistical dependences among random variables, and involves probabilistic reasoning about joint and conditional distributions, $P(X = x, Y = y)$ and $P(Y = y|X = x)$, which can be formalised through *Bayesian Networks* (Pearl, 1988; Cowell et al., 2007) representing a set of variables and their conditional dependencies via a directed acyclic graph (DAG).

Rung 2 (“doing”). This enables us to formalize the concept of actively intervening in the world, and modifying it toward some end (“If I take an aspirin now, will my headache subside?”). Interventions can be formalized using the *do-operator* (Goldschmidt and Pearl, 1992) and *Causal Bayesian Networks* (Pearl, 2009b) to represent, for example, the distribution over Y when intervening on X to set its value to x as $P(Y = y|\text{do}(X = x))$.

Rung 3 (“imagining”). This rung deals with counterfactual reasoning, i.e., reasoning about alternative scenarios in which the world could have been different, possibly even contradicting the factual state (“Would my headache have subsided, if I had taken an aspirin?”). Counterfactual probabilities can be written as $P(Y_x = y)$, representing the probability that “ Y would be y , had X been x ”. Reasoning about Rung 3 quantities requires the introduction of *Structural Causal Models* (SCMs) (Pearl, 2009b). SCMs are especially powerful as they enable any quantity in Rungs 1, 2, and 3 to be formulated precisely (Bareinboim et al., 2022).

3.2.2 Causal Inference

Identification. Causal inference is especially difficult since we typically only have measurements from *lower* rungs, but want to reason about *higher* ones. A crucial question is then under what conditions are such inferences possible, i.e., what assumptions and measurements are required to unambiguously answer a causal query of interest: this is the question of *identification*. As argued in (Bareinboim et al., 2022), “it is generically impossible to draw higher-layer inferences using only lower-layer information”. One may be able to draw inferences at a higher layer given a combination of partial knowledge of the underlying SCM, in the form of a causal graph, and data at lower layers. The graphical structure therefore plays a crucial role in bridging the rungs of the Ladder of Causation, and many prior works have been dedicated to exploiting properties of the graph to transform higher-rung queries into expressions which can be estimated based on lower-rung quantities (Huang and Valtorta, 2006; Shpitser and Pearl, 2006a; Pearl and Bareinboim, 2022).

Causal Inference Engine. An overarching objective of this research is the construction of a *Causal Inference Engine* (CI Engine) (Pearl, 1995; Pearl and Mackenzie, 2018; Hünemann and Bareinboim, 2019), which takes as input a query, a graph, and some available data (typically from lower rungs than the query); and outputs whether a solution exists,

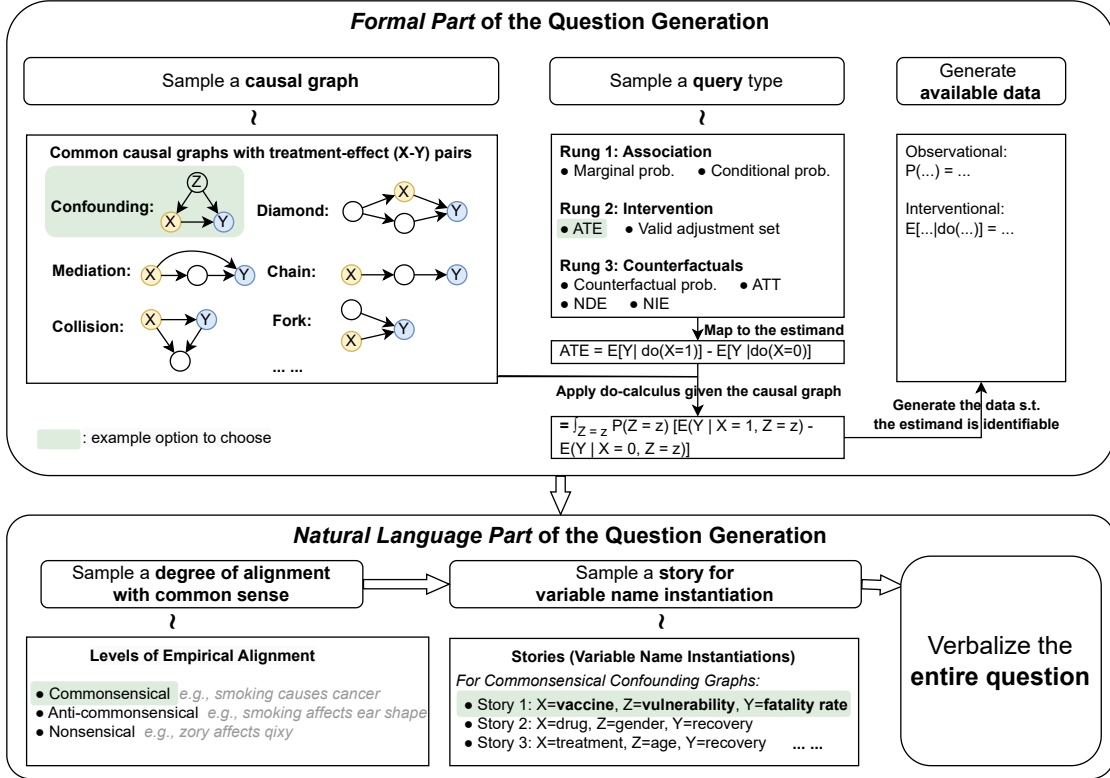


Figure 3.2: The data-generating process of the CLADDER dataset. The upper part of the figure describes the *formal part* of the question generation, which samples inputs for the CI Engine and derives a ground truth answer. The bottom part describes the *natural language part* of the question generation—i.e., its verbalization, based on multiple stories and different degrees of alignment with commonsense knowledge.

and, if so, an equivalent expression of the query which is estimable from the available data. While some previous works refer to the CI engine in the context of Rung 2 queries, where it corresponds to the *do*-calculus (Shpitser and Pearl, 2006a; Huang and Valtorta, 2006), here we refer to it in a more general sense, encompassing all three rungs.

3.3 Composing the CLadder Dataset

Task Formulation. Like in the example of Figure 3.1, our dataset $\mathcal{D} := \{(q_i, a_i, e_i)\}_{i=1}^N$ consists of N triples, each containing a question q_i , binary answer $a_i \in \{\text{Yes}, \text{No}\}$, and an explanation e_i . Our main task is to test the accuracy of the prediction function $f : q \mapsto a$, i.e., a LLM which maps a natural language causal question to an answer. Apart from directly evaluating the answer, we also compose the ground-truth explanations e to evaluate the reasoning steps of LLMs.

Design Principles. In the composition of our dataset, we adhere to the following design principles. First, we ensure broad coverage of all rungs of the ladder of causation. Second, we avoid settings that involve continuous variables and use binary variables instead: this is partly due to the large availability of identifiability results for binary and categorical variables, and partly because queries involving binary variables lend themselves to more

natural-sounding verbalization. Moreover, since LLMs struggle with calculation-heavy tasks (Hendrycks et al., 2021; Stolfo et al., 2023), and we are chiefly interested in causal reasoning abilities, we focus on graphs with few (three to four) variables, in various common configurations, to produce questions which are identifiable from the outset. Lastly, we carefully design a rich set of templates to translate the abstract formulas into grammatically correct and natural-sounding, fluent prompts.

Overall Pipeline. The generation pipeline for CLADDER, depicted in Figure 3.2, consists of two parts:

1. In the *Formal Part* (which we illustrate in Section 3.3.1), we specify all the required inputs (query, model, data) and the ground truth answer generated by the CI Engine.
2. In the *Natural Language Part* (in Section 3.3.2), we verbalize the formal queries and specification of the causal model and data by associating them to a story or narrative, using a rich set of templates.

3.3.1 Formal Part of the Question Formulation

The first step of our data generating process is to construct a set of inputs to the CI Engine such that *by design* there exists a well-defined ground truth answer: i.e., we construct triples of causal queries, graphs, and data such that the query can be unambiguously answered based on the available data (ensuring *identifiability* by construction).⁴ The ground truth causal models, which specify all quantities which are considered measurable in our questions, are causal Bayesian networks (CBNs), where each causal mechanism (i.e., conditional probability of a variable given its parents in the factorization according to the causal graph G) corresponds to a Bernoulli distribution. We compile a selection of graphs G based on examples drawn from multiple sources from the literature (Spirtes et al., 2001; Peters et al., 2017; Pearl, 2009b; Pearl and Mackenzie, 2018), where suitable graph structures are used to illustrate toy problems in causal inference. The complete list of structures we consider can be found in Appendix A.2.1.3; the complete list of sources in Appendix A.2.1.1.

Selecting Query Types. We again draw from the causal inference literature to collect common query types in each rung. As illustrated in the “*Sample a query type*” box in Figure 3.2, for Rung 1, we can ask about probability distributions such as marginal probabilities and conditional probabilities. For Rung 2 questions, we can enquire *average treatment effects* (ATE) (“*how will Y change if X changes from x to x' ?*”), or what constitutes a valid adjustment set that can block all backdoor spurious correlations between X and Y . Lastly, for Rung 3, we include *counterfactuals* (“*what would happen to Y had X been x' instead of x ?*”), *average treatment effect on the treated* (ATT) (“*for the subpopulation whose X changed from x to x' , how does their Y change on average?*”), *natural direct effect* (NDE) (“*what is the direct effect of X in Y , but not through the mediators?*”), and *natural indirect effect* (NIE) (“*what is the effect from X to Y through the mediators?*”).

⁴We use the term “data” to denote numerical values of conditional or *do*-probabilities, and not as collections of data samples. This is in line with how the term is used in other descriptions of the CI Engine (Pearl and Mackenzie, 2018; Hünemund and Bareinboim, 2019).

Applying the Causal Inference Engine for the Ground-truth answer. By construction, the causal processes we define encapsulates all necessary information to make the causal quantities of the query types identifiable. This allows us to apply the rules of causal inference to obtain an estimand for each causal graph and query type, and evaluate the estimand to get a ground truth answer. The Rung 2 queries simplify to Rung 1 terms using the rules of *do*-calculus (Pearl, 1995), and, for the Rung 3 queries, we apply methods of counterfactual causal inference (Pearl, 2009b) (with details in Appendix A.2.3.3). The estimand also specifies exactly which terms are necessary to include in the prompt as “*available data*” in order to ensure that enough information is provided to answer the question correctly (i.e., for identifiability), provided the correct causal reasoning is applied. Our entire code base of the data generation process can be found at our GitHub repository, <https://github.com/causalNLP/cladder>.

3.3.2 Natural Language Part of the Question Formulation

While Section 3.3.1 describes a way to generate the ground-truth causal model, query and answers, computed through a causal inference engine, real-world causal reasoning problems are expressed in natural language rather than symbolic expressions. The next part of the data generation pipeline therefore focuses on the verbalization of all these components with a plausible narrative in natural language.

Generating the Stories. For each causal graph, we collect a set of two to five *stories* which consist of a list of variable names for each node in the graph. The stories are primarily selected from examples in commonly cited causal inference books and papers (see Appendix A.2.1.1), which ensures that the stories and corresponding causal graph structures adhere to empirical common sense (e.g., the drug-gender-recovery example of Pearl and Mackenzie (2018)). However, it is very likely that at least some of the stories appear in the training data of many LLMs. Therefore, we also generate various *anti-common sense* and *nonsensical* variants of the stories, meant to isolate the effects of memorization. For the anti-commonsensual stories, we randomly do one of the actions: (1) replace the effect variable Y with an unusual attribute, that would not be an effect variable in any of the stories (e.g., “ear shape”); or (2) create an irrelevant treatment variable X that does not play a causal role in any of our commonsensual stories, such as “playing card games” (see Appendix A.2.1.7). For the nonsensical variants, we invent artificial words as variable names such as “zory” and “qixy” (see Appendix A.2.1.6).

Verbalizing the Prompts. The verbalization procedure applies the mapping of symbolic variables to semantic concepts to form a plausible narrative for the underlying causal process and then translates the symbolic expressions from the underlying causal process to natural language using carefully designed templates.

Specifically, we use several different grammatical forms for each semantic concept t in the story to make the resulting prompt sound natural and grammatically correct. We first have the overall variable name $v_{\text{overall}}(t)$ (e.g., the recovery status), and, then, for each binary value $i \in \{0, 1\}$, we compose its noun $v_{\text{noun}}(t = i)$ (e.g., recovery), verb (e.g., to recover), sentence $v_{\text{sent}}(t = i)$ (e.g., the patients recover), noun with attributive clause $v_{\text{attr}}(t = i)$ (e.g., patients who recover), and third conditional $v_{\text{cond}}(t = i)$ (e.g., if the patient had recovered).

Using these elements, we first verbalize the causal graph by iterating through each node and its outgoing edges, using the template “ t has a direct effect on $\mathbf{CH}(t)$.”, where $\mathbf{CH}(\cdot)$ denotes the set of direct effects (children) of a variable. Then, for the available data d , we verbalize each conditional probability by “For $v_{\text{attr}}(t_m = i)$, the probability of $v_{\text{noun}}(t_n = 1)$ is p .”, and each marginal probability by “The overall probability of $v_{\text{attr}}(t = 1)$ is p .” Note that our distributions are Bernoulli, so it is adequate to just introduce the parameter p , which is the likelihood of $t = 1$. For example, we generate sentences such as “The overall probability of recovery is 60%.” and “For patients who have small kidney stones, the probability of recovery is 70%.” Finally, for the query q , we instantiate each query type in our dataset following our question templates in Appendix A.2.1.5 such that the questions can always be answered with “yes” or “no”.

Generating the Explanations. Apart from the question-answer pairs, we also generate the step-by-step explanations. Our goal is to provide all intermediate reasoning steps a student of causal inference would use to answer the questions, so that each necessary subskill necessary for causal inference can be evaluated individually. We identify the following six subskills: ① causal graph extraction; ② correct query type interpretation; ③ symbolic formalization of the query; ④ semantic parsing to compile the available data; ⑤ estimand derivation; and ⑥ arithmetic calculation to solve the estimand, as in the colored boxes in Figure 3.1. Our explanation e verbalizes all the elements ①–⑥ as sequential steps using our template in Appendix A.2.1.8.

3.3.3 Dataset Statistics

Our data-generating procedure has the potential to algorithmically generate a vast large number of questions. In practice, we pick a dataset size that is large enough to be representative, and at the same time not too large to be problematic given the expensive inference costs of LLMs. We therefore set our dataset size to be 10K, and report the statistics in Table 3.1.

The dataset roughly balance across the query types, graph structures, stories, and ground truth answers (as seen in Figure 3.3). Note that some causal queries are only compatible with a subset of the graphs, thereby resulting in a slightly lower representation of those queries (such as the NDE and NIE). More details on our design choices can be found in Appendix A.2.1.4.

3.3.4 Data Quality Check

Our dataset is generated through an algorithmic procedure, which has the following potential benefits: formal correctness; zero human annotation cost; and, most importantly, controllability—e.g., for the question distribution, as well as for making it more unlikely that the data was previously seen by the model. However, since the dataset is different from common NLP datasets collected from human natural language writing, we also need to perform additional data quality checks. We therefore checked for a list of non-formal, natural language properties: grammaticality; human readability; naturalness/perplexity; and how well humans perform on this task.

For grammaticality, we ran a grammatical error check on our dataset using the Language-Tool package (Naber et al., 2003), and got on average 1.26 grammatical errors per 100

	Total	Rung 1	Rung 2	Rung 3
Size				
# Samples	10,112	3,160	3,160	3,792
Question				
# Sentences/Sample	6.01	5.88	5.37	6.65
# Words/Sample	80.9	73.43	76.95	90.42
# Nodes/Graph	3.52	3.5	3.5	3.54
# Edges/Graph	3.38	3.3	3.3	3.5
Answer				
Positive Class (%)	50	50	50	50
Explanations				
# Sentences/Sample	9.11	9.1	8.1	9.96
# Words/Sample	47.95	49.87	32.8	58.97

Table 3.1: Statistics of our CLADDER dataset v1.5.

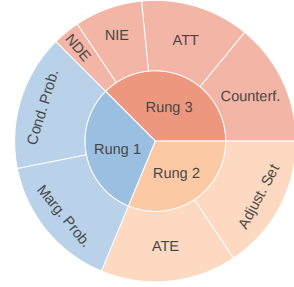


Figure 3.3: Distributions of query types in our 10K data.

words (i.e., 98.74% correctness), which shows that most of the language in our dataset follows English grammar. For human readability, we checked how comprehensible the questions are to students who have taken causality courses. We selected a random subset of 50 questions from the dataset, and let a graduate student annotator go through the questions to judge whether they could understand them or not: 96% of the questions were deemed readable. Next, for the naturalness/perplexity score, we used the open-sourced GPT-2 model and obtained a perplexity score of 21.17 on our dataset, which is substantially lower (i.e., closer to the distribution of natural human-written text) than the one of MATH (Hendrycks et al., 2021), a commonly used dataset of maths questions. Lastly, we conducted a sanity check where one expert evaluator tried to solve a random sample of 50 questions from the dataset, and we recorded an accuracy of 82% on this task.

3.4 Our CausalCoT Model

In order to guide LLMs in correctly answering the questions in CLADDER, we draw inspiration from the ideal functioning of the CI engine (Pearl and Mackenzie, 2018), which breaks down a causal reasoning problem into multiple symbolically-grounded, simpler steps. We develop CAUSALCoT, a multi-step causal chain-of-thought prompt in Figure 3.4, which combines formal causal reasoning skills with the idea of chain-of-thought prompting (Wei et al., 2022b) and the use of scratch pads for solving more complicated problems requiring a long list of steps (Nye et al., 2021) for LLMs.

We base our prompt design on the multi-step reasoning process of causal inference as shown in Figure 3.4, first starting with four preparation steps: ① identifying the causal graph structure; ② determining the causal query type;⁵ ③ formulating the query symbolically precisely; and ④ extracting relevant data from the prompt. Then, given all the information collected in the preparation stage, we introduce the formal solution: ⑤ correctly deducing the estimand using causal inference techniques; and finally ⑥ evaluating the estimand to answer the question. This set of steps require both *natural language understanding* to parse the question (as in most steps in the preparation phase), as well as *formal causal reasoning* to derive the correct estimand (as in the solution phase).

⁵This step amounts to a multi-class classification problem, where each class is a different causal query.

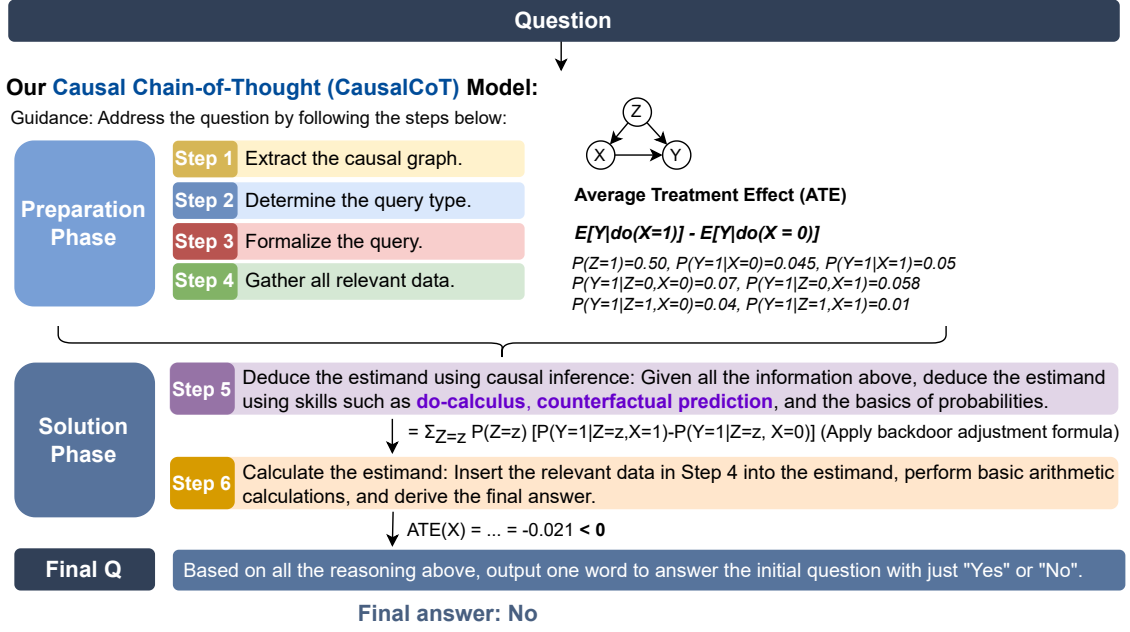


Figure 3.4: Illustration of our CAUSALCOT prompting strategy, which designs a chain of subquestions inspired by the idea of a CI engine (Pearl and Mackenzie, 2018).

We build our CAUSALCOT prompting strategy using GPT-4 (OpenAI, 2023), a recent autoregressive LLM that achieves state-of-the-art performance on many tasks. This latest model builds upon the previous series of general pretrained models (GPT) (Radford et al., 2019; Brown et al., 2020) and adds reinforcement learning with human feedback, or instruction-tuning (Ziegler et al., 2019; Ouyang et al., 2022; Bai et al., 2022), to align the model responses to free-form questions with human preferences. It has achieved human-competitive performance over a list of tasks (OpenAI, 2023; Bubeck et al., 2023; Nori et al., 2023; Katz et al., 2023; Ziems et al., 2024), among which the more formal tasks unseen in the training data still remain elusive (Razeghi et al., 2022; Stolfo et al., 2023; Jin et al., 2024).

Given a causal question q , we provide the LLM a list of instructions $\ell := (s_1, \dots, s_6)$ consisting of the detailed descriptions of the six steps s_1, \dots, s_6 in Figure 3.4. As the model $f_{\text{LLM}} : s_i \mapsto r_i$ autoregressively produces responses r_1, \dots, r_6 sequentially corresponding to the six steps, we concatenate all the above before asking the final question “Based on all the reasoning above, output one word to answer the initial question with just ‘Yes’ or ‘No’.” See the complete prompt in Appendix A.2.2.1. In the end, we obtain the binary answer $a \in \{\text{Yes}, \text{No}\}$ as the final result.

Compared with the standard strategy of directly prompting the LLMs a question, we impose an *inductive bias* upon LLMs by using the causal inference framework, thus incorporating some of the powerful, principled insights of the causal inference community for NLP tasks. In this way, we enhance the strong natural language ability of LLMs with formal causal reasoning skills.

3.5 Testing LLMs with CLadder

3.5.1 Experimental Setup

Our empirical investigation focuses on some of the most recent language models. We include the latest GPT-4 (OpenAI, 2023) with 1T parameters by the time we conduct the experiments (i.e., gpt-4-1106-preview), the previous ChatGPT (i.e., GPT-3.5) with 175B parameters, and then a series of earlier models with instruction-tuning on the 175B GPT-3 (text-davinci-001, -002, and -003) (Ouyang et al., 2022). As baselines, we also include the non-instruction-tuned GPT-3 (davinci). We use the OpenAI API with temperature 0 when querying these models. We also include open-source, more efficient models like LLaMa (Touvron et al., 2023) and its instruction-tuned version Alpaca (Taori et al., 2023a), both with the same number of parameters, 6.7B.

3.5.2 Main Results

	Overall Acc.	Acc. by Rung			Acc. by Commonsense Alignment		
		1	2	3	Comm.	Nonsens.	Anti-C.
Random	49.27	50.28	48.40	49.12	49.01	49.69	49.12
LLaMa	44.03	48.23	29.46	52.66	45.14	44.22	42.67
Alpaca	44.66	52.03	29.53	51.13	44.86	44.40	44.77
GPT-3 Non-Instr. (davinci)	49.92	50.00	49.75	50.00	49.06	49.97	50.72
GPT-3 Instr. (text-davinci-001)	51.40	51.30	52.63	50.47	54.31	50.13	50.05
GPT-3 Instr. (text-davinci-002)	53.15	50.85	56.96	51.90	55.33	52.47	51.81
GPT-3 Instr. (text-davinci-003)	56.26	51.11	62.97	54.96	56.83	54.79	57.49
GPT-3.5	52.18	51.80	54.78	50.32	54.09	50.68	52.09
GPT-4	62.03	63.01	62.82	60.55	62.27	63.09	60.47
+ CAUSALCoT	70.40	83.35	67.47	62.05	69.25	71.58	70.12

Table 3.2: Performance of all models on our CLADDER dataset v1.5. We report the overall accuracy (Acc.), and also fine-grained accuracy by rung, and by degree of commonsense alignment, from commonsensical (Comm.), nonsensical (Nonsens.), to anti-commonsensical (Anti-C.).

We compare the performance of all models in Table 3.2. First, we can see that the causal reasoning task in CLADDER is in general very challenging for all models. Models such as the earlier, non-instruction-tuned GPT-3, and both LLaMa and Alpaca are around random performance. With instruction-tuning, models start to show some improvement. And amongst all, our CAUSALCoT achieves the highest performance of 70.40%, which is substantially better than the vanilla GPT-4 by 8.37 points. Moreover, CAUSALCoT also achieve the best performance across all three rungs of causal questions, with a monotonically decreasing performance as the rungs get higher, i.e., the questions get more difficult. See Appendix A.2.4 for experiments on our earlier dataset v1.0.

3.5.3 Isolating the Effect of Data Contamination

A well-known problem with evaluating LLMs on question-answering tasks is the data contamination problem, i.e., that LLMs perform well on a test set because the test set is (unintentionally) contained partially or even entirely in the training data (OpenAI, 2023; Brown et al., 2020). We address this problem by creating not only the commonsensical subset of our dataset, but also anti-commonsensical and nonsensical, both of which, by

construction, are very likely not in the training data of LLMs. From the accuracy by commonsense alignment degree in Table 3.2, we can see the original GPT-4 model performs the worst on the anti-commonsensical subset (1.8 points lower than that on the commonsensical subset). However, our CAUSALCoT enhances the reasoning ability across all levels, with substantial improvement on anti-commonsensical data by 9.65 points, highlighting the strength of CAUSALCoT on unseen data.

3.5.4 Error Analysis by Subquestions

Step ①			Step ②				Step ③ & ⑤	Step ④	Step ⑥
Node	Edge	Dist. (↓)	Overall F1	Rung 1	Rung 2	Rung 3	Estimand	F1	Arithmetic
99.34	97.01	1.69	50.65	69.99	59.14	42.12	53	47.53	99

Table 3.3: Performance for each step in CAUSALCoT. For Step ①, we report the F1 score of node prediction, edge prediction, and also the graph edit distance (Dist.) with the true graph. See more details in Appendix A.2.5.1.

We conduct a fine-grained error analysis by looking into the performance of different steps of CAUSALCoT in Table 3.3.⁶ We can see that the model is good at Step ① to extract causal graph \mathcal{G} , achieving high F1 scores for predicting both the nodes and the edges correctly, although not perfect, still leaving a graph edit distance of 1.69 between the ground truth causal graph and the model-identified graph. The other steps are more challenging for the model. Among those, Steps ②, ③ and ⑤ require careful and correct application of causal inference, where the model struggles. This reveals a notable weakness of current LLMs to perform formal causal reasoning, which is an important direction for future work on improving and enhancing LLMs. To better understand the reasoning abilities of LLMs, we also perform an extensive analysis taking the entire reasoning chain of our CAUSALCoT and the ground-truth explanations, to produce 20 fine-grained scores about the multi-step reasoning quality using the ROSCOE framework (Golovneva et al., 2022), and show detailed results in Appendix A.2.5.2.

3.5.5 Effect of In-Context Learning

As an additional analysis, we look into the effect of in-context learning (ICL) by providing an example solution before asking the question. The interesting question to us is whether models can generalize across different query types. Namely, we keep our CAUSALCoT framework, and prepend a reasoning example of query type i , and then calculate how much improvement it can bring when models answer new questions of query type j . In Figure 3.5, we can see that conditional probability and NIE are the questions that benefit the most from ICL, and showing examples of marginal probability and ATT are among the most helpful to all questions in general.

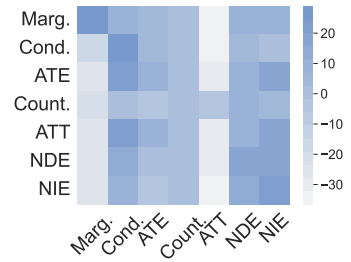


Figure 3.5: Heatmap showing the how helpful each query type is to solving subsequent query types.

⁶We experienced some rate-limiting in the fine-grained analysis of LLMs that are only accessible through a web API. As a result, we occasionally had to evaluate on a subset of 2K random samples.

3.6 Related Work

Skill evaluation for LLMs. Our work may be seen as part of the literature aimed at evaluating the performance of current LLMs (Radford et al., 2019; Devlin et al., 2019; Brown et al., 2020; Zhang et al., 2022; OpenAI, 2023, *inter alia*), focusing on understanding their strengths and weaknesses. Various studies into the capabilities of LLMs (Bubeck et al., 2023; Qin et al., 2023; OpenAI, 2023; Ignat et al., 2024) change people’s perception of domains such as education (Baidoo-Anu and Owusu Ansah, 2023; Rudolph et al., 2023), medicine (Singhal et al., 2022; Nori et al., 2023), law (Katz et al., 2023), and computational social science (Ziems et al., 2024). However, most work evaluates new models on existing datasets from previously-curated large-scale benchmarks (Wang et al., 2019, 2022; Srivastava et al., 2022), or human exams (Katz et al., 2023; OpenAI, 2023; Jin et al., 2022b) which is becoming increasingly unreliable due to training set contamination.

Causality-related skills for NLP. With the increasing attention on LLMs and causality (Zečević et al., 2023; Zhang et al., 2023), we review several formulations of causality-related skills for NLP, which we summarize into (1) causality as knowledge, (2) causality as language comprehension, and (3) causality as reasoning. In the *causality-as-knowledge* line of work, many existing studies investigate how well NLP models understand commonsense causality, such as the cause and effect of an agent’s action (Sap et al., 2019a), motivation and emotional reaction in a social context (Sap et al., 2019b), correspondence of a set of steps with a high-level goal (Zhang et al., 2020a), development of a story given a different beginning (Qin et al., 2019), and how in general LLMs serve as a knowledge base of causality (Zečević et al., 2023). Concurrent work (Kıcıman et al., 2023) focuses on evaluating LLMs on various causality related tasks by leveraging the conceptual knowledge accrued from the training data, rather than formal causal inference, except for their causal sufficiency analysis which is close to our counterfactual questions. Importantly, most work in this line does not define explicit causal graphs, making it difficult to quantitatively define the ground-truth causal relationships in a principled way. The *causality-as-language-comprehension* line of work stems from traditional linguistic studies on causal connectives and causal language usage (Stede, 2008; Cao et al., 2022; Yu et al., 2019), to the recent causal relation extraction (Bethard et al., 2008; Hidey and McKeown, 2016; Xu et al., 2020) to identify cause-effect pairs as a subtask of information extraction from text.

Finally, for *causality as formal reasoning*, our CLADDER work formulates the task of causal inference for NLP, and our other work, CORR2CAUSE (Jin et al., 2024), addresses the causal discovery problem to infer causation from correlation. Together, they cover the two major branches of causal reasoning investigated in existing technical literature on causality. See a comprehensive comparison of literature in Appendix A.2.6.

3.7 Discussion of Limitations and Future Work

A Natural Language “Mini Turing Test” for Causality. Pearl and Mackenzie (2018) describe an ideal “mini-Turing test” to assess understanding of causal inference, and argue that if a machine can answer all possible questions correctly, then it “understands” causality. According to the authors, this is because there are no possible shortcuts when you consider all possible combinations of queries, graphs and data in this ideal test: due to

their combinatorial explosion, the machine can only answer all questions right if it correctly applies causal reasoning. From this point of view, our work constitutes a *first step towards a mini-Turing test formulated in natural language*. However, we cover only some of the commonly studied causal queries spanning all three rungs. Future work may extend this to further queries, such as, e.g., path-specific effects other than NDE and NIE (Nabi and Shpitser, 2018), thereby increasing the number of potential questions and moving closer to the ideal test.

LLMs and Causal Reasoning. It has been claimed that LLMs understand causality well (e.g., (Kıcıman et al., 2023) report high performance, such as 97% and 92%). In contrast, our work suggests that LLMs may still be far from reasoning reliably about causality (reaching only 60+% on CLADDER). As argued in Section 3.1, we believe that investigating this aspect may be of particular importance, since causal inference is crucial in many policy-relevant scenarios, where reliable AI systems could assist decision-making: from epidemiology (Glass et al., 2013; Rothman and Greenland, 2005) to economics (Card, 1999; Hünermund and Bareinboim, 2019) to fairness (Loftus et al., 2018; Plecko and Bareinboim, 2022). Testing the abilities of these systems in semi-realistic scenarios is therefore crucial, motivating some of the design choices in our dataset: e.g., the example in Figure 3.1 was inspired by similar questions which arose in the context of the COVID-19 pandemic, where incorrect causal reasoning resulted in a fallacy where vaccinations were considered to be harmful instead of beneficial (Morris, 2021; Ellenberg, 2021). Further work may be dedicated to making the questions and verbalizations even closer to realistic instances of causal inference problems.

A CI Engine Plug-in for LLMs. An interesting direction for future research could be to provide the LLM access to an actual implementation of the CI engine. For example, Davis and Aaronson (2023) tested the improvement of math abilities in LLMs augmented with plug-ins (i.e., external modules that extend the model’s capabilities by adding specific functionality or customizing its behaviour for particular tasks, like a calculator), suggesting that they significantly enhance the model’s ability to solve these problems. However, even with plug-ins, there are still often “*interface*” failures: that is, “[the LLM] often has trouble formulating problems in a way that elicits useful answers from the plug-ins”. We hypothesise that something similar would happen for causal inference: even once suitable plug-ins are built, the language-to-tool interface may still be a non-trivial research question.

3.8 Conclusion

We proposed formal causal reasoning as a new task to evaluate LLMs, and created the CLADDER benchmark, covering several aspects of causal inference across all rungs of the ladder of causation and verbalizations involving semi-realistic scenarios. To address the task, we proposed a prompting strategy, CAUSALCoT, inspired by the principles of formal causal inference, which introduces multistep chain-of-thought reasoning for causal questions. Extensive experiments indicate that this dataset is highly challenging, thus offering a principled tool to gain a better understanding of the reasoning abilities of LLMs and to develop better models for causal reasoning in natural language.

Part II

Causal Understanding of How LLMs Work

Competition of Mechanisms: Tracing How LLMs Handle Facts and Counterfactuals

Interpretability research aims to bridge the gap between empirical success and our scientific understanding of the inner workings of large language models (LLMs). However, most existing research focuses on analyzing a single mechanism, such as how models copy or recall factual knowledge. In this work, we propose a formulation of *competition of mechanisms*, which focuses on the interplay of multiple mechanisms instead of individual mechanisms and traces how one of them becomes dominant in the final prediction. We uncover how and where mechanisms compete within LLMs using two interpretability methods: logit inspection and attention modification. Our findings show traces of the mechanisms and their competition across various model components and reveal attention positions that effectively control the strength of certain mechanisms.

Our code is available at <https://github.com/francescortu/comp-mech>, and our data at <https://huggingface.co/datasets/francescortu/comp-mech>.

4.1 Introduction

Recent advancements in large language models (LLMs) have brought unprecedented performance improvements to NLP (Brown et al., 2020; Touvron et al., 2023; OpenAI, 2023; Anil et al., 2023, *inter alia*). However, the black-box nature of these models obfuscates our scientific understanding of *how these models achieve certain capabilities*, and *how can we trace the problem when they fail at other tasks*. This has brought an increasing focus on interpretability research to help us understand the inner workings of LLMs.

Existing interpretability research has been largely focused on discovering the *existence* of single mechanisms, such as the copy mechanism in induction heads of LLMs Elhage et al. (2021); Olsson et al. (2022), and factual knowledge recall in the MLP layers (Geva et al., 2021; Meng et al., 2022; Geva et al., 2023). However, different from discovering *what mechanisms exist in LLMs*, we propose a more fundamental question: *how do different mechanisms interact in the decision-making of LLMs?* We show a motivating example in Figure 4.1, where the model fails to recognize the correct mechanism when it needs to judge between two possible mechanisms: whether to recall the factual knowledge on who developed the iPhone (i.e., Mechanism 1) or to follow its counterfactual redefinition in the new given context (i.e., Mechanism 2).

We propose a novel formulation of *competition of mechanisms*, which focuses on tracing each mechanism in the model and understanding how one of them becomes dominant in the final prediction by winning the “competition”. Specifically, we build our work

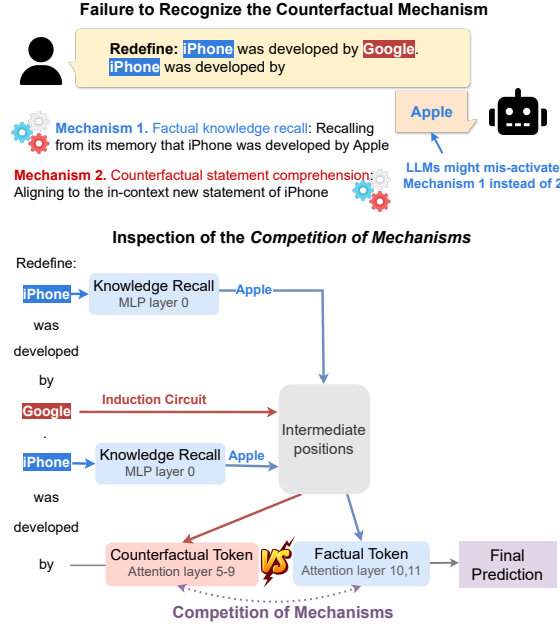


Figure 4.1: Top: An example showing that LLMs can fail to recognize the correct mechanism when multiple possible mechanisms exist. Bottom: Our mechanistic inspection of where and how the competition of mechanisms takes place within the LLMs.

on two single mechanisms that are well-studied separately in literature: (1) the factual knowledge recall mechanism, which can be located in the MLP layers (Geva et al., 2021; Meng et al., 2022; Geva et al., 2023); and (2) the in-context adaptation to a counterfactual statement, which is enabled by the copy mechanism conducted by induction heads of attention layers Elhage et al. (2021); Olsson et al. (2022). Based on the latest tools to inspect each of these two mechanisms Nostalgebraist (2020); Wang et al. (2023); Geva et al. (2023), we then unfold *how and where* the competition of the two mechanisms happen, and how it leads to the overall success or failure of LLMs.

Technically, we deploy two main methods: logit inspection Nostalgebraist (2020); Geva et al. (2022) by projecting the outputs of each model component by an unembedding matrix, and attention modification Geva et al. (2023); Wang et al. (2023). Using these methods, we assess the contributions of various model components, both from a macroscopic view (e.g., each layer) and a microscopic view (e.g., attention heads), and identify critical positions and attention heads involved in the competition of the two mechanisms. Moreover, we locate a few localized positions of some attention head matrices that can significantly control the strength of the factual mechanism. We summarize our main findings as follows:

1. In early layers, the factual attribute is encoded in the subject position, while the counterfactual is in the attribute position (Section 4.6.1);
2. The attention blocks write most of the factual and counterfactual information to the last position (Section 4.6.2);
3. All the highly activated heads attend to the attribute position regardless of the specific type of information they promote. The factual information flows by penalizing the counterfactual attribute rather than promoting the factual one (Section 4.6.3);

4. We find that we can up-weight the value of a few very localized values of the attention head matrix to strengthen factual mechanisms substantially (Section 4.6.4).

4.2 Related Work on Interpretability

As deep learning approaches show increasingly impressive performance in NLP, their black-box nature has hindered the scientific understanding of these models and their effective future improvements. To this end, interpretability research has been a rising research direction to understand the internal workings of these models.

Interpreting the Representations. One major type of work in interpretability has focused on understanding what has been encoded in the representations of deep learning models. This is usually achieved by a *probe*, namely by training a supervised classifier to predict features from its representations (Alain and Bengio, 2016; Conneau et al., 2018; Hupkes et al., 2018; Hewitt and Liang, 2019; Tenney et al., 2019; Jiang et al., 2020; Elazar et al., 2021, *inter alia*), or by geometric methods Doimo et al. (2020); Valeriani et al. (2024); Park et al. (2024); Cheng et al. (2024). Example features of interest include part of speech Belinkov et al. (2017), verb tense Conneau et al. (2018), syntax Hewitt and Manning (2019), and factual knowledge Petroni et al. (2019).

Interpreting the Mechanisms/Functions. Beyond interpreting the representations in the hidden states of the black-box models, another research direction is to interpret the mechanisms or functions that the models have learned, giving rise to the field of mechanistic interpretability (Olah et al., 2020; Elhage et al., 2021; Olsson et al., 2022; Nanda et al., 2023, *inter alia*). Some example mechanisms decoded in recent work include mathematical operations such as modular addition Nanda et al. (2023) and the greater-than operation (Hanna et al., 2023); natural language-related operations such as the copy mechanism achieved by induction heads in LLMs Olsson et al. (2022) and factual knowledge recall achieved by MLP layers (Geva et al., 2021; Meng et al., 2022; Geva et al., 2023), which we describe below.

The Single Mechanism of Copy: One of the basic actions in LLMs is the copy mechanism, which is found to be operationalized by attending to the copied token in the attention heads and passing it on to the next token prediction Elhage et al. (2021); Olsson et al. (2022). This foundational mechanism enables further research to decode more complex mechanisms, such as indirect object identification Wang et al. (2023).

The Single Mechanism of Factual Knowledge Recall: Another major direction is understanding how LLMs mechanistically recall factual information (Geva et al., 2021; Meng et al., 2022; Geva et al., 2023). For example, Meng et al. (2022) develop the *causal tracing* method to show that the factual information is found in the mid-layer MLP units in GPT-2. A followup work Geva et al. (2023) shows that MLPs of early layers enrich the subject embeddings with related attributes, and late attention blocks select and write the correct factual information to the sentence’s last position.

Interplay of Multiple Mechanisms: In the final stage of our project in December 2023, we noticed a related study by Yu et al. (2023), which also investigates the role of two different mechanisms in LLMs. Specifically, they inspect a type of prompt whose subjects are the

capital cities and whose attributes are the countries, examine the dynamics of the factual recall mechanism and the effect of the in-context counterfactual statement, and find that the subject and attribute frequency in the pre-training set can affect the ability of factual recall. Differently, the methods in our work are applied to a broader set of prompts; moreover, we also establish novel analyses of the underlying mechanistic details of the competition, and precisely localize the path where the information flows at the level of single attention map activations, based on which we discover new findings that are unique to our study.

4.3 Problem Setup

Following the setup of many existing interpretability studies (Olah et al., 2020; Elhage et al., 2021; Olsson et al., 2022; Nanda et al., 2023, *inter alia*), we look into the next token prediction behavior of autoregressive LLMs in their inference mode, namely

$$P(t_k | t_{<k}), \quad (4.1)$$

which predicts the k -th token t_k given all the previous tokens in the context.

Next, we design the task to incorporate the competition of mechanisms as in Figure 4.1. Specifically, for each factual statement $f := (t_1^f, \dots, t_k^f)$ consisting of k tokens (e.g., “iPhone was developed by Apple.”), we compose a corresponding counterfactual statement $c := (t_1^c, \dots, t_k^c)$ (e.g., “iPhone was developed by Google.”). Then, we compose a prompt connecting the two statements as “Redefine: c . $f_{1:k-1}$.”, such as “Redefine: iPhone was developed by Google. iPhone was developed by ___”.

The two mechanisms can be traced by inspecting the rise and fall of the factual token t_k^f and the counterfactual token t_k^c . For the simplicity of notation, we take the tokens out of the context of their exact position and denote them as t_{fact} and t_{cofa} , respectively, in the rest of the paper.

4.4 Method and Background

Method 1: Logit Inspection. To inspect the inner workings of the two mechanisms, we trace the *residual stream* Elhage et al. (2021), or logits of each component in the LLM. Given a text sequence of k tokens, LLMs map it into the residual stream, namely a matrix $x \in \mathbb{R}^{d \times k}$, where d is the dimension of the internal states of the model. We use the term x_i^l to specify the residual stream at position i and layer l .

An LLM produces the initial residual stream x_i^0 by applying an embedding matrix $W_E \in \mathbb{R}^{|V| \times d}$ to each token t_i , where $|V|$ is the size of the vocabulary. Then, it modifies the residual stream by a sequence of L layers, each consisting of an attention block a^l and MLP m^l . Finally, after the last layer, it projects the internal state of the residual stream back to the vocabulary space with an unembedding matrix $W_U \in \mathbb{R}^{d \times |V|}$. Formally, the update of the residual stream at the l^{th} layer is:

$$x^l = x^{l-1} + a^l + m^l, \quad (4.2)$$

where both the attention and the MLP block take as input the x after layer normalization

norm:

$$\mathbf{a}^l = \mathbf{a}^l(\text{norm}(\mathbf{x}^{l-1})) , \quad (4.3)$$

$$\mathbf{m}^l = \mathbf{m}^l(\text{norm}(\mathbf{x}^{l-1} + \mathbf{a}^l)) . \quad (4.4)$$

To understand which token the residual stream \mathbf{x}^l favors, we follow the common practice in previous work (Halawi et al., 2023; Geva et al., 2023; Dar et al., 2023; Geva et al., 2022) to project it to the vocabulary space using the aforementioned unembedding matrix W_U which maps the latent embeddings to actual tokens in the vocabulary, enabling us to obtain the logits of the factual t_{fact} and counterfactual token t_{cofa} .

Known as the *Logit Lens* (Nostalgebraist, 2020), this method is broadly adopted due to its consistent success in yielding interpretable results, demonstrating its effectiveness through broad empirical usage. However, it is important to note that it can occasionally fail to reflect the actual importance of vocabulary items, especially in the early layers of the network (Belrose et al., 2023).

Method 2: Attention Modification. Modifying or ablating the activation of a specific model component is also a strategy used to improve the understanding of the information flow within LLMs, including techniques such as causal tracing Meng et al. (2022) and attention knockout (Wang et al., 2023; Geva et al., 2023).

In our work, we focus on modifying a small number of entries in the attention matrix. Namely, in the attention matrix A^{hl} of the h -th head of the l -th attention layer \mathbf{a}^l , we focus on a certain entry, e.g., at the (i, j) position, where $j < i$, which is the attention value of the token \mathbf{x}_i^l attending to one of its earlier tokens \mathbf{x}_j^l . Following recent work Yu et al. (2023), the modification is after the softmax layer, so the other attention values of the matrix stay unchanged. For the target entry A_{ij}^{hl} , we scale it up by a multiplier of α :

$$A_{ij}^{hl} \leftarrow \alpha \cdot A_{ij}^{hl}, \quad \text{where } j < i . \quad (4.5)$$

4.5 Experimental Setup

Data Creation To compose the factual and counterfactual statements as introduced in Section 4.3, we adopt COUNTERFACT¹ (Meng et al., 2022), commonly used dataset to interpret models’ ability of factual knowledge recall. We select 10K data points by considering only examples where the attributes are represented by a single token and where the model completes the sentence in a factually accurate manner.

Each instance of COUNTERFACT expresses a relation r between a subject s and an attribute a : (s, r, a) . For example, in the sentence “*iPhone was developed by Apple*”, $s = \text{“iPhone”}$, $r = \text{“was developed by”}$, $a = \text{“Apple”}$. Moreover, each (s, r) instance is provided two values of the attribute a , namely a factual token t_{fact} , and a counterfactual token t_{cofa} , representing a false fact.

Using this source data, we compose each instance of our test set in the format of (“*Redefine:*”, $s, r, t_{\text{cofa}}, s, r, _$), such as “*Redefine: iPhone was developed by Google. iPhone was developed by __*”. We preprocess the original dataset by keeping only the data points whose attribute is a single

¹<https://rome.baulab.info/data/>

token (for the simplicity of our implementation), and where the model correctly predicts the factual token t_{fact} when completing the sentence $(s, r, _)$. We randomly select 10K test samples into our test set from 219,180 such samples. We open-source our dataset at <https://huggingface.co/datasets/francescortu/comp-mech>.

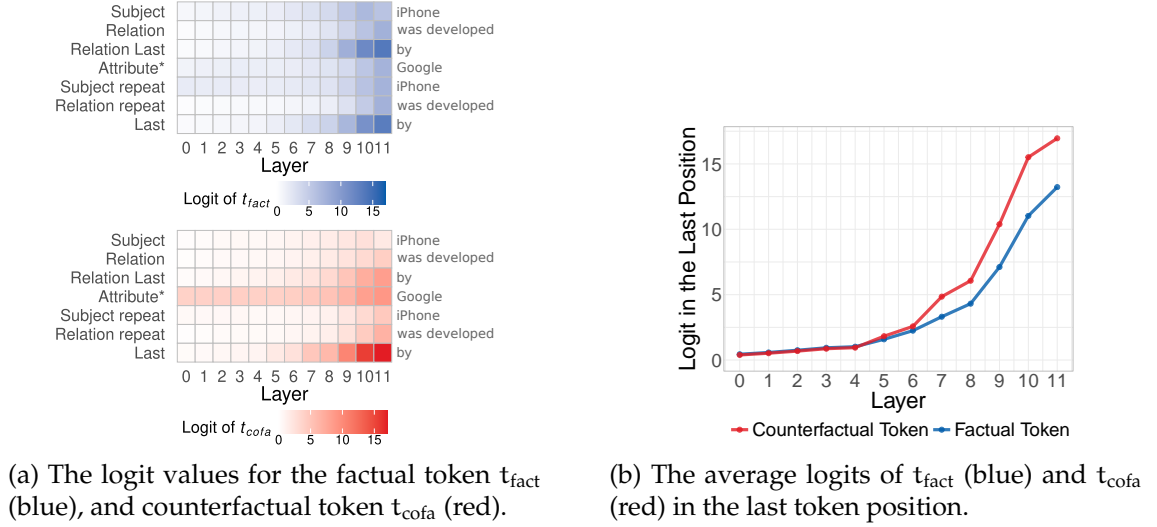


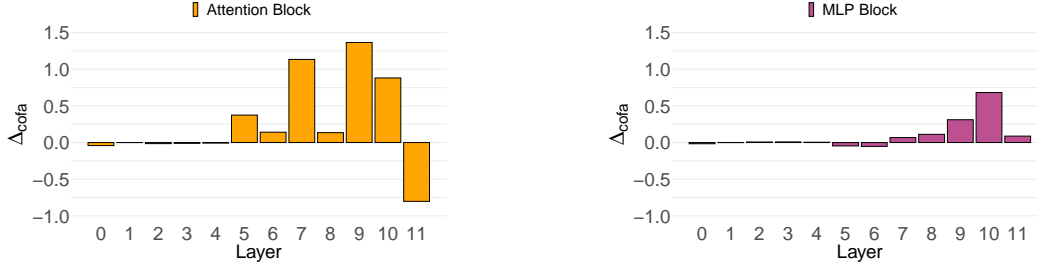
Figure 4.2: Logits of the factual token t_{fact} and counterfactual token t_{cofa} across different positions and layers in GPT-2. The logit of t_{fact} is higher in the subject position in the initial layers and in the last position of the premise and second sentence in the final layers. The logit of t_{cofa} is higher in the attribute position in the first layers and in the last position of the second sentence at the end of the network.

Models For this work, we first choose the GPT-2 small (Radford et al., 2019) model as it is the most commonly used one in previous interpretability studies (e.g., Meng et al., 2022; Wang et al., 2023; Conmy et al., 2023; Hanna et al., 2023). Aligning the same model with those studies can communicate the findings of this work better in the context of existing literature. Then, in addition to GPT-2, we check the generalizability of our work by provide supplemental results of Pythia-6.9B (Biderman et al., 2023) in Appendix A.3.1, to show the robustness of our findings across the two LLMs of different architectures and scales. In this way, having similar results across the two very diverse models makes the finding stronger than existing studies, most of which are only on GPT-2.

Implementation Details As for experimental details, GPT-2 small has 117M parameters, consisting of 12 layers with 12 self-attention heads each and a residual stream of 768 dimensions. Pythia-6.9B has 32 layers with 32 self-attention heads each and a model dimension of 4,096, with a 30x increase in the number of parameters. For all our experiments, we deploy the pre-trained models from the Huggingface Hub Wolf et al. (2020), and inspect the residual streams by the LogitLens tool in the TransformerLens library (Nanda and Bloom, 2022).

4.6 Results and Findings

In this section, we trace the competition of the mechanisms within the LLM via the two methods introduced in Section 4.4, i.e., inspecting the residual stream and intervening on



(a) Logit difference Δ_{cofa} of the last token position after the attention block in each layer of GPT-2.

(b) Logit difference Δ_{cofa} of the last token position after the MLP block in each layer of GPT-2.

Figure 4.3: Contributions of the attention and MLP blocks to the competition of the mechanisms. The attention blocks (left) contribute more to the marginal win of the counterfactual mechanism than the MLP blocks (right).

the attention. We provide mechanistic analyses on five research questions in the following subsections:

1. Macroscopic view: Which layers and token positions contribute to the two mechanisms? (Section 4.6.1)
2. Intermediate view: How do we attribute the prediction to attention and MLP blocks? (Section 4.6.2)
3. Microscopic view: How do individual attention heads contribute to the prediction? (Section 4.6.3)
4. Intrinsic intervention: Can we edit the model activations to modify the strength of a certain mechanism? (Section 4.6.4)
5. Behavioral analysis: What word choice varies the strength of the counterfactual mechanism in the given context? (Section 4.6.5)

4.6.1 Macroscopic Inspection across Layers and Token Positions

In the main model that we inspect, GPT-2, we find that it can usually identify the counterfactual mechanism in 96% of the 10K test examples. This means that, in the last sequence position, at the output of the network, the counterfactual token, t_{cofa} , gets most of the times a higher probability than t_{fact} . In the following, we will inspect how the “winning” of the counterfactual mechanism happens across the layers of the LLM.

Method. We study how t_{fact} and t_{cofa} are encoded in the residual stream using the logit inspection method described in Section 4.4. Specifically, for a given token position i and a layer l , we project the embedding x_i^l , i.e., the residual stream in Eq. (4.2), to the vocabulary space by $\mathbf{A}_i^l = W_U \cdot \text{norm}(x_i^l)$, where W_U is the unembedding matrix and norm is the normalization of the last layer of the model. By varying l , we measure the values of the logits of t_{fact} and t_{cofa} as they evolve in the residual stream after the first attention block.

Results. Our results reveal the prevalence of each mechanism by varying the layer l and position i .

Finding 1: Information flows from different tokens for different mechanisms. We analyze the

role of previous context at various token positions with respect to different depths of the layers. In Figure 4.2a, the blue heatmap above shows the logits of the factual token t_{fact} , and the red heatmap below shows those of the counterfactual token t_{cofa} .

Looking at the blue heatmap, we see that the *subject* position is the main contributor to the logits of t_{fact} in early layers, which is consistent with a previous finding Geva et al. (2023). Specifically, we also locate the factual attribute in the subject positions by the first MLP layer, and find they increase on average the value of t_{fact} from 0.38 to 0.74 in the premise and from 0.9 to 1.93 in the second sentence. Then, in the later layers, the strongest contributor is the last tokens before the attribute, as the last token position is used to predict the attribute. From the red heatmap, we see the evolution of t_{cofa} 's logits. The observations of later layers are similar across two mechanisms, in that the last token contributes the most. However, in early layers, the counterfactual mechanism's t_{cofa} token is best encoded in the *attribute* position instead of the subject position for the factual mechanism.

Such information flow between different token positions suggests a major role played by the attention mechanism in moving such information to the last position, resonating with observations in Geva et al. (2023).

Finding 2: Both the individual mechanisms and competition take place in late, but not early layers. We trace the competition of the two mechanisms across the layers by plotting in Figure 4.2b the scale of the logits corresponding to the two mechanisms in the last token position. The first observation is that the strength of each individual mechanism increases monotonically across the layers, from a relatively small logit below 1 in early layers to large values of around 15 in the final layer.

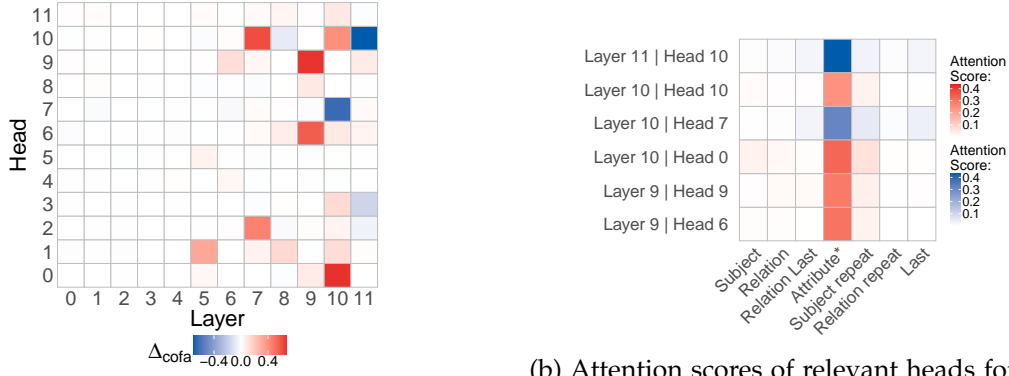
Another observation is that, although both mechanisms increase in strength, stronger signals of the competition (where the counterfactual mechanism prevails the factual one) start after the fifth layer, and this prevalence gradually grows in later layers. The logits of the counterfactual mechanism are, in most of the examples, the highest in the 50K-dimensional vocabulary of GPT-2, making t_{fact} dominant in 96% of the examples.

4.6.2 Intermediate Inspection of Attention and MLP Blocks

Behind the overall win of the counterfactual mechanism, we want to trace the contributions from the attention and MLP blocks in each layer.

Method. For each attention or MLP block, it processes the input embedding and outputs the logits of t_{fact} and t_{cofa} to be added to the residual stream. We can consider the contribution of each block as its added logit values to the residual stream. Intuitively, if the added logit value for t_{cofa} is higher than that of t_{fact} , then this block pushes the overall prediction to lean towards the counterfactual mechanism; otherwise, this block suppresses the counterfactual mechanism.

Hence, we inspect the margin of the added logit of t_{cofa} over that of t_{fact} in each block, represented by $\Delta_{\text{cofa}} := \text{BlockLogit}(t_{\text{cofa}}) - \text{BlockLogit}(t_{\text{fact}})$. To this end, we apply the logit inspection method to analyze the logit distribution at $W_U \mathbf{a}_N^l$ and $W_U \mathbf{m}_N^l$, where N denotes the last token position in the sequence. The logit contribution of the attention block is the sum over that of all the attention heads. As for the result, a positive value of



(a) Direct contribution to Δ_{cofa} of all heads in GPT-2. Heads favoring t_{fact} are colored in blue, and those favoring t_{cofa} in red.

(b) Attention scores of relevant heads for the last position. A large attention score in the attribute position is found in all highly activated heads.

Figure 4.4: Attention pattern for relevant attention heads.

Δ_{cofa} for a block means that it supports the counterfactual mechanism in the competition, and a negative value indicates suppression.

Results. We quantify the contribution of each block in each layer by plotting the Δ_{cofa} values in Figure 4.3.

Finding 1: The attention blocks play a larger role in the competition of mechanisms than the MLP blocks. Contrasting the Δ_{cofa} margin of the added logits of the attention blocks in Figure 4.3a and MLP blocks in Figure 4.3b, we see that the size of Δ_{cofa} is almost always larger in the attention blocks than in MLP blocks. This is consistent with the work of Geva et al. (2023) showing that the attention blocks adds most of the information about t_{cofa} to the residual stream.

Finding 2: Only late but not early layers contribute to the competition of the mechanisms. We find that the early layers have almost no contribution to the competition of the mechanisms, reflected by the close-to-zero margin Δ_{cofa} in Layer 0-4 for both types of blocks. However, later layers contribute to substantially to the increase of the margin Δ_{cofa} , by a relatively smaller rate for the MLP blocks, and a larger overall rate for the attention blocks, together with a large variance.

Note that we observe a negative Δ_{cofa} around -0.8 in the last attention block, somewhat favoring t_{fact} , which might be since the factual information is moved to the last position in the last layers, as already noted by Geva et al. 2023.

4.6.3 Microscopic Inspection of Individual Attention Heads

Beyond the overall contributions of the attention block, we further study the contribution of each individual attention head in this section.

Method. We analyze the effect of each individual attention head with the logit inspection method by projecting the outputs of each attention head to the last sequence position N in the vocabulary space. Formally, we consider $\Delta_{\text{cofa}} = \text{HeadLogit}(t_{\text{cofa}}) - \text{HeadLogit}(t_{\text{fact}})$ with the logits from the projection $W_U a_N^{\text{hl}}$ of each head h . Here a_N is

the output of the attention head h after it has been processed by the output matrix of the Attention Block but before its sum to the residual stream.

Results. We plot the contributions of individual attention heads to Δ_{cofa} in Figure 4.4, and introduce the main findings as follows.

Finding 1: A few specialized attention heads contribute the most to the competition. As we can see from the overall contributions of all attention heads across all the layers in Figure 4.4a, several attention heads (e.g., L9H6, L9H9, L10H0, and L10H10) strongly promote the counterfactual mechanism, i.e., with a positive value of Δ_{cofa} colored in dark red, and two attention heads (L10H7 and L11H10) strongly support the factual mechanism instead, reflected by the large negative Δ_{cofa} in dark blue.

For example, the sum of L7H2 and L7H10 equals 75% of the large positive Δ_{cofa} contribution of Layer 7. The sum of L9H6 and L9H9 explains 65% of the Δ_{cofa} at Layer 9.

On the other hand, the two attention heads, L10H7 and L11H10, explain almost the 70% of the total negative contribution to Δ_{cofa} in the entire network (33% and 37% respectively). This also explains the reason behind the negative Δ_{cofa} in Figure 4.3a of the previous section. Our study is consistent with McDougall et al. (2023) showing that these two heads are responsible for suppressing the copy mechanisms in GPT-2 small. In our setting, the joint ablation of these two heads decreases the factual recall of GPT-2 small from 4.13% to 0.65%.

Finding 2: All the highly activated heads attend to the same position – the attribute token. Focusing on the heads with large absolute values of Δ_{cofa} , we show the attention scores of the last position N to different tokens in Figure 4.4b. Expectedly, the major heads supporting the counterfactual mechanism (those in red) attend to the attribute position because they need to copy this token for the prediction, which echoes the findings in Section 4.6.1.

However, it is surprising to see the other heads supporting the factual mechanism (those in blue) also mainly attend to the counterfactual attribute token. We find that those heads read from the attribute position to give a lower value to the logit of t_{cofa} , which might be an easier operation for it to learn than increasing the logit of the factual token. The evidence is that, in these two heads, the logit of t_{fact} is smaller than the mean of the two layers, but the logit of t_{cofa} (which is -1.13 for L10H7 and -1.05 for L11H10) are the lowest of all the heads in the network.

We include supplementary analyses showing the consistency of Finding 2 on Pythia in Appendix A.3.1, and provide the full attention maps with attention scores between every pair of tokens for these heads in Appendix Appendix A.3.2.2.

4.6.4 Intrinsic Intervention by Attention Modification

After *understanding* where the two mechanisms take place, we use the insights to *intervene* on the model internal states. Specifically, we perform model editing to alter the factual mechanism, which concentrates on a few strongly activated attention heads (L10H7 and L11H10 in GPT-2, and mostly L17H28, L20H18, and L21H8 in Pythia, see Appendix A.3.1.3), and has most of the information flowing from the attribute position (see Figure 4.4-right

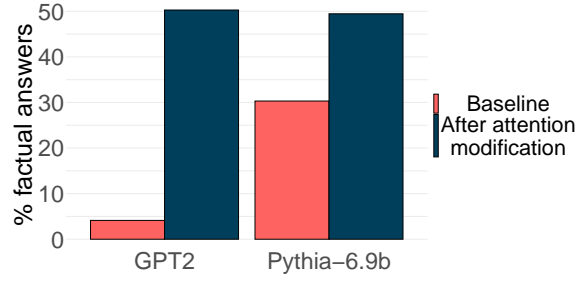


Figure 4.5: The factual recall mechanism increases substantially across GPT-2 and Pythia after attention modification.

and Section 4.6.2). In the following, we show that enlarging the value of a few well-localized attention values can largely improve the factual recall of the model.

Method. We utilize the attention modification method in Eq. (4.5) to apply a multiplier of α to the attention weights of the last token to the attribute position in L10H7 and L11H10 for GPT-2, and L17H28, L20H18, and L21H8 in Pythia. To choose the α value, we perform a grid search over $[2, 5, 10, 100]$ to maximize the factual recall rate of the model. We find that $\alpha = 5$ is the best value for both GPT-2 for Pythia.

Results. We highlight the effect of our model editing method on the strength of the factual recall mechanism in Figure 4.5. Originally, GPT-2 has only 4% of the cases where the factual mechanism prevails the counterfactual one, and Pythia only 30%. However, after modifying the attention weights of the entries mentioned above, the strength of the factual mechanism increases drastically that it wins over the other mechanism in 50% of the cases for both models. This result is remarkable since we modify only two entries in the attention map out of the 33,264 attention values of GPT-2 (117M parameters) and three entries out of the 270,848 attention values of Pythia (6.9B parameters). This highlights the importance of the interpretability analysis in Sections 4.6.2 and 4.6.3, which enables us to find the detailed role played by the individual units of each transformer layer.

4.6.5 What Word Choices Intensify the Competition?

After the intrinsic intervention to edit the internal states of the model, we explore how the similarity between t_{fact} and t_{cofa} in our dataset affects the mechanism described in the previous sections.

Method. We divided the dataset into 10 equal bins based on the similarity between the vectors for t_{fact} and t_{cofa} , with each bin containing 1000 items. Starting from the lowest, each group represents a 10% segment of the dataset, arranged by increasing similarity scores. For our word similarity metric, we calculate the cosine similarity of the 300-dimensional word embeddings from the pre-trained Word2Vec model (Mikolov et al., 2013) implemented in the Gensim Python package (Řehůřek and Sojka, 2010).

Results. As a result of the varying similarity of the two tokens, we see a drastic change in the dominance of the factual mechanism in Figure 4.6.

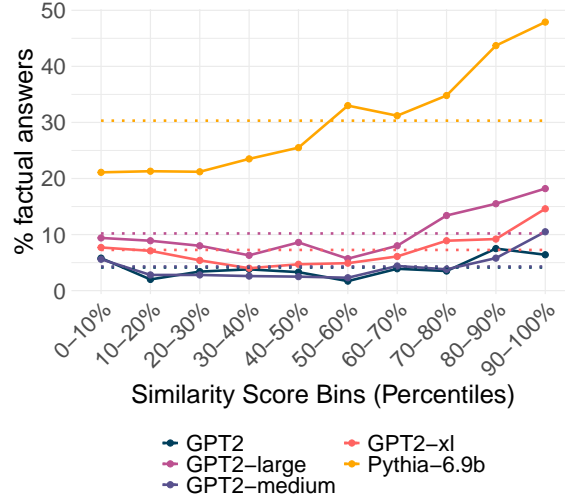


Figure 4.6: Prediction frequency of factual token by similarity level. We show the percentage of t_{fact} predictions within each bin compared to the entire dataset (represented by a dotted line) across various model sizes. We can notice that more similar t_{fact} and t_{cofa} are, and the factual mechanism is stronger.

Finding 1: Similar tokens confuse the model more easily. Consistently across all the models, the more similar the two tokens are, the more likely the model is to be confused and mistakenly let the factual mechanism dominate, to predict t_{fact} as its output.

Finding 2: Larger models suffer more from such confusion. For example, the largest one, Pythia-6.9B, demonstrates a very strong attachment to the factual information, letting the factual mechanism win almost 45% of the cases when the token similarity reaches 90%. Even when the similarity is low, larger models are still more likely to confuse and lean towards the factual mechanism. This finding resonates with the observations from the inverse scaling prize McKenzie et al. (2023) that larger models have a greater capacity to store and retrieve factual information, thus more influenced by the factual mechanism.

4.7 Discussion and Future Work

Situating Our Findings in Related Work. Our findings about the late attention blocks are consistent with Geva et al. (2023), showing that late attention blocks write most of the information to the last layer when adding a counterfactual premise. Surprisingly, however, we find that the largest contribution to the factual prediction of the network mostly comes from the suppression of the counterfactual token read from the attribute position rather than the promotion of the factual token from the subject position.

Consistently with McDougall et al. (2023), we find that few highly specialized heads suppress the counterfactual information. Moreover, we make a unique contribution up-weighting only two or three attention entries of these heads to increase substantially the number of factual responses of the model.

With an approach similar to ours, Yu et al. (2023) find that more heads can promote the factual mechanism, also in early layers, but found it challenging to improve the factual responses by scaling up the weights of the attention maps. This discrepancy can be due

to the broader set of topics we include in our prompts, which allowed us to select fewer, more specialized heads, to the different ways the prompts are framed, or also to our more focused modification of the attention maps.

Future Work For future research directions, we aim to analyze more in depth how our findings depend on the prompt structure and whether the promotion of factual responses by suppressing the counterfactuals generalizes to larger models and a more comprehensive variety of datasets.

4.8 Conclusion

In this work, we have proposed the formulation of the *competition of mechanisms* as a powerful interpretation when LLMs need to handle multiple mechanisms, only one of which leads to the correct answer. We deployed two mechanistic interpretability tools, logit inspection and attention modification, and identified critical positions and model components involved in competing for the mechanisms. Finally, we discovered a few localized positions in the attention map, which largely control the strength of the factual mechanism. Our study sheds light on future work on interpretability research for LLMs.

Limitations

Limited models: Our study aligns with most existing work in mechanistic interpretability to use GPT-2 small. However, we understand that this is a small model with far fewer parameters than current state-of-the-art LLMs. Future work is welcome to extend to larger-sized models, which might generalize our conclusion to a certain extent, and also reveal interesting behavior once the models get beyond a specific size, maybe also seeing a U-shaped curve Wei et al. (2023) for the dominance of the counterfactual mechanism.

Interpretability method: Furthermore, our experiments and insights are heavily grounded in the interpretability within the embedding space of the model’s inner components. This approach is reliable and extensively employed in mechanistic interpretability research (Dar et al., 2023; Geva et al., 2022; Halawi et al., 2023). The logit inspection method, although commonly employed in previous work, can occasionally fail to reflect the actual importance of some vocabulary items, especially in the early layers of the network (Belrose et al., 2023).

Simplicity of the prompts: Our prompts have a relatively simple structure for the controllability of the counterfactual information tracing, as it is very challenging to follow the information flow in a more diversified set of prompts. We welcome future work to explore methodological advances to enable analyses over more diverse prompts.

Ethical Considerations

The aim of our study is to enhance comprehension of the interplay among mechanisms within language models that may yield unforeseen and undesirable outcomes. Additionally, our research serves as a conceptual demonstration of methods to guide model behavior under such conditions. We believe that recognizing and dissecting the mechanisms by which LLMs produce unpredictable responses is crucial for mitigating biases

and unwanted results. Moreover, understanding the competitive dynamics under investigation is critical for improving the safety of LLMs. Specifically, inputting a prompt with an inaccurate redefinition may lead the model to inadvertently reveal sensitive factual information.

A Causal Framework to Quantify the Robustness of Mathematical Reasoning in LLMs

We have recently witnessed a number of impressive results on hard mathematical reasoning problems with language models. At the same time, the robustness of these models has also been called into question; recent works have shown that models can rely on shallow patterns in the problem description when generating a solution. Building on the idea of behavioral testing, we propose a novel framework, which pins down the causal effect of various factors in the input, e.g., the surface form of the problem text, the operands, and math operators on the output solution. By grounding the behavioral analysis in a causal graph describing an intuitive reasoning process, we study the behavior of language models in terms of robustness and sensitivity to direct interventions in the input space. We apply our framework on a test bed of math word problems. Our analysis shows that robustness does not appear to continuously improve as a function of size, but the GPT-3 Davinci models (175B) achieve a dramatic improvement in both robustness and sensitivity compared to all other GPT variants. Our code and data are available at <https://github.com/alestolfo/causal-math>.

5.1 Introduction

Many natural language understanding situations, such as understanding the financial news, require reasoning with text that includes numbers. However, such mathematical reasoning is challenging for NLP models (Cobbe et al., 2021; Mishra et al., 2022b). Mathematical reasoning for text has been an active area of research for a while (Seo et al., 2015; Sachan and Xing, 2017; Sachan et al., 2017, 2018, *inter alia*), and has also emerged as a key task to track the capabilities of large language models (LLMs) in recent years (Brown et al., 2020; Ouyang et al., 2022; Wei et al., 2022a, *inter alia*).

However, despite the impressive performance of LLMs on various math reasoning benchmarks (e.g., Ouyang et al., 2022; Chowdhery et al., 2022), it remains unclear whether these models have learned mere artifacts in the data or have truly mastered the mathematical concepts needed to consistently solve all variations of the same problem (Patel et al., 2021; Razeghi et al., 2022; Welleck et al., 2022). In sharp contrast with a large number of papers on improving the performance of LLMs on various types of math-based problems, there has been little effort on behavioral analysis of LLMs for these tasks. Existing methods for understanding the robustness of these models (Patel et al., 2021) rely on manually constructing variations of math problems, and we do not yet have a principled, comprehensive framework for quantifying such robustness.

Thus, in this work, we propose a formal framework based on causal inference, to quantify

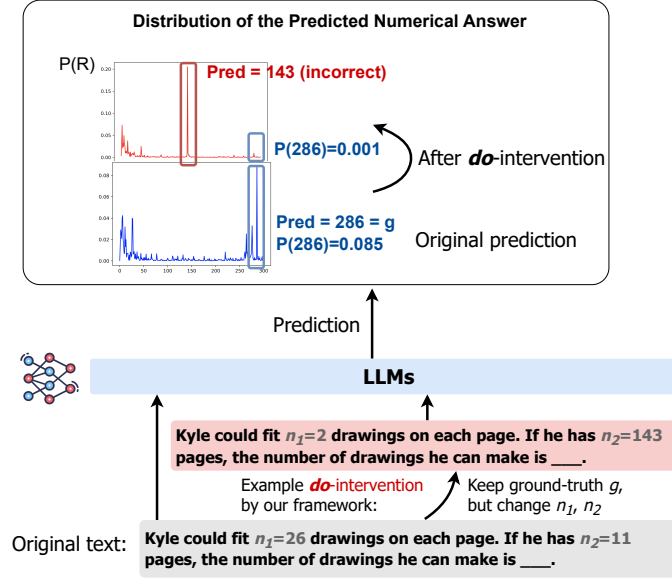


Figure 5.1: Through our framework, we conduct do-interventions on the input and evaluate the change in the distribution $\mathbb{P}(R)$ of the prediction R by LLMs, in this figure, GPT-J. This allows us to measure the causal effect of each factor in the input on the model’s response.

the robustness of NLP models’ math reasoning abilities. Specifically, we describe a causal graph formulation of math reasoning, where the graph allows us to measure the difference in the structural causal models of human reasoning and model judgment. We consider various causal factors such as the textual framing of the question, numerical operands, and operation types. Then, we identify a set of interventions in the context of math word problems (an example of which is illustrated in Figure 5.1), and provide a causal inference framework to obtain causal effects of each factor via direct do-interventions (Pearl, 1995) and causal mediation analysis (Pearl, 2001). While our approach is reminiscent of recent studies using causal analysis for LLMs (Finlayson et al., 2021; Vig et al., 2020b; Meng et al., 2022), in this work, we provide a new theoretical analysis framework specifically suitable for math reasoning. Using our framework, we disentangle factors affecting the model’s predictions and measure their influences. This way, we are able to provide insights into the model’s reasoning in terms of *robustness* and *sensitivity* with respect to changes in these factors.

We apply our framework to study a set of thirteen GPT models with various sizes and training procedures (i.e., instruction-tuned and non-instruction-tuned). We observe that, among non-instruction-tuned language models, the larger ones tend to be more sensitive to changes in the ground-truth result of a math word problem, but not necessarily more robust. However, we observe a different behavior in the instruction-tuned GPT-3 models (Ouyang et al., 2022), which show a remarkable improvement in both sensitivity and robustness, although the robustness reduces when problems get more complicated. We additionally investigate the role of size and instruction tuning on the model’s performance with three models of the LLaMA family (Touvron et al., 2023) and Stanford Alpaca (Taori et al., 2023a).

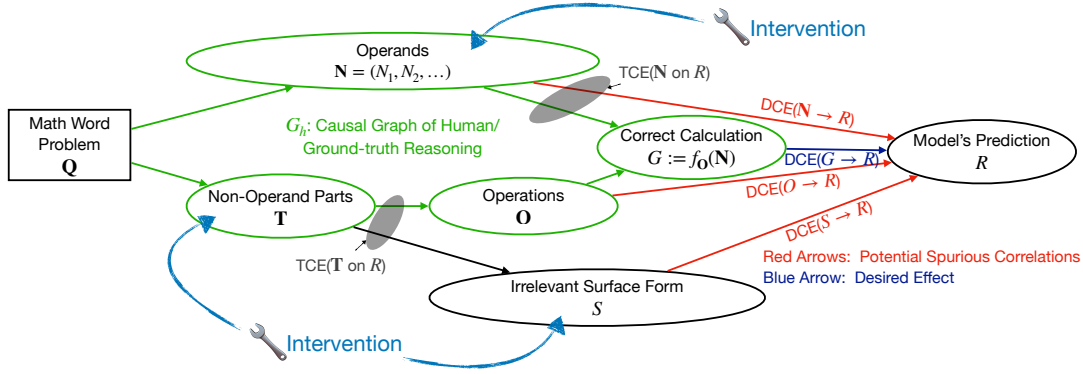


Figure 5.2: Causal graph of model predictions on math questions. We highlight the difference between a cognitively-inspired correct reasoning path (\mathcal{G}_h) and the undesired effects that some factors might have on the model’s prediction (red arrows). By performing controlled interventions of the numerical values (N) and on the textual framing of the problem (T, S), we are able to quantify the causal effects of each factor.

5.2 Problem Setup

We consider a dataset \mathcal{D} of math word problems (MWP), where each MWP is denoted as a question Q . Q is a list (T, N) consisting of a question template T and an ordered list of operands $N = (N_1, N_2, \dots, N_m)$. Each question template $T := (O, S)$ further contains two types of information: a set of arithmetic operations O implicitly expressed in the question, and the text surface form S irrelevant to the arithmetic operations. O incorporates the information relative to the operations as a collection of tuples $\{(O_1, i_1, j_1), (O_2, i_2, j_2), \dots\}$, where $O_k \in \{+, -, \times, \div\}$ ($k \in \mathbb{N}$) and $i_k, j_k \in \mathbb{N}$ represent the indices of the operands to which operator O_k should be applied.¹ The ground-truth result $G = f_O(N)$ is calculated by computing the function f_O , which represents the application of all the operators in O to the respective operands. We illustrate the factors in Q and their inter-dependency in the causal graph in Figure 5.2. A two-operand instance q of Q in this form from Patel et al. (2021) is:

Template t : Mark has n_1 trees in his backyard. If he plants n_2 more, how many trees will he have?

Operands n : ($n_1 = 12, n_2 = 13$)

Operations o : $\{("+", 1, 2)\}$

Result: $g = f_o(n) = n_1 + n_2 = 25$

Our goal is to quantify the robustness of a model \mathcal{M} on the set of problems $q \in \mathcal{D}$. Ideally, \mathcal{D} should be a dataset not seen by the model during training. We assume that a model takes q as input and predicts a probability distribution of the result R : $\mathbb{P}(R | t, n)$. Our formulation below will be easier to understand using this finite discrete set and can be generalized to any kind of data pairing a natural language template with a function that maps a set of operands to a result (e.g., a Python program; Mishra et al. 2022a).

¹The intermediate result of operation O_1 is indicated by $i_k = m + 1$.

5.3 A Causal Framework

In this section, we describe our framework in three steps. First, we define the idea of model robustness on MWPs. Then, we identify possible do-interventions (Pearl, 1995) that we can perform. Finally, we describe the causal effects that we measure to quantify the robustness of various models.

5.3.1 Step 1. Question Reformulation

We address the research question “*Is a model reasoning robustly on MWPs?*” by comparing the causal mechanisms of the model’s decisions to a hypothesized human reasoning mechanism. Note that we do not claim to know how humans reason about these problems. We simply propose a reasonable and intuitive way to judge model robustness given a reasonable and intuitive human reasoning mechanism inspired by findings regarding the independence of language and mathematical reasoning in humans (Brannon, 2005; Monti et al., 2012).

Human Reasoning Mechanisms. The causal mechanisms of how humans might solve q include

$$o = f_{\text{abstract}}(q), \quad (5.1)$$

$$g = f_o(n), \quad (5.2)$$

where they first abstract the arithmetic operations o from the problem q by some cognitive process f_{abstract} , and then apply the operation to the operands to obtain the result g . We show these mechanisms in the green subgraph \mathcal{G}_h of Figure 5.2.

Model Reasoning Mechanisms. In contrast, the causal mechanisms of how a model might solve q are as follows:

$$r = f_{\text{blackBox}}(t, n), \quad (5.3)$$

where we are unsure about (1) *what* part(s) of t the model takes into account, and (2) *how* it operates over the relevant variables.

Thus, we draw all possible causal mechanisms that might take place in the black-box model f_{blackBox} in the complete causal graph in Figure 5.2. Some possible fine-grained causal mechanisms are

1. The model might attend over the question template t in two ways: paying attention to the text surface form s via the causal path $T \rightarrow S \rightarrow R$, or text relevant to the math operations o via the causal path $T \rightarrow O \rightarrow R$.
2. The model might also attend to the operands $n := (n_1, n_2, \dots)$ via a causal path $N \rightarrow R$.
3. If the model learns the correct causal mechanisms as in the human cognitive process, it should capture how the operator and the operands matter to the ground-truth result g (via $O \rightarrow G$ and $N \rightarrow G$) and then the model prediction should be sensitive to any changes in the ground truth, namely $G \rightarrow R$. No spurious correlations can directly affect R without going through the mediator G .

Hence, to answer the question “How robust is the mathematical reasoning of a model on MWP?” we can answer the following subquestions:

1. How does R change in response to G ? By quantifying this, we assess the *sensitivity* (correct responsiveness) of the model to changes in the problem. In other words, does the model correctly adjust its prediction in response to a change in the correct solution of the problem?
2. What is the (unwanted) direct causal effect size of $S \rightarrow R$, and $N \rightarrow R$? We see the quantities as a measure of the *brittleness* (i.e., wrong responsiveness) of the model to result-preserving changes in the input. The lower the direct causal effect of S and N , the more *robust* the model is.

5.3.2 Step 2. Causal Intervention List

After formulating the cognitively-inspired subgraph \mathcal{G}_h and defining the undesired causal paths in Figure 5.2, we list all feasible limited actions that allow us to perform our causal analysis. In the context of MWPs, we use the following interventions:

1. Direct intervention on all possible n_1, n_2, \dots ;
2. Partially controllable interventions on T . We can replace the template T in two ways:
 - (a) both S and O are affected, or
 - (b) S is affected but O is not affected.

5.3.3 Step 3. Turning Limited Actions into Causal Effect Sizes

Next, we explain how we can obtain the causal effect sizes we want (listed in Step 1) from the limited set of interventions we can do (listed in Step 2). Specifically, we first start from all the feasible interventions, and for variables that we cannot directly intervene on, we apply deductions from do-calculus (Pearl, 1995) to obtain or approximate the direct causal effect sizes. In the following, we describe a list of causal effect sizes that we need.

General Formulation. Let us consider an intervention $\text{do}(X : x \rightarrow x')$, where $X \in \{T, S, N\}$ and a problem $Q = \{T, N\}$. The support of the numerical values N_i 's and R is $\mathcal{I} \subseteq \mathbb{N}$, and we consider N to be distributed uniformly over the set $\{n \in \mathcal{I}^2 \mid f_O(n) \in \mathcal{I}\}$. We denote the distribution before the intervention $\mathbb{P}(R \mid T, N)$ as P and the distribution after the intervention as P' .

Following the distributional definition of causal effect by Pearl (1995), we quantify the effect of factor X in our causal graph using a distance metric δ between the distributions P and P' . That is,

$$\text{CE} = \delta(P, P'), \quad (5.4)$$

where CE can refer to the **total causal effect** (TCE, i.e., the joint effect through all the directed causal paths from a variable to another), or the **direct causal effect** (DCE, i.e., the effect from the directed causal path from a variable to another that does not go through any intermediate variables) (Pearl, 2001). We describe our choices for δ in Section 5.3.4.

Causal Effects of the Operands. When intervening on the operands $N := (N_1, N_2, \dots)$, we can obtain the size of the total causal effect of N on R , namely

$$\text{TCE}(N \text{ on } R) := \mathbb{E}_{\mathbf{n}' \sim \mathbb{P}(N)}[\delta(P, P')], \quad (5.5)$$

$$\text{where } P' = \mathbb{P}(R|T, \text{do}(N = \mathbf{n}')) . \quad (5.6)$$

Note that this TCE is not the exact desired quantity, because we want to separate two different paths of how N affects R : (1) the path $N \rightarrow G \rightarrow R$, which is the correct decision path that we want the model to pick up (where the model reacts to the change in the ground-truth answer), and (2) the path $N \rightarrow R$, which is the spurious correlation that the model might have learned (where the model relies on some spurious correlations with certain numerical values, which could be traced to perhaps their frequencies in the training corpus).

We can quantify the **direct causal effect** (DCE, i.e., the effect from the directed causal path from a variable to another that does not go through any intermediate variables) (Pearl, 2001) of N on R , namely the strength of the direct causal path $N \rightarrow R$, by controlling for G to be fixed every time we intervene on N :

$$\text{DCE}(N \rightarrow R) := \mathbb{E}_{\mathbf{n}' \sim \mathbb{P}(N|G)}[\delta(P, P')], \quad (5.7)$$

$$\text{where } P' = \mathbb{P}(R|T, \text{do}(N = \mathbf{n}')) . \quad (5.8)$$

For example, if we observe a model doing $100 + 100 = 200$ correctly, we want to separate the math ability here into (1) the model's sensitivity towards the ground-truth answer, and (2) the model's decisions based on its familiarity with just the operand 100. Here, the overall effect is the calculable $\text{TCE}(N \text{ on } R)$ by Eq. (5.5), and one of the subeffects is the calculable $\text{DCE}(N \rightarrow R)$ by Eq. (5.7).

Causal Effects of the Text Surface Form. As for the operands, we can compute both the direct and indirect effects of the surface form representing the math problem. In particular, intervening on T without controlling for O (intervention 2a in Section 5.3.2), we can compute the total effect, i.e.,

$$\text{TCE}(T \text{ on } R) := \mathbb{E}_{\mathbf{t}' \sim \mathbb{P}(T)}[\delta(P, P')], \quad (5.9)$$

$$\text{where } P' = \mathbb{P}(R|N, \text{do}(T = \mathbf{t}')) . \quad (5.10)$$

Controlling for the operations O (intervention 2b in Section 5.3.2) will instead allow us to obtain the direct causal effect of the surface text:

$$\text{DCE}(S \rightarrow R) := \mathbb{E}_{\mathbf{t}' \sim \mathbb{P}(T|O)}[\delta(P, P')], \quad (5.11)$$

$$\text{where } P' = \mathbb{P}(R|N, \text{do}(T = \mathbf{t}')) . \quad (5.12)$$

Note that since there is no mediator between S and R , the $\text{DCE}(S \rightarrow R)$ is also TCE of S on R . The only adaptation that we need to make with regard to the MWP is that it is not feasible to enumerate all possible perturbations of S . Therefore, the practical results that researchers can achieve are over a certain subset of S . In practice, we obtain this by intervening on T without affecting O .

Causal Effects of the Operators. The ideal way to obtain the TCE of O on R is through some careful human annotation that minimally changes the templates as Kaushik et al. (2020) do for sentiment classification. The challenge for MWP in our case is that with all our possible interventions, we cannot *only* intervene on O without introducing changes to the irrelevant surface form. However, we might get some information about $\text{TCE}(O \text{ on } R)$ because, on the causal graph, the total causal influence of T on R actually flows into two directed paths, one through S to R (which is the $\text{DCE}(S \rightarrow R)$), and the other from O to R , which is our interested quantity $\text{TCE}(O \text{ on } R)$. Therefore, we compare the two quantities we know, $\text{TCE}(T \rightarrow R)$ and $\text{DCE}(S \rightarrow R)$, to get a sense of the causal influence of O on R that we cannot obtain in any other way.

5.3.4 Step 4. Quantifying the Causal Influence

Consider a realization of problem Q with operands n and ground-truth result $g = f_o(n)$, and denote by g' the result after the intervention $\text{do}(X : x \rightarrow x')$. We quantify the causal effect of factor X on the model's prediction R in two ways: by assessing the change in the predicted result, and by measuring the change in the probability assigned by the model to the correct result g (or g').

Change in the Prediction To account for the inability of LMs to capture the continuous property of numbers (Jin et al., 2021a), we measure the change in the model's prediction using an indicator of the "change result" event:

$$\delta_{\text{cp}}(P, P') := \mathbf{1}_{r \neq r'}, \quad (5.13)$$

where $r = \arg\max_{x \in \mathcal{I}} P(x)$, and $r' = \arg\max_{x \in \mathcal{I}} P'(x)$.

Relative Change in Confidence Inspired by Finlayson et al. (2021), we also highlight the change in terms of the relative difference in the probability assigned to g and g' . We formulate two types of relative change, one quantifying the relative change in the confidence of g , and the other quantifying the relative change in the confidence of g' :

$$\Delta_{\text{rel}} = \frac{P(g) - P'(g)}{P'(g)} \quad (5.14)$$

$$\Delta'_{\text{rel}} = \frac{P'(g') - P(g')}{P(g')}. \quad (5.15)$$

We quantify the overall relative change in confidence (RCC) as the average of the two relative changes above:

$$\delta_{\text{rcc}}(P, P') = \frac{1}{2} \left(\Delta_{\text{rel}} + \Delta'_{\text{rel}} \right). \quad (5.16)$$

A Unified Form We are interested in the average causal effect of the intervention across all problems in \mathcal{D} . Thus, we measure the average of the effects over all instances $q \in \mathcal{D}$. We denote by the subscripts $\text{TCE}_{\text{cp}}/\text{DCE}_{\text{cp}}$ and $\text{TCE}_{\text{rcc}}/\text{DCE}_{\text{rcc}}$ the causal effects computed using the change in prediction metric and the relative change in confidence, respectively. We describe how we construct the dataset \mathcal{D} in Section 5.4.2.

5.4 Experimental Setup

In this section, we describe the data used to perform the interventions and to measure the causal effects.

5.4.1 Datasets

For our analyses, we use instances of math word problems from three popular datasets: ASDiv-A (Miao et al., 2020), MAWPS (Koncel-Kedziorski et al., 2016), and SVAMP (Patel et al., 2021). The examples contained in these collections are pairs (t, o) consisting of a question template t with its annotated operations o . Each of these pairs can be instantiated multiple times into problems $q = (t, n)$ by filling the template with numerical values (n_1, n_2, \dots) and computing the ground-truth result $g = f_o(n)$ (most problems involve two to three operands, i.e., $|n| \in \{2, 3\}$). We select a set of 437 two-operand and 307 three-operand template-expression pairs that we use to generate pairs of prompts representing an intervention. More details about the prompt generation procedure are in Appendix A.4.1. We use (t, n) to refer to an instantiated template that we use as a prompt.

5.4.2 Intervention Data

Given an MWP $q = (t, n)$ and its solution g , we generate a second problem-solution instance (q', g') depending on the type of causal effect CE we want to measure and on the considered variable. When intervening on the operands of the problem, the text of the problem is kept unaltered and a set of new operands n is sampled in such a way that the result g is affected or not depending on the effect that is being measured. When changing the textual description of the problem, we change t such that either $o' = o$, or $o' \neq o$. In the former case, we sample a different template $t' = (s', o)$ from the set of templates describing the same operations o , in the latter case we sample a new t' describing a different operation. In Appendix A.4.2.1 we report some examples of (q, q') pairs representing the different types of interventions.

Given a model, we use the question pair (q, q') to obtain a pair of answer distributions $\mathbb{P}(R|t, n)$ and $\mathbb{P}(R|t', n')$, which we use to measure the causal effect of the intervention. We consider the space for the numerical values to be $\mathcal{I} = \{1, 2, \dots, C\}$ consisting of integer values, following the setup of several existing MWP datasets (Miao et al., 2020; Koncel-Kedziorski et al., 2016; Patel et al., 2021). To control our experimental costs and make sure the models keep the number as one token, we set $C = 300$. From all the tokens in a model’s vocabulary, we focus on the probability assigned to the numbers in our numerical space \mathcal{I} , and thus we use $\mathbb{P}(R = r)$ to denote the normalized probability $\mathbb{P}_{\text{raw}}(R = r)/Z$, where $Z = \sum_{r=1}^C \mathbb{P}_{\text{raw}}(R = r)$, and $\mathbb{P}_{\text{raw}}(x)$ is the raw probability score assigned to the vocabulary token x . For each intervention type, we generate a dataset \mathcal{D} consisting of (q, q') pairs. Unless otherwise specified, for our experiments we generate 500 intervention pairs for each template, and results are averaged over three seeds.

5.4.3 Models to Evaluate

We use our framework to assess the robustness of reasoning in thirteen pre-trained language models. We consider five sizes of the GPT-2 model (Radford et al., 2019): distilled (Sanh et al., 2019), small, medium, large, and XL. We evaluate four models from EleutherAI that were pre-trained on the Pile (Gao et al., 2020): GPT-Neo 1.3B and 2.7B

(Black et al., 2021), GPT-J-6B (Wang and Komatsuzaki, 2021), and GPT-NeoX-20B (Black et al., 2022). We use HuggingFace Transformers (Wolf et al., 2020) to access the models. Additionally, we experiment with a set of instruction-tuned versions of GPT-3 (Brown et al., 2020): Instruct (Ouyang et al., 2022), Curie, Davinci-002, and Davinci-003.² Experiments with GPT-3 are carried out under the constraints set by the OpenAI APIs³, which prevent us from computing the causal effect using the same procedure as for the other models. We report the details about how the metrics were computed for GPT-3 in Appendix A.4.3. In the reported results, we indicate with an asterisk (*) the metrics that were influenced by this limitation.

5.5 Results

Our analyses focus primarily on two-operand problems (Sections 5.5.1 and 5.5.2) and later extend to more complex problems that involve three operands (Section 5.5.5) for the models that perform best on the two-operand test bed. We compare the direct causal effect DCE and the total causal effect TCE of N and T on R . DCE represents the undesired effect for a model to being mistakenly responsive to a change in N or T not leading to a change in the result g (low robustness), whereas higher values of TCE indicate a higher ability of the model to correctly adjust the probability weight assigned to the new solution g' after the intervention (high sensitivity).

5.5.1 Effect of N on R

From the results in Figure 5.3, we notice that larger models exhibit a larger TCE_{rcc}/DCE_{rcc} ratio. In particular, in GPT-J-6B and NeoX, the TCE is, respectively, 30x and 1000x larger than the DCE. However, this improvement in sensitivity is not manifested in terms of change of prediction (δ_{cp}), for which the models show to be affected by result-preserving changes almost as equally as by result-altering interventions. This behavior changes significantly in instruction-tuned models. In particular, for the 175B-parameter GPT-3, performance varies depending on the type of supervision, with the PPO-trained Davinci-003 exhibiting an 84% difference between direct and total effect.

In Figure 5.4, we present a different visualization of the direct causal effect of N on the model’s prediction. We report the heatmaps showing the probability assigned by the model to the result g of a problem $(t, (n_1, n_2), g) \mid g = n_1 + n_2, \forall g \in \{0, 1, \dots, 50\}, \forall (n_1, n_2) \in \{0, 1, \dots, 50\}^2$. For Distil-GPT-2 we observe low overall probability assigned to g and diagonal patterns indicating consistency in assigning higher probability to specific results (e.g., 10, 20, 30, 40, 50). For the two larger models we notice a higher probability mass assigned to the problem’s result, but less consistency on the prediction of the same result with different sets of operands (this is true for GPT-J in particular). This result is consistent with the observed higher DCE and TCE in larger models: $P(g)$ might vary more considerably when intervening on N without affecting g , but overall the model assigns higher probability weight to the correct result, which correlates with higher sensitivity.

²The OpenAI ids for these models are, respectively, `davinci-instruct-beta`, `text-curie-001`, `text-davinci-002`, and `text-davinci-003`.

³<https://openai.com/api/>

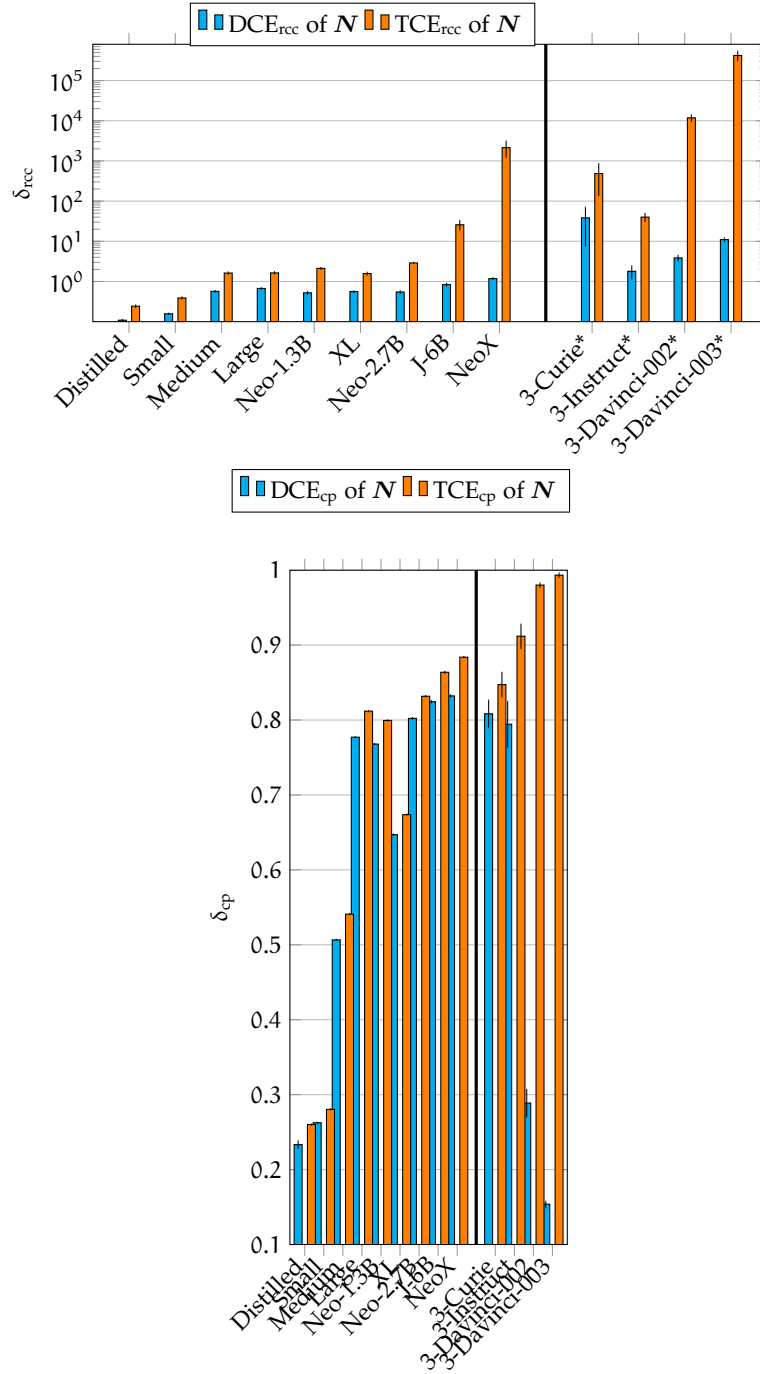


Figure 5.3: Comparison of $\text{DCE}(N \rightarrow R)$ and $\text{TCE}(N \text{ on } R)$. *approx values, see Appendix A.4.3.

5.5.2 Effect of T on R

In Figure 5.5, we report the total causal effect of the textual framing T and the direct causal effect of the irrelevant text elements S on the model’s prediction. For the instruction-tuned models, the improvement in terms of prediction change (δ_{cp}) follows a similar trend as for N , with GPT-3 Davinci-003 showing a 76% difference between direct and total effect. An interesting observation is that the irrelevant textual information S appears to have a lower direct effect than N for all non-instruction-tuned models. However, in the GPT-3 Davinci-00x models, we observe the opposite (i.e., $DCE(N \rightarrow R) \leq DCE(S \rightarrow R)$). This suggests that large instruction-based models tend to be more susceptible to variation in the textual framing of a problem, while smaller models are more responsive to changes in the numerical values (though not necessarily correctly).

5.5.3 Overall Insights

In comparison to other models, GPT-3 Davinci shows the highest DCE_{rcc} , but low DCE_{cp} . This discrepancy is related to the quantities that the two metrics consider. δ_{rcc} takes into account the probability assigned to g , while δ_{cp} does not consider the ground truth solution. One interpretation of this result is that GPT-3 Davinci consistently predicts the same answer $r = r'$ when $g = g'$, but the probabilities $P(g)$ and $P'(g)$ might vary significantly.

The results observed for the two kinds of intervention $do(T : t \rightarrow t')$ and $do(N : (n_1, n_2) \rightarrow (n'_1, n'_2))$ show similar trends. Small models (Distilled and Small GPT-2) exhibit low sensitivity to interventions. Larger models (from GPT-2 Medium to GPT-Neo) appear to be more influenced by changes in both N and T . However, they display similar sensitivity to both result-altering and result-preserving interventions. An improvement in sensitivity is noticeable in GPT-J and NeoX, though not accompanied by an improvement in robustness. Remarkably different behavior is instead shown by the GPT-3 Davinci models, which demonstrate substantially higher sensitivity to result-altering interventions (high TCE), and higher robustness (in terms of prediction change). In Appendix A.4.2.2, we report the accuracy of the models on the generated instances of MWPs, which exhibits a similar trend as the robustness/sensitivity changes we observed.

Possible explanations for the improved robustness and sensitivity demonstrated by the large GPT-3 models might be the dramatic size increase and extension/enhancement of the training procedure involving instructions. The former idea is aligned with the *emergent abilities* hypothesis (Wei et al., 2022a), which postulates the existence of skills that are displayed by large-scale models but are not present in smaller-scale models. However, our observations show different performances in versions of GPT-3 Davinci that differ in the training procedure.⁴ This raises the question of whether the capability of LLMs to reason about math problems benefits from instruction-based tuning. We address this question in the following section.

5.5.4 Extending to LLaMA-Based Models

To further investigate the roles played by size and training method in the model’s performance, we carry out our experimental procedure on three versions with different sizes (7B, 13B, and 30B) of the LLaMA model (Touvron et al., 2023), and on Stanford Alpaca

⁴A high-level description of the training procedures for the models is provided at <https://beta.openai.com/docs/model-index-for-researchers>.

(which applies instruction tuning on LLaMA 7B) (Taori et al., 2023a). We present these results separately, as the LLaMA tokenization makes the prediction setup different from the one used from the other models, and prevents us from computing the relative change in confidence (δ_{rcc}).⁵

From the results (Figure 5.6), two notable observations emerge. Firstly, the increased difference between TCE and DCE observed with the increasing size of the LLaMA models suggests that a larger number of parameters can be a significant driver behind robustness/sensitivity improvement. However, this is not necessarily the case across different models: GPT-NeoX-20B shows a smaller $\text{TCE}_{\text{cp}}\text{-DCE}_{\text{cp}}$ gap compared to LLaMA 7B (5.2% vs 9.0%). Secondly, the instruction tuning procedure of Alpaca does not seem to help significantly with mathematical computation: the decrease in both TCE and DCE shows that robustness improves at the expense of sensitivity. Nonetheless, overall, when comparing Alpaca compared to its base model, LLaMA 7B, we observe an increase in the gap between TCE and DCE, although this difference is minimal (9.5% vs 9.0%).

The limited improvement of Alpaca might be attributed to its instruction tuning procedure consisting of “a list of user-oriented instructions including email writing, social media, and productivity tools” (Taori et al., 2023a), which differs from reasoning-intensive tasks. We suggest future work to examine different types of instruction tuning (e.g., focused on reasoning procedures or reinforcement learning from human feedback), which might help the model answer more complex types of questions in a step-by-step manner and more accurately. We hypothesize that the different performances in versions of GPT-3 Davinci might be produced by the specific type of instructions used for training, by the reinforcement learning component (Ouyang et al., 2022), or simply by an extension of the language modeling pre-training. It is challenging to pinpoint the exact factor in the training procedure that contributes to this improvement, as specific methodological details are not available.

5.5.5 Moving to Three-Operand Problems

We extend our evaluation to consider the three-operand problems in the dataset. In these experiments, we consider only the GPT-3 175B-parameter models, as they are the only models performing well on the simpler bivariate problems. The results regarding the effects of N are reported in Figure 5.7. We notice that the large difference between the desired (TCE) and undesired (DCE) effects observed on simpler problems shrinks significantly for both metrics. In particular, for Davinci-003, the direct effect of N (measured as δ_{cp}) grows from 0.17 to 0.87. That is, GPT-3 Davinci-003 predicts a different result 87% of the time after an intervention that does not affect the ground-truth solution. The increase in direct effect indicates a performance degradation in terms of brittleness: even the models that show good performance on two-operand problems, now display an unstable behavior after result-preserving interventions.

5.6 Related Work

Causal NLP Causal inference aims to study the cause and effect from observational and interventional data (Pearl, 2009b; Peters et al., 2017). Traditionally, researchers usually

⁵The LLaMA tokenizer considers each digit as an independent token in the vocabulary. This makes it problematic to compare the probability value assigned by the model to multi-digit numbers.

apply causal techniques to phenomena in nature and human society. With the rise of powerful models in NLP, recent research has started to explore the intersection of causal inference and NLP, forming the study of Causal NLP (Jin et al., 2022a; Feder et al., 2021a).

There are several formulations for Causal NLP: the *causality for NLP* thread involves using the causal framework for data collection and task formulation (Jin et al., 2021c), inspecting the (path-specific) causal effect of certain neurons on predictions (Vig et al., 2020b; Meng et al., 2022), understanding the causal effect of data and learning paradigm for model performance (Ni et al., 2022), and as a way to frame prompts (Lyu et al., 2024); and *NLP for causality* involves testing the pure causal inference skills of LLMs (Jin et al., 2023a, 2024), and use text as a variable for causal effect estimation (Roberts et al., 2020; Veitch et al., 2020; Jin et al., 2021b, 2023b).

The most similar line of research to our work is the application of causal effect estimation on interpreting models' behavior, such as how models understand syntactic agreement (Finlayson et al., 2021), and how interventions in the representations and weights affect the model prediction (Feder et al., 2021b). To the best of our knowledge, our work is the first to formulate a causal framework for robustness behavioral tests, and also we are the first to introduce the idea to quantify the differences in the causal mechanisms of human reasoning and model decisions.

Math Reasoning in NLP A growing body of work tries to improve the math reasoning capability in NLP models (Zhang et al., 2020b; Geva et al., 2020; Spokoyny et al., 2021), and prompting techniques for LLMs (Cobbe et al., 2021; Shen et al., 2021b; Kojima et al., 2022; Wei et al., 2022b; Chowdhery et al., 2022). For analysis, significant attention has been given to models' ability to understand numerical quantities (Wallace et al., 2019; Thawani et al., 2021) and numerical operations (Pal and Baral, 2021; Berg-Kirkpatrick and Spokoyny, 2020; Piękos et al., 2021; Razeghi et al., 2022).

5.7 Conclusion

We developed a framework to disentangle and separately measure the effect of different factors influencing the predictions of LLMs for math reasoning. Our results indicate that a drastic increase in both robustness and sensitivity emerges in the GPT-3 Davinci models. Additionally, we study the contribution of size and instruction tuning in the models of the LLaMA family, observing that the Alpaca instruction tuning, while increasing the model's robustness, does not significantly improve the overall performance. Our framework provides a formalized theory of behavioral testing for math reasoning models and opens new future directions to design behavioral tests of models in a principled way.

Ethical Considerations

As for the ethical practice in this work, the data involved are from existing MWP datasets with no private user information, and available under the MIT license. As for the ethical impact of the use of this work, the study is about providing a metric and analyzing existing models' robustness, so there is less concern over harmful usage. Rather, it is more about putting checks on existing AI models and helping humans understand them better before use. Potential stakeholders that could benefit from this research include NLP

researchers working on math models, practitioners working on various applications involving mathematical reasoning with text, and e-learning design.

Limitations

A key limitation in our work is that LLMs might have seen these math problems. Our work theoretically assumes this is not the case. Another limitation is that for the sake of simplicity, our work makes some assumptions. For example, we assume all numbers in the range of integers 0 to $C = 300$. This would not cover every MWP out there. And future work is needed to generalize our framework to other forms of MWPs. In this work, we are also constrained by the limitations of the OpenAI policy on the GPT-3 API. This limits the number of perturbations we consider in this work as well as the accuracy with which we can estimate our causal distributions. Finally, our work is restricted to English, and extending it to other languages will require us to create an MWP dataset in that language.

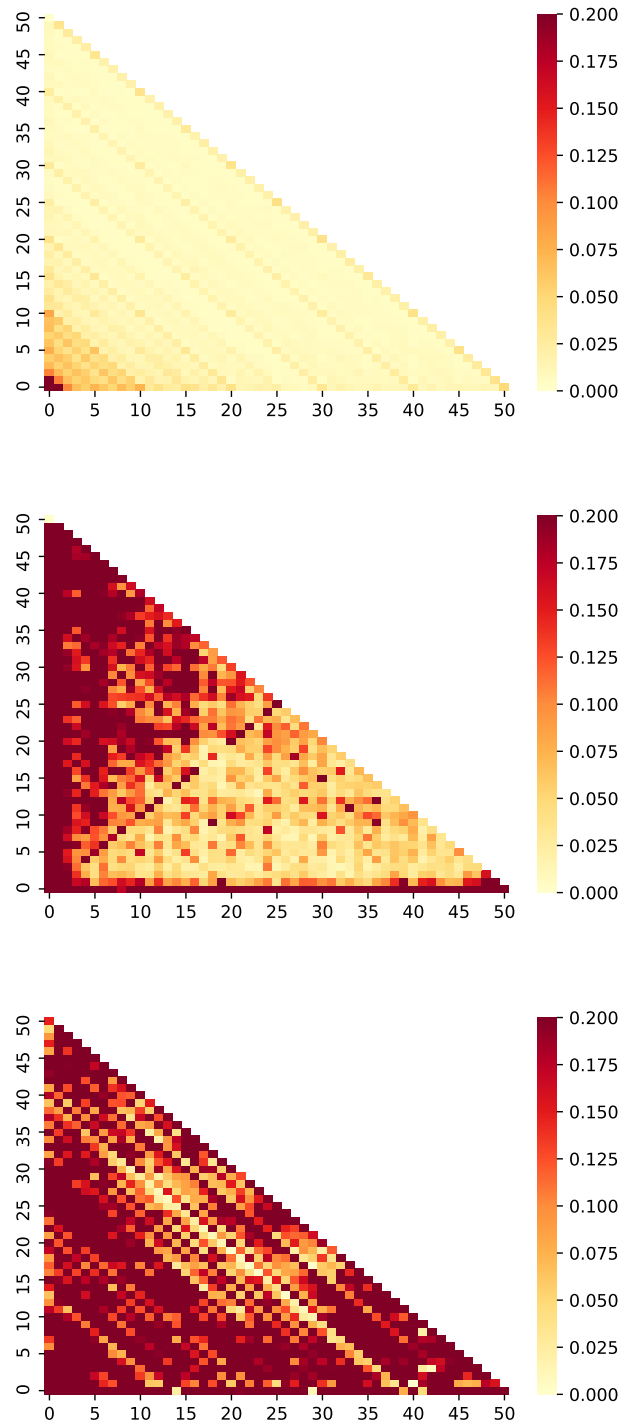


Figure 5.4: Heatmaps displaying $P(g)$ for Distil-GPT-2 (left), GPT-J-6B (center), and GPT-3 Davinci-002 (right). g is the ground-truth result $g = n_1 + n_2$ (n_1 and n_2 are represented by the x and y axes, respectively). The probability values for each combination of $((n_1, n_2), g)$ are averaged over 20 different templates. Probability values over 0.2 are displayed with the darkest color.

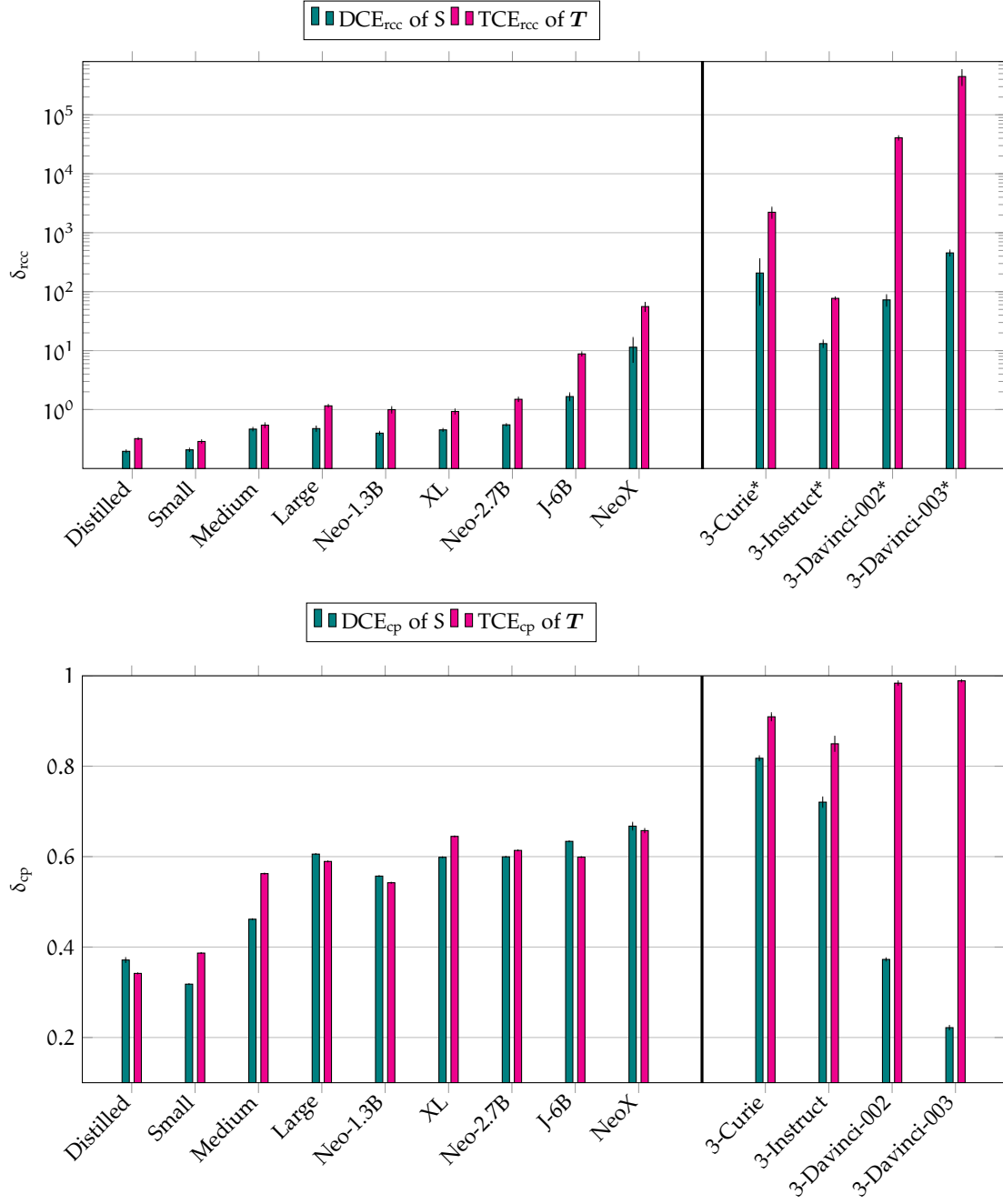


Figure 5.5: Comparison of DCE($S \rightarrow R$) and TCE(T on R). We use * to denote approximated values, explained in Appendix A.4.3.

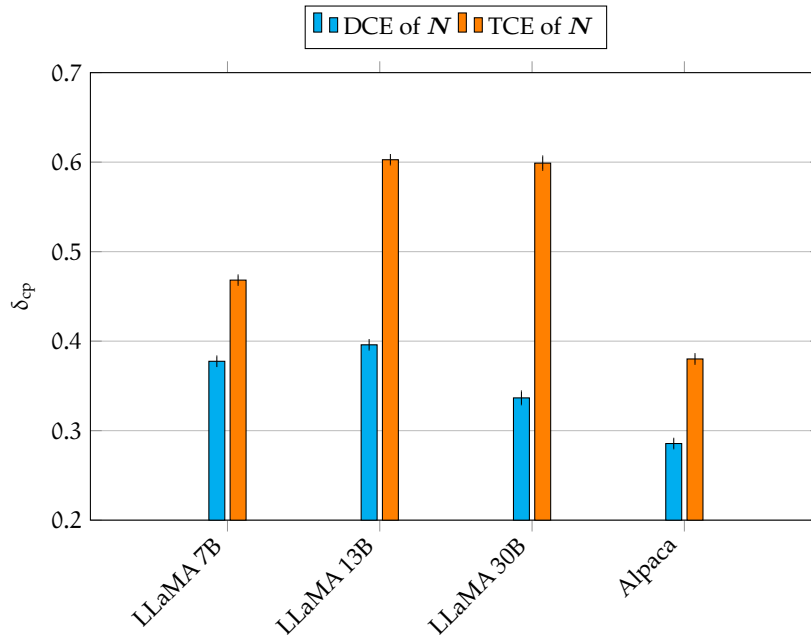


Figure 5.6: Comparison of direct and total effects of N on R for LLaMA and Alpaca.

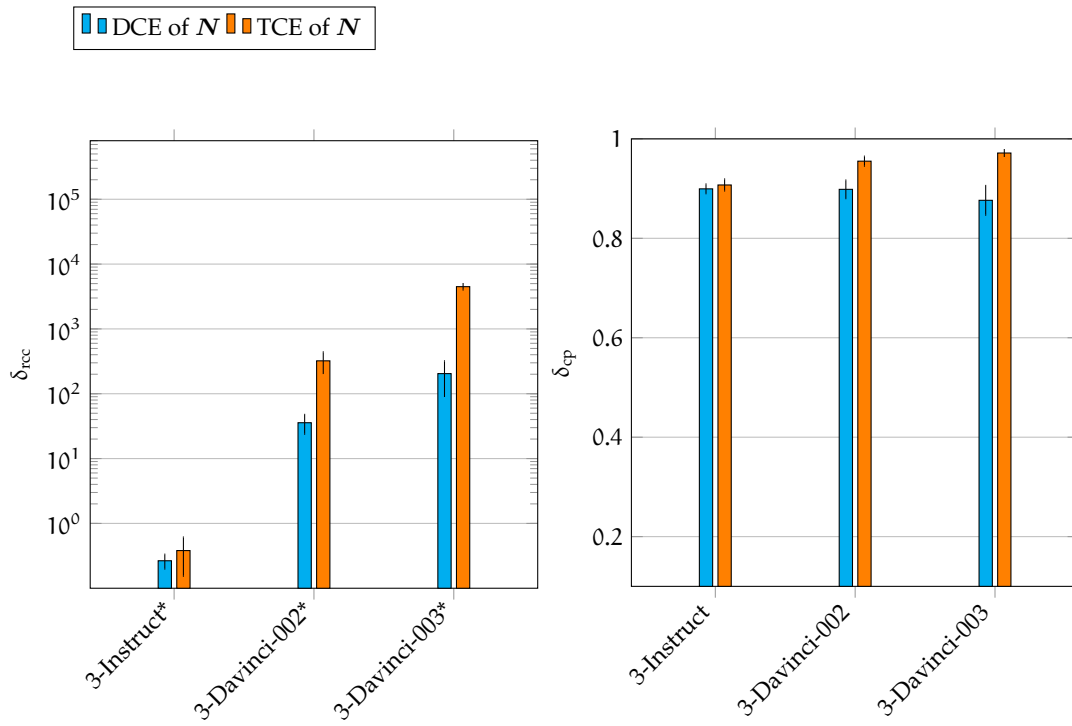


Figure 5.7: Comparison of direct and total effects of N on R for three-operand problems.

Part III

Causality among the Learning Variables

Causal Direction in Data Matters: Implications of Causal and Anticausal Learning in NLP

The principle of independent causal mechanisms (ICM) states that generative processes of real world data consist of independent modules which do not influence or inform each other. While this idea has led to fruitful developments in the field of causal inference, it is not widely-known in the NLP community. In this work, we argue that the causal direction of the data collection process bears nontrivial implications that can explain a number of published NLP findings, such as differences in semi-supervised learning (SSL) and domain adaptation (DA) performance across different settings. We categorize common NLP tasks according to their causal direction and empirically assay the validity of the ICM principle for text data using minimum description length. We conduct an extensive meta-analysis of over 100 published SSL and 30 DA studies, and find that the results are consistent with our expectations based on causal insights. This work presents the first attempt to analyze the ICM principle in NLP, and provides constructive suggestions for future modeling choices. Our code is available at <https://github.com/zhijing-jin/icm4nlp>.

6.1 Introduction

NLP practitioners typically do not pay great attention to the causal direction of the data collection process. As a motivating example, consider the case of collecting a dataset to train a machine translation (MT) model to translate from English (En) to Spanish (Es): it is common practice to mix all available En-Es sentence pairs together and train the model on the entire pooled data set (Bahdanau et al., 2015; Cho et al., 2014). However, such mixed corpora actually consist of two distinct types of data: (i) sentences that originated in English and have been translated (by human translators) into Spanish (En→Es); and (ii) sentences that originated in Spanish and have subsequently been translated into English (Es→En).¹

Intuitively, these two subsets are qualitatively different, and an increasing number of observations by the NLP community indeed suggests that they exhibit different properties (Freitag et al., 2019; Edunov et al., 2020; Riley et al., 2020; Shen et al., 2021a). In the case of MT, for example, researchers find that training models on each of these two types of data separately leads to different test performance, as well as different performance improvement by semi-supervised learning (SSL) (Bogoychev and Sennrich, 2019; Graham et al., 2020; Edunov et al., 2020). Motivated by this observation that the data collection process seems to matter for model performance, in this work, we provide an explanation of this phenomenon from the perspective of causality (Pearl, 2009b; Peters et al., 2017).

¹There is, in principle, a third option: both could be translations from a third language, but this occurs less frequently.

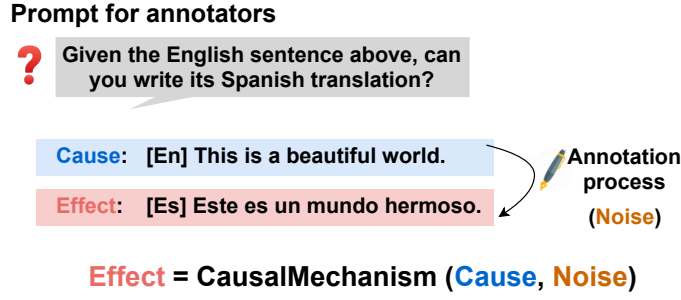


Figure 6.1: Annotation process for NLP data: the random variable that exists first is typically the cause (e.g., a given prompt), and the one generated afterwards is typically the effect (e.g., the annotated answer).

First, we introduce the notion of the *causal direction* for a given NLP task, see Figure 6.1 for an example. Throughout, we denote the input of a learning task by X and the output which is to be predicted by Y . If, during the data collection process, X is generated first, and then Y is collected based on X (e.g., through annotation), we say that X causes Y , and denote this by $X \rightarrow Y$. If, on the other hand, Y is generated first, and then X is collected based on Y , we say that Y causes X ($Y \rightarrow X$).²

Based on whether the direction of prediction aligns with the causal direction of the data collection process or not, Schölkopf et al. (2012) categorize these types of tasks as *causal learning* ($X \rightarrow Y$), or *anticausal learning* ($Y \rightarrow X$), respectively; see Figure 6.2 for an illustration. In the context of our motivating MT example this means that, if the goal is to translate from English ($X = \text{En}$) into Spanish ($Y = \text{Es}$), training *only* on subset (i) of the data consisting of $\text{En} \rightarrow \text{Es}$ pairs corresponds to *causal learning* ($X \rightarrow Y$), whereas training *only* on subset (ii) consisting of $\text{Es} \rightarrow \text{En}$ pairs is categorised as *anticausal learning* ($Y \rightarrow X$).

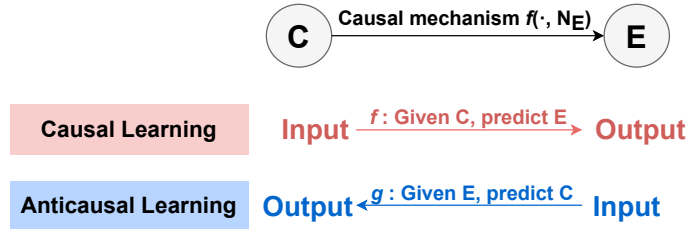


Figure 6.2: (Top) A causal graph $C \rightarrow E$, where C is the cause and E is the effect. The function $f(\cdot, N_E)$ denotes the causal process, or mechanism, $P_{E|C}$ by which the effect E is generated from C and unobserved noise N_E . (Bottom) Based on whether the direction of prediction aligns with the direction of causation or not, we distinguish two types of tasks: (i) causal learning, i.e., predicting the effect from the cause; and (ii) anticausal learning, i.e., predicting the cause from the effect.

Based on the principle of independent causal mechanisms (ICM) (Janzing and Schölkopf, 2010; Peters et al., 2017), it has been hypothesized that the causal direction of data collection (i.e., whether a given NLP learning task can be classified as causal or anticausal) has

²This corresponds to an *interventional* notion of causation: if one were to manipulate the cause, the annotation process would lead to a potentially different effect. A manipulation of the effect, in contrast, would not change the cause.

Category	Example NLP Tasks
Causal learning	Summarization, parsing, tagging, data-to-text generation, information extraction
Anticausal learning	Author attribute classification, review sentiment classification
Other/mixed (depending on data collection)	Machine translation, question answering, question generation, text style transfer, intent classification

Table 6.1: Classification of typical NLP tasks into causal (where the model takes the cause as input and predicts the effect), and anticausal (where the model takes the effect as input and predicts the cause) learning problems, as well as other tasks which do not have a clear causal interpretation of the data collection process, or where a mixture of both types of data is typically used.

implications for the effectiveness of commonly used techniques such as SSL and domain adaptation (DA) (Schölkopf et al., 2012). We will argue that this can explain performance differences reported by the NLP community across different data collection processes and tasks. In particular, we make the following contributions:

1. We categorize a number of common NLP tasks according to the causal direction of the underlying data collection process (Section 6.2).
2. We review the ICM principle and its implications for common techniques of using unlabelled data such as SSL and DA in the context of causal and anticausal NLP tasks (Section 6.3).
3. We empirically assay the validity of ICM for NLP data using minimum description length in a machine translation setting (Section 6.4).
4. We verify experimentally and through a meta-study of over respectively 100 (SSL) and 30 (DA) published findings that the difference in SSL (Section 6.5) and domain adaptation (DA) (Section 6.6) performance on causal vs anticausal datasets reported in the literature is consistent with what is predicted by the ICM principle.
5. We make suggestions on how to use findings in this paper for future work in NLP (Section 6.7).

6.2 Categorization of Common NLP Tasks

We start by categorizing common NLP tasks which use an input variable X to predict a target or output variable Y into causal learning ($X \rightarrow Y$), anticausal learning ($Y \rightarrow X$), and other tasks that do not have a clear underlying causal direction, or which typically rely on mixed (causal and anticausal) types of data, as summarised in Table 6.1.

Key to this categorization is determining whether the input X corresponds to the cause or the effect in the data collection process. As illustrated in Figure 6.1, if the input X and output Y are generated at two different time steps, then the variable that is generated first is typically the cause, and the other that is subsequently generated is typically the effect, provided it is generated based on the previous one (rather than, say, on a common confounder that causes both variables). If X and Y are generated jointly, then we need to distinguish based on the underlying generative process whether one of the two variables is causing the other variable.

Learning Effect from Cause (Causal Learning) Causal ($X \rightarrow Y$) NLP tasks typically aim to predict a post-hoc generated human annotation (i.e., the target Y is the effect) from a given input X (the cause). Examples include: summarization (*article*→*summary*) where the goal is to produce a summary Y of a given input text X ; parsing and tagging (*text*→*linguists' annotated structure*) where the goal is to predict an annotated syntactic structure Y of a given input sentence X ; data-to-text generation (*data*→*description*) where the goal is to produce a textual description Y of a set of structured input data X ; and information extraction (*text*→*entities/relations/etc*) where the goal is to extract structured information from a given text.

Learning Cause from Effect (Anticausal Learning) Anticausal ($Y \rightarrow X$) NLP tasks typically aim to predict or infer some latent target property Y such as an unobserved prompt from an observed input X which takes the form of one of its effects. Typical anticausal NLP learning problems include, for example, author attribute identification (*author attribute*→*text*) where the goal is to predict some unobserved attribute Y of the writer of a given text snippet X ; and review sentiment classification (*sentiment*→*review text*) where the goal is to predict the latent sentiment Y that caused an author to write a particular review X .

Other/Mixed Some tasks can be categorized as either causal or anticausal, depending on how exactly the data is collected. In Section 6.1, we discussed the example of MT where different types of (causal and anticausal) data are typically mixed. Another example is the task of intent classification: if the *same* author reveals their intent before the writing (i.e., *intent*→*text*), it can be viewed as an anticausal learning task; if, on the other hand, the data is annotated by *other* people who are not the original author (i.e., *text*→*annotated intent*), it can be viewed as a causal learning task. A similar reasoning applies to question answering and generation tasks which respectively aim to provide an answer to a given question, or vice versa: if first a piece of informative text is selected and annotators are then asked to come up with a corresponding question (*answer*→*question*) as, e.g., in the SQuAD dataset (Rajpurkar et al., 2016), then question answering is an anticausal and question generation a causal learning task; if, on the other hand, a question such as a search query is selected first and subsequently an answer is provided (*question*→*answer*) as, e.g., in the Natural Questions dataset (Kwiatkowski et al., 2019), then question answering is a causal and question generation an anticausal learning task. Often, multiple such datasets are combined without regard for their causal direction.

6.3 Implications of ICM for Causal and Anticausal Learning

Whether we are in a causal or anticausal learning scenario has important implications for semi-supervised learning (SSL) and domain adaptation (DA) (Schölkopf et al., 2012; Sgouritsa et al., 2015; Zhang et al., 2013, 2015a; Gong et al., 2016; von Kügelgen et al., 2019, 2020), which are techniques also commonly used in NLP. These implications are derived from the principle of independent causal mechanisms (ICM) (Schölkopf et al., 2012; Lemeire and Dirkx, 2006) which states that “*the causal generative process of a system’s variables is composed of autonomous modules that do not inform or influence each other*” (Peters et al., 2017).

In the bivariate case, this amounts to a type of independence assumption between the distribution P_C of the cause C , and the causal process, or mechanism, $P_{E|C}$ that generates

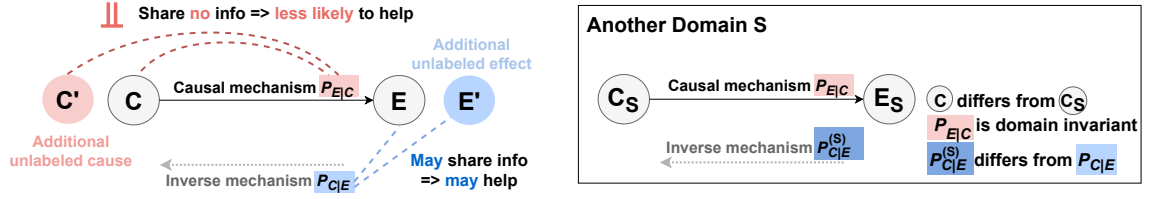


Figure 6.3: The ICM principle assumes that *the generative process P_C of the cause C is independent of the causal mechanism $P_{E|C}$* : the two distributions share no information and each may be changed or manipulated without affecting the other. In the anticausal direction, on the other hand, the effect distribution P_E is (in the generic case) *not independent of the inverse mechanism $P_{C|E}$* : they may share information and change dependently. (Left) SSL, which aims to improve an estimate of the target conditional $P_{Y|X}$ given additional unlabelled input data from P_X , should therefore not help for causal learning ($X \rightarrow Y$), but may help in the anticausal direction ($Y \rightarrow X$). (Right) DA, which aims to adapt a model of $P_{Y|X}$ from a source domain to a target domain (e.g., fine-tuning on a smaller dataset), should work better for causal learning settings where a change in P_C is not expected to lead to a change in the mechanism $P_{E|C}$, whereas in the anticausal direction P_E and $P_{C|E}$ may change in a dependent manner.

the effect from the cause. For example, for a question answering task, the generative process P_C by which one person comes up with a question C is “independent” of the process $P_{E|C}$ by which another person produces an answer E for question C .³

Here, “independent” is not meant in the sense of *statistical* independence of random variables, but rather as *independence at the level of generative processes or distributions* in the sense that P_C and $P_{E|C}$ *do not share information* (the person asking the question and the one answering may not know each other) and *can be manipulated independently of each other* (we can swap either of the two for another participant without the other one being influenced by this). Crucially, this type of independence is generally violated in the opposite, i.e., *anticausal*, direction: P_E and $P_{C|E}$ may share information and change dependently (Daniušis et al., 2010; Janzing et al., 2012). This has two important implications for common learning tasks (Schölkopf et al., 2012) which are illustrated in Figure 6.3.

Implications of ICM for SSL First, if P_C shares no information with $P_{E|C}$, SSL—where one has additional unlabelled input data from P_X and aims to improve an estimate of the target conditional $P_{Y|X}$ —should not work in the causal direction ($X \rightarrow Y$), but may work in the anticausal direction ($Y \rightarrow X$), as P_E and $P_{C|E}$ may share information. Causal NLP tasks should thus be less likely to show improvements over a supervised baseline when using SSL than anticausal tasks.

Implications of ICM for DA Second, according to the ICM principle, the causal mechanism $P_{E|C}$ should be invariant to changes in the cause distribution P_C , so domain—specifically, covariate shift (Shimodaira, 2000; Sugiyama and Kawanabe, 2012)—adaptation, where P_X changes but $P_{Y|X}$ is assumed to stay invariant, should work in the causal direction, but not necessarily in the anticausal direction. Hence, DA should be easier for causal

³The validity of this is meant in an approximate sense, and one can imagine settings where it is questionable. E.g., if the person asking the question has prior knowledge of the respondent (e.g., in a classroom setting), then she might adjust the question accordingly which would violate the assumption.

NLP tasks than for anticausal NLP tasks.

6.4 Validity of ICM for NLP Data Using MDL

Traditionally, the ICM principle is thought of in the context of *physical* processes or mechanisms, rather than *social* or *linguistic* ones such as language. Since ICM amounts to an independence assumption that—while well motivated in principle—may not always hold in practice,⁴ we now assay its validity on NLP data.

Recall, that ICM postulates a type of independence between P_C and $P_{E|C}$. One way to formalize this uses Kolmogorov complexity $K(\cdot)$ as a measure of algorithmic information, which can be understood as the length of the shortest program that computes a particular algorithmic object such as a distribution or a function (Solomonoff, 1964; Kolmogorov, 1965). ICM then reads (Janzing and Schölkopf, 2010):⁵

$$\begin{aligned} K(P_{C,E}) &\stackrel{+}{=} K(P_C) + K(P_{E|C}) \\ &\stackrel{+}{\leq} K(P_E) + K(P_{C|E}). \end{aligned} \quad (6.1)$$

In other words, the shortest description of the joint distribution $P_{C,E}$ corresponds to describing P_C and $P_{E|C}$ separately (i.e., they share no information), whereas there may be redundant (shared) information in the non-causal direction such that a separate description of P_E and $P_{C|E}$ will generally be longer than that of the joint distribution $P_{C,E}$.

6.4.1 Estimation by MDL

Since Kolmogorov complexity is not computable (Li et al., 2008), we adopt a commonly used proxy, the minimum description length (MDL) (Grünwald, 2007), to test the applicability of ICM for NLP data. Given an input, such as a collection of observations $\{(c_i, e_i)\}_{i=1}^n \sim P_{C,E}$, MDL returns the shortest codelength (in bits) needed to compress the input, as well as the parameters needed to decompress it. We use MDL to approximate Eq. (6.1) as follows:

$$\begin{aligned} \text{MDL}(c_{1:n}, e_{1:n}) &= \text{MDL}(c_{1:n}) + \text{MDL}(e_{1:n}|c_{1:n}) \\ &\leq \text{MDL}(e_{1:n}) + \text{MDL}(c_{1:n}|e_{1:n}), \end{aligned} \quad (6.2)$$

where $\text{MDL}(\cdot|\cdot)$ denotes a conditional compression where the second argument is treated as “free parameters” which do not count towards the compression length of the first argument. Eq. (6.2) can thus be interpreted as a comparison between two ways of compressing the same data $(c_{1:n}, e_{1:n})$: either we first compress $c_{1:n}$ and then compress $e_{1:n}$ conditional on $c_{1:n}$, or vice versa. According to the ICM principle, the first way should tend to be more “concise” than the second.

6.4.2 Calculating MDL Using Machine Translation as a Case Study

To empirically assess the validity of ICM for NLP data using MDL as a proxy, we turn to MT as a case study. We choose MT because the input and output spaces of MT are

⁴E.g., due to confounding influences from unobserved variables, or mechanisms which have co-evolved to be dependent

⁵Here, $\stackrel{+}{=}$ and $\stackrel{+}{\leq}$ hold up a constant due to the choice of a Turing machine in the definition of algorithmic information.

Dataset	Size	Note
En→Es	81K	Original English, Translated Spanish
Es→En	81K	Original Spanish, Translated English
En→Fr	16K	Original English, Translated French
Fr→En	16K	Original French, Translated English
Es→Fr	15K	Original Spanish, Translated French
Fr→Es	15K	Original French, Translated Spanish

Table 6.2: Details of the CausalMT corpus.

relatively symmetric, as opposed to other NLP tasks such as text classification where the input space is sequences, but the output space is a small set of labels.

There are only very few studies which calculate MDL on NLP data, so we extend the method of Voita and Titov (2020) to calculate MDL using online codes (Rissanen, 1984) for deep learning tasks (Blier and Ollivier, 2018). Since the original calculation method for MDL by Voita and Titov (2020) was developed for classification, we extend it to sequence-to-sequence (Seq2Seq) generation. Specifically, given a translation dataset $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ of n pairs of sentences x_i with translation y_i , denote the size of the vocabulary of the source language by V_x , and the size of the vocabulary of the target language by V_y . In order to assess whether Eq. (6.2) holds, we need to calculate four different terms: two marginal terms $\text{MDL}(x_{1:n})$ and $\text{MDL}(y_{1:n})$, and two conditional terms $\text{MDL}(y_{1:n}|x_{1:n})$ and $\text{MDL}(x_{1:n}|y_{1:n})$.

Codelength of the Conditional Terms To calculate the codelength of the two conditional terms, we extend the method of Voita and Titov (2020) from classification to Seq2Seq generation. Following the setting of Voita and Titov (2020), we break the dataset D into 10 disjoint subsets with increasing sizes and denote the end index of each subset as t_i .⁶ We then estimate $\text{MDL}(y_{1:n}|x_{1:n})$ as

$$\begin{aligned} \widehat{\text{MDL}}(y_{1:n}|x_{1:n}) &= \sum_{i=1}^{t_1} \text{length}(y_i) \cdot \log_2 V_y \\ &\quad - \sum_{i=1}^{n-1} \log_2 p_{\theta_i}(y_{1+t_i:t_{i+1}} | x_{1+t_i:t_{i+1}}), \end{aligned} \quad (6.3)$$

where $\text{length}(y_i)$ refers to the number of tokens in the sequence y_i , θ_i are the parameters of a translation model h_i trained on the first t_i data points, and $\text{seq}_{\text{idx}_1:\text{idx}_2}$ refers to the set of sequences from the idx_1 -th to the idx_2 -th sample in the dataset D , where $\text{seq} \in \{x, y\}$ and $\text{idx}_i \in \{1, \dots, n\}$. Similarly, when calculating $\text{MDL}(x_{1:n}|y_{1:n})$, we simply swap the roles of x and y .

Codelength of the Marginal Terms When calculating the two marginal terms, $\text{MDL}(x_{1:n})$ and $\text{MDL}(y_{1:n})$, we make two changes from the above calculation of conditional terms: first, we replace the *translation models* h_i with *language models*; second, we remove the conditional distribution. That is, we calculate $\text{MDL}(x_{1:n})$ as

$$\begin{aligned} \widehat{\text{MDL}}(x_{1:n}) &= \sum_{i=1}^{t_1} \text{length}(x_i) \cdot \log_2 V_x \\ &\quad - \sum_{i=1}^{n-1} \log_2 p_{\theta_i}(x_{1+t_i:t_{i+1}}), \end{aligned} \quad (6.4)$$

⁶The sizes of the 10 subsets are 0.1, 0.2, 0.4, 0.8, 1.6, 3.2, 6.25, 12.5, 25, and 50 percent of the dataset size, respectively. E.g., $t_1 = 0.1\%n$, $t_2 = (0.1\% + 0.2\%)n, \dots$

Data ($X \rightarrow Y$)	MDL(X)	MDL(Y)	MDL(Y X)	MDL(X Y)	MDL(X)+MDL(Y X) vs. MDL(Y)+MDL(X Y)
En \rightarrow Es	46.54	105.99	2033.95	2320.93	2080.49 < 2426.92
Es \rightarrow En	113.42	55.79	3289.99	3534.09	3403.41 < 3589.88
En \rightarrow Fr	20.54	53.83	503.78	535.88	524.32 < 589.71
Fr \rightarrow En	53.83	21.6	705.28	681.12	759.11 > 702.72
Es \rightarrow Fr	58.26	55.66	701.04	755.5	759.30 < 811.16
Fr \rightarrow Es	56.14	54.34	665.26	706.53	721.40 < 760.87

Table 6.3: Codelength (in kbits) of MDL(X), MDL(Y), MDL(Y|X), and MDL(X|Y) on six CausalMT datasets.

where θ_i are the parameters of a language model h_i trained on the first t_i data points. We apply the same method to calculate $\text{MDL}(\mathbf{y}_{1:n})$.

For the language model, we use GPT2 (Radford et al., 2019), and for the translation model, we use the Marian neural machine translation model (Junczys-Dowmunt et al., 2018) trained on the OPUS Corpus (Tiedemann and Nygaard, 2004). For fair comparison, all models adopt the transformer architecture (Vaswani et al., 2017), and have roughly the same number of parameters. See Appendix A.5.2 for more experimental details.

6.4.3 CausalMT Corpus

For our MDL experiment, we need datasets for which the causal direction of data collection is known, i.e., for which we have ground-truth annotation of which text is the original and which is a translation, instead of a mixture of both. Since existing MT corpora do not have this property as discussed in Section 6.1, we curate our own corpus, which we call the CausalMT corpus.

Specifically, we consider the existing MT dataset WMT’19,⁷ and identify some subsets that have a clear notion of causality. The subsets we use are the EuroParl (Koehn, 2005) and Global Voices translation corpora.⁸ For EuroParl, each text has meta information such as the speaker’s language; for Global Voices, each text has meta information about whether it is translated or not. We regard text that is in the same language as the speaker’s native language in EuroParl (and non-translated text in Global Voices) as the original (i.e., the cause). We then retrieve a corresponding effect by using the cause text to match the parallel pairs in the processed dataset. In this way, we compile six translation datasets with clear causal direction as summarized in Table 6.2. For each dataset, we use 1K samples each as test and validation sets, and use the rest for training.

6.4.4 Results

The results of our MDL experiment on the six CausalMT datasets are summarised in Table 6.3. If ICM holds, we expect the sum of codelengths to be smaller for the causal direction than for the anticausal one, see Eq. (6.2). As can be seen from the last column, this is the case for five out of the six datasets. For example, on one of the largest datasets (En \rightarrow Es), the MDL difference is 346 kbits.⁹

⁷Link to WMT’19.

⁸Link to Global Voices.

⁹As far as we know, determining statistical significance in the investigated setting remains an open problem. While, in theory, one may use information entropy to estimate it, in practice, this may be inaccurate

Comparing the dataset sizes in Table 6.2 and results in Table 6.3, we observe that the absolute MDL values are roughly proportional to dataset size, but other factors such as language and task complexity also play a role. This is inherent to the nature of MDL being the sum of codelengths of the model and of the data given the model. Since we use equally-sized datasets for each language pair in the CausalMT corpus (i.e., in both the $X \rightarrow Y$ and $Y \rightarrow X$ directions, see Table 6.2), numbers for the same language pair in Table 6.3, including the most important column “MDL(X)+MDL(Y | X) vs. MDL(Y)+MDL(X | Y)”, form a valid comparison. That is, En&Es experiments are comparable within themselves, so are the other language pairs.

For some of the smaller differences in the last column in Table 6.3, and, in particular the reversed inequality in row 4, a potential explanation may be the relatively small dataset size, as well as the fact that text data may be confounded (e.g., through shared grammar and semantics).

6.5 SSL for Causal vs. Anticausal Models

In semi-supervised learning (SSL), we are given a typically-small set of k labeled observations $D_L = \{(x_1, y_1), \dots, (x_k, y_k)\}$, and a typically-large set of m unlabeled observations of the input $D_U = \{x_1^{(u)}, \dots, x_m^{(u)}\}$. SSL then aims to use the additional information about the input distribution P_X from the unlabeled dataset D_U to improve a model of $P_{Y|X}$ learned on the labeled dataset D_L .

As explained in Section 6.3, SSL should only work for anticausal (or confounded) learning tasks, according to the ICM principle. Schölkopf et al. (2012) have observed this trend on a number of classification and regression tasks on small-scale numerical inputs, such as predicting Boston housing prices from quantifiable neighborhood features (causal learning), or breast cancer from lab statistics (anticausal learning). However, there exist no studies investigating the implications of ICM for SSL on NLP data, which is of a more complex nature due to the high dimensionality of the input and output spaces, as well as potentially large confounding. In the following, we use a sequence-to-sequence decipherment experiment (Section 6.5.1) and a meta-study of existing literature (Section 6.5.2) to showcase that the same phenomenon also occurs in NLP.

6.5.1 Decipherment Experiment

To have control over causal direction of the data collection process, we use a synthetic decipherment dataset to test the difference in SSL improvement between causal and anticausal learning tasks.

Dataset We create a synthetic dataset of encrypted sequences. Specifically, we (i) adopt a monolingual English corpus (for which we use the English corpus of the En→Es in the CausalMT dataset, for convenience), (ii) apply the ROT13 encryption algorithm (Schneier, 1996) to obtain the encrypted corpus, and then (iii) apply noise on the corpus that is chosen to be the effect corpus.

since (i) MDL is only a proxy for algorithmic information; and (ii) ICM may not hold exactly, but only approximately. We evaluate on six different datasets, so that the overall results can show a general trend.

Causal Data	Learning Task	Sup. BLEU	Δ SSL (BLEU)
En→Cipher	Causal	19.20	+1.84
	Anticausal	7.75	+38.02
Cipher→En	Causal	17.08	+4.05
	Anticausal	7.97	+38.01

Table 6.4: SSL improvements (Δ SSL) in BLEU score across causal vs. anticausal learning tasks on the synthetic decipherment datasets.

In the encryption step (ii), for each English sentence x , its encryption $\text{ROT13}(x)$ replaces each letter with the 13th letter after it in the alphabet, e.g., “A”→“N,” “B”→“O.” Note that we choose ROT13 due to its invertibility, since $\text{ROT13}(\text{ROT13}(x)) = x$. Therefore, without any noises, the corpus of English and the corpus of encrypted sequences by ROT13 are symmetric.

In the noising step (iii), we apply noise either to the English text or to the ciphertext, thus creating two datasets Cipher→En, and En→Cipher, respectively. When applying noise to a sequence, we use the implementation of the Fairseq library.¹⁰ Namely, we mask some random words in the sequence (word masking), permute a part of the sequence (permuted noise), randomly shift the endings of the sequence to the beginning (rolling noise), and insert some random characters or masks to the sequence (insertion noise). We set the probability of all noises to $p = 5\%$.

Results For each of the two datasets En→Cipher and Cipher→En, we perform SSL in the causal and anticausal direction by either treating the input X as the cause and the target Y as the effect, or vice versa. Specifically, we use a standard Transformer architecture for the supervised model, and for SSL, we multitask the translation task with an additional denoising autoencoder (Vincent et al., 2008) using the Fairseq Python package. The results are shown in Table 6.4. It can be seen that in both cases, anticausal models show a substantially larger SSL improvement than causal models.

We also note that there is a substantial gap in the supervised performance between causal and anticausal learning tasks on the same underlying data. This is also expected as causal learning is typically easier than anticausal learning since it corresponds to learning the “natural” forward function, or causal mechanism, while anticausal learning corresponds to learning the less natural, non-causal inverse mechanism.

6.5.2 SSL Improvements in Existing Work

After verifying the different behaviour in SSL improvement predicted by the ICM principle on the decipherment experiment, we conduct an extensive meta-study to survey whether this trend is also reflected in published NLP findings. To this end, we consider a diverse set of tasks, and SSL methods. The tasks covered in our meta-study include machine translation, summarization, parsing, tagging, information extraction, review sentiment classification, text category classification, word sense disambiguation, and chunking. The SSL methods include self-training, co-training (Blum and Mitchell, 1998), tri-training (Zhou and Li, 2005), transductive support vector machines (Joachims, 1999), ex-

¹⁰Link to the Fairseq implementation.

Task Type	Mean Δ SSL (\pm std)	According to ICM
Causal	+0.04 (\pm 4.23)	Smaller or none
Anticausal	+1.70 (\pm 2.05)	Larger

Table 6.5: Meta-study of SSL improvement (Δ SSL) across 55 causal and 50 anticausal NLP tasks.

Task Type	Mean Δ DA (\pm std)	According to ICM
Causal	5.18 (\pm 6.57)	Larger
Anticausal	1.26 (\pm 1.79)	Smaller

Table 6.6: Meta-study of DA improvement (Δ DA) across 22 causal and 11 anticausal NLP tasks.

pectation maximization (Nigam et al., 2006), multitasking with language modeling (Dai and Le, 2015), multitasking with sentence reordering (as used in Zhang and Zong (2016)), and cross-view training (Clark et al., 2018). Further details on our meta study are explained in Appendix A.5.1.

We covered 55 instances of causal learning and 50 instances of anticausal learning. A summary of the trends of causal SSL and anticausal SSL are listed in Table 6.5. Echoing with the implications of ICM stated in Section 6.3, for causal learning tasks, the average improvement by SSL is only very small, 0.04%. In contrast, the anticausal SSL improvement is larger, 1.70% on average. We use Welch’s t-test (Welch, 1947) to assess whether the difference in mean between the two distributions of SSL improvment (with unequal variance) is significant and obtain a p-value of 0.011.

6.6 DA for Causal vs. Anticausal Models

We also consider a supervised domain adaptation (DA) setting in which the goal is to adapt a model trained on a large labeled data set from a source domain, to a potentially different target domain from which we only have a small labeled data set. As explained in Section 6.3, DA should only work well for causal learning, but not necessarily for anticausal learning, according to the ICM principle.

Similar to the meta-study on SSL, we also review existing NLP literature on DA. We focus on DA improvement, i.e., the performance gain of using DA over an unadapted baseline that only learns from the source data and is tested on the target domain. Since the number of studies on DA that we can find is smaller than for SSL, we cover 22 instances of DA on causal tasks, and 11 instances of DA on anticausal tasks.

The results are summarised in Table 6.6. We find that the observations again echo with our expectations (according to ICM) that DA should work better for causal, than for anticausal learning tasks. Again, we use Welch’s t-test (Welch, 1947) to verify that the DA improvements of causal learning and anticausal learning are statistically different, and obtain a p-value of 0.023.

6.7 How to Use the Findings in this Study

Data Collection Practice in NLP Due to the different implications of causal and anti-causal learning tasks, *we strongly suggest annotating the causal direction when collecting new NLP data*. One way to do this is to only collect data from one causal direction and to mention this in the meta information. For example, summarization data collected from the TL;DR of scientific papers SciTldr (Cachola et al., 2020) should be *causal*, as the TL;DR summaries on OpenReview (some from authors when submitting the paper, others derived from the beginning of peer reviews) were likely composed after the original papers or reviews were written. Alternatively, one may allow mixed corpora, but label the causal direction for each (x, y) pair, e.g., which is the original vs. translated text in a translation pair. Since more data often leads to better model performance, it is common to mix data from both causal directions, e.g., training on both $\text{En} \rightarrow \text{Es}$ and $\text{Es} \rightarrow \text{En}$ data. Annotating the causal direction for each pair allows future users of the dataset to potentially handle the causal and anticausal parts of the data differently.

Causality-Aware Modeling When building NLP models, the causal direction provides additional information that can potentially be built into the model. In the MT case, since causal and anticausal learning can lead to different performance (Ni et al., 2022), one way to take advantage of the known causal direction is to add a prefix such as “[Modeling-Effect-to-Cause]” to the original input, so that the model can learn from causally-annotated input-output pairs. For example, Riley et al. (2020) use labels of the causal direction to elicit different behavior at inference time. Another option is to carefully design a combination of different modeling techniques, such as limiting self-training (a method for SSL) only to the anticausal direction and allowing back-translation in both directions, as preliminarily explored by Shen et al. (2021a).

Causal Discovery Suppose that we are given measurements of two types of NLP data X and Y (e.g., text, parse tree, intent type) whose collection process is unknown, i.e., which is the cause and which the effect. One key finding of our study is that there is typically a causal footprint of the data collection process which manifests itself, e.g., when computing the description length in different directions (Section 6.4) or when performing SSL (Section 6.5) or DA (Section 6.6). Based on which direction has the shorter MDL, or allows better SSL or DA, we can thus infer one causal direction over the other.

Prediction of SSL and DA Effectiveness Being able to predict the effectiveness of SSL or DA for a given NLP task can be very useful, e.g., to set the weights in an ensemble of different models (Søgaard, 2013). While predicting SSL performance has previously been studied from a non-causal perspective (Nigam and Ghani, 2000; Asch and Daelemans, 2016), our findings suggest that a simple qualitative description of the data collection process in terms of its causal direction (as summarised for the most common NLP tasks in Table 6.1) can also be surprisingly effective to evaluate whether SSL or DA should be expected to work well.

6.8 Limitations and Future Work

We note that ICM—when taken strictly—is an idealized assumption that may be violated and thus may not hold exactly for a given real-world data set, e.g., due to confounding,

i.e., when both variables are influenced by a third, unobserved variable. In this case, one may observe less of a difference between causal and anticausal learning tasks.

We also note that, while we have made an effort to classify different NLP tasks as *typically* causal or anticausal, our categorization should not be applied blindly without regard for the specific generative process at hand: deviations are possible as explained in the Mixed/Other category.

Another limitation is that the SSL and DA settings considered in this paper are only a subset of the various settings that exist in NLP. Our study does not cover, for example, SSL that uses additional output data (e.g., Jean et al. (2015); Gülçehre et al. (2015); Sennrich and Zhang (2019)), or unsupervised DA (as reviewed by Ramponi and Plank (2020)). In addition, in our meta-study of published SSL and DA findings, the improvements of causal vs. anticausal learning might be amplified by the scale of research efforts on different tasks and potentially suffer from selection bias.

Finally, we remark that, in the present work, we have focused on bivariate prediction tasks with an input X and output Y . Future work may also apply ICM-based reasoning to more complex NLP settings, for example, by (i) incorporating additional (sequential/temporal) structure of the data (e.g., for MT or language modeling) or (ii) considering settings in which the input X consists of both cause X_{CAU} and effect X_{EFF} features of the target Y (von Kügelgen et al., 2019, 2020).

6.9 Related Work

NLP and Causality Existing work on NLP and causality mainly focuses on the extracting text features for causal inference. Researchers first propose a causal graph based on domain knowledge, and then use text features to represent some elements in the causal graph, e.g., the cause (Egami et al., 2018; Jin et al., 2021b), effect (Fong and Grimmer, 2016), and confounders (Roberts et al., 2020; Veitch et al., 2020; Keith et al., 2020). Another line of work mines causal relations among events from textual expressions, and uses them to perform relation extraction (Do et al., 2011; Mirza and Tonelli, 2014; Dunietz et al., 2017; Hosseini et al., 2021), question answering (Oh et al., 2016), or commonsense reasoning (Sap et al., 2019a; Bosselut et al., 2019). For a recent survey, we refer to Feder et al. (2021a).

Usage of MDL in NLP Although MDL has been used for causal discovery for low-dimensional data (Budhathoki and Vreeken, 2017; Mian et al., 2021; Marx and Vreeken, 2021), only very few studies adopt MDL on high-dimensional NLP data. Most existing uses of MDL on NLP are for probing and interpretability: e.g., Voita and Titov (2020) use it for probing of a small Bayesian model and network pruning, based on the method proposed by Blier and Ollivier (2018) to calculate MDL for deep learning. We are not aware of existing work using MDL for causal discovery, or to verify causal concepts such as ICM in the context of NLP.

Existing Discussions on SSL and DA in NLP SSL and DA has long been used in NLP, as reviewed by Søgaard (2013) and Ramponi and Plank (2020). However, there have been a number of studies that report negative results for SSL (Clark et al., 2003; Steedman et al.,

2003; Reichart and Rappoport, 2007; Abney, 2007; Spreyer and Kuhn, 2009; Søgaard and Rishøj, 2010) and DA (Plank et al., 2014). Our work constitutes the first explanation of the ineffectiveness of SSL and DA on certain NLP tasks from the perspective of causal and anticausal learning.

6.10 Conclusion

This work presents the first effort to use causal concepts such as the ICM principle and the distinction between causal and anticausal learning to shed light on some commonly observed trends in NLP. Specifically, we provide an explanation of observed differences in SSL (Tables 6.4 and 6.5) and DA (Table 6.6) performance on a number of NLP tasks: DA tends to work better for causal learning tasks, whereas SSL typically only works for anticausal learning tasks, as predicted by the ICM principle. These insights, together with our categorization of common NLP tasks (Table 6.1) into causal and anticausal learning, may prove useful for future NLP efforts. Moreover, we empirically confirm using MDL that the description of data is typically shorter in the causal than in the anticausal direction (Table 6.3), suggesting that a causal footprint can also be observed for text data. This has interesting potential implications for discovering causal relations between different types of NLP data.

Ethical Considerations

Use of Data This paper uses two types of data, a subset of an existing machine translation dataset, and synthetic decipherment data. As far as we know, there are no sensitive issues such as privacy regarding the data usage.

Potential Stakeholders This research focuses on meta properties of two commonly applied methodologies, SSL and DA in NLP. Although this research is not directly connected to specific applications in society, the usage of this study can benefit future research in SSL and DA.

On the Causal Nature of Sentiment Analysis

Sentiment analysis (SA) aims to identify the sentiment expressed in a text, such as a product review. Given a review and the sentiment associated with it, this work formulates SA as a combination of two tasks: (1) a causal discovery task that distinguishes whether a review “primes” the sentiment (Causal Hypothesis C1), or the sentiment “primes” the review (Causal Hypothesis C2); and (2) the traditional prediction task to model the sentiment using the review as input. Using the peak-end rule in psychology, we classify a sample as C1 if its overall sentiment score approximates an average of all the sentence-level sentiments in the review, and C2 if the overall sentiment score approximates an average of the peak and end sentiments. For the prediction task, we use the discovered causal mechanisms behind the samples to improve LLM performance by proposing *causal prompts* that give the models an inductive bias of the underlying causal graph, leading to substantial improvements by up to 32.13 F1 points on zero-shot five-class SA. Our code is available at <https://github.com/cogito233/causal-sa>.

7.1 Introduction

Sentiment analysis (SA) is the task of identifying the sentiment y given a piece of text x . The field has a rich history originating from subjectivity analysis (Wiebe, 1994; Hatzivassiloglou and Wiebe, 2000), and developed rapidly with the availability of large opinionated online data such as reviews with ratings (Turney, 2002; Nasukawa and Yi, 2003; Zhang et al., 2015b; Keung et al., 2020, *inter alia*).

Despite recent advances in large language models (LLMs), it is still challenging to address the fine-grained five-class SA (which corresponds to the five star ratings in most datasets) for document-level classification (Choi et al., 2020; Fei et al., 2023; Truică et al., 2021), due to the subtle nature of the task including aspects such as inter-aspect relations, commonsense reasoning, among others (Poria et al., 2023; Venkit et al., 2023).

In this paper, we propose a causally-informed solution for the SA task. Different from the approach of naïvely applying up-to-date LLMs, we leverage insights from causal inference to propose a reformulation for SA into two tasks, as in Figure 7.1: (1) a causal discovery task to identify the cause-effect relation between the review X and the sentiment Y , and (2) the traditional prediction task $f : x \mapsto y$ to model the sentiment using the review as input.

We first look into the causal discovery task. In the study of affect science (Salovey and Mayer, 2004; Barrett, 2006; Feinstein, 2013), language can be the cause of emotion (Satpute et al., 2013; Kassam and Mendes, 2013) – namely a review priming the following sentiment, i.e., the Causal Hypothesis C1 of $X \rightarrow Y$; or emotion can affect the use of language (Barrett, 2006) – namely sentiment priming the review as an ad-hoc justification for

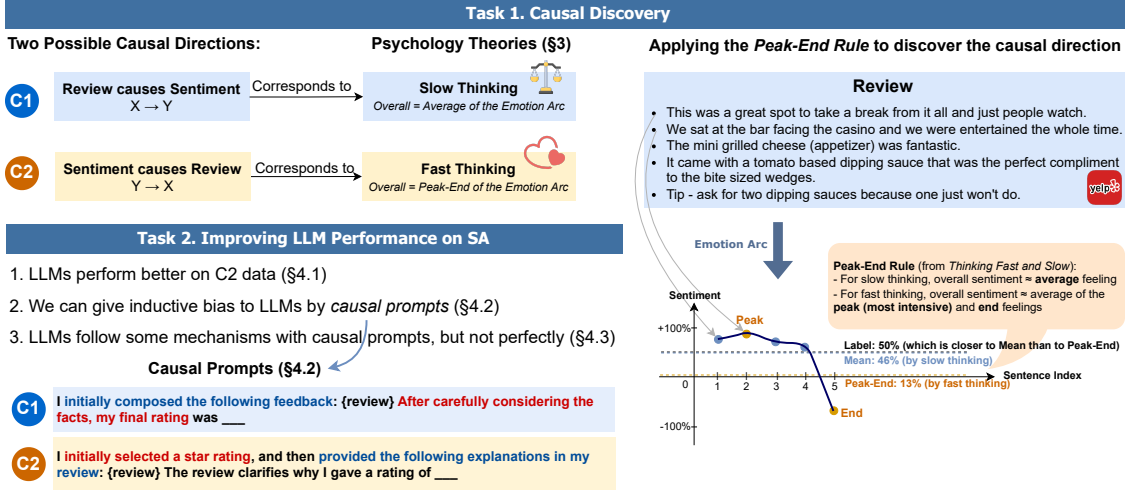


Figure 7.1: An overview of the paper structure, where we first investigate the causal discovery task, and then use it to improve LLM performance. For each document-level text review, we parse its *emotion arc* consisting of the sentiment of each sentence in the review, and then use the peak-end rule (Kahneman et al., 1993; Kahneman, 2011) to identify whether the overall sentiment is an average of the arc (corresponding to *Slow Thinking*), or an average of the peak and end sentiments (corresponding to *Fast Thinking*).

the emotion, i.e., the Causal Hypothesis C2 of $Y \rightarrow X$. These two processes might arise from the data annotation process (Jin et al., 2021c), but is hard to discover post-hoc in existing datasets.

Given the possibility of both causal directions $X \rightarrow Y$ or $Y \rightarrow X$ in the SA data, we identify the actual underlying mechanism based on insights from psychology (Kahneman, 2011; Epstein, 1994). Specifically, we identify the correspondence of the above two causal mechanisms with the *Fast* and *Slow Thinking* systems (Kahneman, 2011): (1) a review-driven sentiment (as in C1) largely resembles the Slow Thinking process applying reasoning based on evidence, and (2) the process of first coming up with the sentiment and then justifying it by a review (as in C2) conforms to Fast Thinking. Given this correspondence, we apply the peak-end rule from psychology (Kahneman et al., 1993; Kahneman, 2011). As shown in the right part of Figure 7.1, we classify a sample as C1 if its overall sentiment score approximates an average of all the sentence-level sentiments in the review, and as C2 if the overall sentiment score approximates an average of the peak and end sentiments.

Based on the identified causal mechanism behind SA data from the causal discovery task, we further explore how it can improve prediction performance in the era of LLMs. Existing literature highlights “causal alignment,” namely to align the prediction direction along the underlying causal direction (Jin et al., 2021c; Schölkopf, 2022; Schölkopf et al., 2021), but to our knowledge we are the first to explore how causal alignment improves model performance of SA in the era of LLMs. Specifically, we answer three subquestions: (Q1) If using the standard SA prompt, do models perform differently on C1/C2 data? (Q2) Does it help if we make the prompt aware of the underlying causality, i.e., use *causal prompts*? And (Q3) When prompted causally, do LLMs mechanistically understand the corresponding causal processes?

Our empirical results show that under the standard prompt, LLMs perform better on data corresponding to the C2 causal process. Moreover, causal prompts aligned with the causal direction of the data can substantially improve the performance of zero-shot SA. Finally, we apply mechanistic interpretability methods to probe the models, and find that there is still improvement space for LLMs to correctly grasp the essence of the two causal processes. In summary, the contributions of this paper are as follows:

1. We propose the dual nature of SA as a combination of two tasks: a causal discovery task, and a prediction task.
2. For causal discovery, we ground the two possible causal processes in psychology, and use the peak-end rule to identify them.
3. For the prediction task, we inspect existing LLMs' performance on data corresponding to the two underlying causal processes, and design *causal prompts* to improve model performance by up to 32.13 F1 points.

7.2 Problem Formulation of SA

In this section, we formulate SA as a combination of two tasks: the traditional prediction task in NLP and the causal discovery task in statistics, which we will introduce in the following.

7.2.1 The Prediction Task (in NLP)

SA is a prediction task to identify the sentiment y given a piece of text x . We adopt the setup in most existing SA datasets (Maas et al., 2011; Zhang et al., 2015b; Keung et al., 2020), where the text x is a review consisting of n sentences (t_1, \dots, t_n) , and the label y is a sentiment score corresponding to the star rating of the review in 1 (most negative), 2, ..., 5 (most positive).

7.2.2 The Cause-Effect Discovery Task (in Statistics)

As a separate problem, there is an established task in causal discovery, the causal-effect problem (see the review by Janzing, 2019), which aims to tell the cause from effect using only observational data. Its formal formulation is as follows: Suppose we have an i.i.d. dataset $\mathcal{D} := \{(x_i, y_i)\}_{i=1}^n$ containing n observational data pairs of the two variables, X and Y . The task is to *infer whether X causes Y (i.e., $X \rightarrow Y$), or Y causes X (i.e., $Y \rightarrow X$)*, if one out of the two is true. In causality, " \rightarrow " indicates the directional causal relation between two variables. The two hypotheses can also be expressed in their equivalent structural causal models (SCMs; Pearl, 2009b) as introduced in Peters et al. (2017):

$$\text{Causal Hypothesis 1 (C1): } X \rightarrow Y \quad (7.1)$$

$$\Leftrightarrow Y := f_Y(X, N_Y) \text{ with } N_Y \perp X, \quad (7.2)$$

$$\text{Causal Hypothesis 2 (C2): } Y \rightarrow X \quad (7.3)$$

$$\Leftrightarrow X := f_X(Y, N_X) \text{ with } N_X \perp Y, \quad (7.4)$$

where N_i is an unobserved noise term orthogonal to the input distribution.

7.2.3 Causality and NLP Model Performance

For many years, causality and machine learning have been two separate domains on their own. Recently, researchers started to think about how the causal knowledge of the data can improve machine learning performance on the prediction task, especially for the two variable cause-effect case (Schölkopf et al., 2012; Jin et al., 2021c; Ni et al., 2022). The essence of this line of research is that causality makes the two learning tasks $x \mapsto y$ and $y \mapsto x$ asymmetric, as one function’s prediction direction *aligns with* the ground-truth causal direction behind the two random variables, and another contradicts. We call this phenomenon “*causal alignment*,” or “*direction match*,” of the prediction task and the causality.

To contrast the contribution of our work, we review the previous literature on causal alignment, which only shows its effect on the performance of trained-from-scratch machine learning models, without any indications in the era of LLMs:

1. Causal alignment makes a model more robust against **covariant shifts** (Jin et al., 2021c; Schölkopf, 2022; Schott et al., 2018, *inter alia*)
2. Semi-supervised learning (SSL) only works under causal misalignment, as the cause variable contains no information about the mechanism, but the effect variable does. So in the misaligned case, additional P_X (i.e., the effect variable) helps **SSL** (Schölkopf et al., 2012; Jin et al., 2021c).
3. Learning a causally-aligned model induces **less Kolmogorov complexity** (a more minimal description length) than the causally-misaligned model on the same X-Y data (Jin et al., 2021c; Janzing and Schölkopf, 2010)
4. Causal alignment significantly affects model performance in **supervised learning**, in the case of machine translation (Ni et al., 2022).

All the above findings are drawn under the training condition that we can isolate the training data to be only of one causal direction. In the era of LLMs, we have seen substantial differences: (1) the training data can be a mixture of both causal directions, (2) the operationalization of the prediction task is through prompting, but no longer a separate model for each direction, and, (3) in general, research has shifted to designing better prompts for already pre-trained models in their inference mode.

Given these changes, we use the rest of the paper to address the following research questions:

1. What is the causal direction in SA? (Section 7.3)
2. Can causal alignment help us improve SA prompts in the era of LLMs? (Section 7.4)

7.3 Causal Discovery of Sentiment and Review

7.3.1 Problem Setup

As mentioned previously, the setup of the bivariate causal discovery problem is to infer whether X causes Y (C1), or Y causes X (C2), based on a dataset $\mathcal{D} := \{(x_i, y_i)\}_{i=1}^n$ containing only observational data of the joint distribution.

Challenges The common paradigm to check causal discovery results is to generate simulated data, of which the ground truth causal graph is known (Zhang and Hyvärinen, 2009; Spirtes and Zhang, 2016). However, in the context of the established SA datasets, such as Yelp (Zhang et al., 2015b), Amazon (Keung et al., 2020), and App Review (Grano et al., 2017), we would not be able to track each individual user and survey their original causal process when composing the review and the rating. Another solution would also be difficult, as it would require SA to abandon all the above well-established datasets, and meticulously collect new data while surveying the users’ underlying causal process.

Our Approach In the context of our work, we propose that there are still rich findings that we could derive from the observation-only data in the existing datasets, without interviewing or conducting new costly data collection.

The key to our approach is the psychology theories of the two causal processes, as the relation between sentiment and text has been well-studied and verified by randomized control trials (RCTs), among many other experiments. In the rest of the section, we first introduce in Section 7.3.2 the psychology theories of fast and slow thinking, followed by the Peak-End Rule as the quantitative signal. Then, we operationalize the theory with computational techniques in Section 7.3.3, and the report findings on three different SA datasets in Section 7.3.4.

7.3.2 Psychological Processes Underlying Sentiment Processing

Two Systems of Emotional Responses In psychology, the bifurcation into System 1 and System 2 in human decision-making, including sentiment processing, has garnered substantial empirical support (Kahneman, 2011; Epstein, 1994).

System 1, or the “*Fast Thinking*” system, operates involuntarily, effortlessly, and without conscious awareness. It is often optimized in evolution to provide rapid responses to environmental stimuli (LeDoux, 1998), and guides most of our daily cognitive processing (Kahneman, 2011), and emotional responses such as fear or joy (Zajonc, 1980).

Conversely, System 2, often termed as “*Slow Thinking*,” is deliberate, slower, and more rational, requires more conscious effort (Kahneman, 2011), and allows for self-regulation and thoughtful consideration before making decisions (Baumeister et al., 1998). The interplay between these systems influences everything from mundane to critical decisions, highlighting the complexity of human emotional and cognitive processing (Kahneman and Frederick, 2002; Kahneman, 2011).

Correspondence to the Two Causal Processes There is a nice correspondence between the fast/slow thinking systems and our two causal hypotheses. As mentioned previously, the Causal Hypothesis 2 (C2) posits $Y \rightarrow X$, where the sentiment Y causes the review X , which aligns well with the *Fast Thinking* system (Kahneman, 2011; LeDoux, 1998), as it rapidly generates an emotional reaction Y , and then writes text to justify it Y . On the other hand, the Causal Hypothesis 1 (C1) refers to the case where $X \rightarrow Y$, namely the review X causing the sentiment Y . It is an instance of *Slow Thinking* (Baumeister et al., 1998; Kahneman and Frederick, 2002), which deliberately uses conscious efforts to list out the up- and downsides of an experience in the review X , and come up with a thoughtful final decision as the rating Y .

	Yelp			Amazon			App Review		
	All	C1	C2	All	C1	C2	All	C1	C2
# Samples	34,851	19,557 (56%)	15,294 (44%)	2,582	1,393 (54%)	1,189 (46%)	9,696	3,809 (39%)	5,887 (61%)
# Sents/Review	11.11	11.30	10.87	6.70	6.62	6.80	6.34	6.33	6.35
# Words/Sent	15.53	15.55	15.49	11.04	11.28	10.77	10.53	10.90	10.29
Vocab Size	64,864	48,889	44,826	10,271	7,609	7,049	20,248	12,773	15,400
Avg Sentiment	2.93	2.74	3.18	2.94	2.82	3.07	2.9	2.72	3.02
Avg λ_1	3.78	2.97	4.83	3.77	3.10	4.55	6.03	5.18	6.58
Avg λ_2	4.48	6.05	2.48	4.21	5.72	2.45	5.10	7.96	3.26

Table 7.1: Statistics of the entire datasets and their C1 and C2 subsets for Yelp, Amazon, and App Review. We can see that a roughly balanced number of reviews aligning with the C1 and C2 processes.

Quantitative Signals of the Two Processes In sentiment processing, an evidence for the two processes is the famous Kahneman et al. (1993) study illustrating the *Peak-End Rule* of how individuals recall and evaluate past emotional experiences, which we show in Figure 7.1. As we know, fast thinking is prone to systematic biases and errors in the judgment (Tversky and Kahneman, 1974), and the Kahneman et al. (1993) study provides important quantitative results showing that, in the Fast Thinking system, people’s emotional memories of an experience are disproportionately influenced by its most intense point (the “peak”) and its conclusion (the “end”), rather than by the average experience as in the Slow Thinking system. The important role of peak and end for the fast thinking system implies that it is the intensity of specific moments that dominate memory and judgment.

7.3.3 Operationalizing the Theory

We summarize the previous psychological insights in the upper left part of Figure 7.1, where the Causal Hypothesis 1 corresponds to taking the average of all emotional experiences mentioned in the review X for the sentiment Y , and the Causal Hypothesis 2 uses the peak and end emotions in the review X to derive the sentiment Y . In this section, we introduce a formalization of the theory, and suggest signals to distinguish the two causal hypotheses.

Emotion Arc To capture the aforementioned trajectory of emotional experiences, we use the concept of the *emotion arc* (Reagan et al., 2016), an example of which we visualize in Figure 7.1. Contextualizing it in the task of SA, we formally define an emotion arc of the review as follows. Given a review x consisting of n sentences (t_1, \dots, t_n) , we identify the sentiment for each of them, thus obtaining a series of sentiment labels (s_1, \dots, s_n) . We denote this series as the emotion arc $e := (s_1, \dots, s_n)$ of the review.

The Two Causal Processes Provided the notion of the emotion arc $e := (s_1, \dots, s_n)$ for a review x , we formulate the sentiment labels corresponding to the two causal processes

<p style="text-align: center;">Example of a C1-Dominant Review</p> <ul style="list-style-type: none"> • This was a great spot to take a break from it all and just people watch. $s_1 = 4.57$ • We sat at the bar facing the casino and we were entertained the whole time. $s_2 = 4.67$ • The mini grilled cheese (appetizer) was fantastic. $s_3 = 4.53$ • It came with a tomato based dipping sauce that was the perfect compliment to the bite sized wedges. $s_4 = 4.20$ • Tip - ask for two dipping sauces because one just won't do. $s_5 = 1.60$ <p>Stars y: 4</p> <p>Psychology Scores (λ): $\lambda_1 = 0.0884 < \lambda_2 = 0.8683$</p>	<p style="text-align: center;">Example of a C2-Dominant Review</p> <ul style="list-style-type: none"> • I read the reviews and should have steered away... but it looked interesting. $s_1 = 3.72$ • Salad was wilted, menus are on the wall, with no explanation so you are ordering blind, service was NOT with a smile from the bartender to the waitress, to the server who helped the waitress, and the waitress never checked back to see how everything is. $s_2 = 2.20$ • Terribly overpriced for what you get, and as an Italian, this does not even pass for a facsimile thereof! $s_3 = 1.45$ • Stay away for sure. $s_4 = 1.85$ • I only gave them one star, as I had to fill something in, they should get no stars! $s_5 = 1.32$ <p>Stars y: 1</p> <p>Psychology Scores (λ): $\lambda_1 = 1.1827 > \lambda_2 = 0.3647$</p>
---	---

Table 7.2: Examples of C1- and C2-dominant reviews.

as follows:

Slow Thinking (Causal Process 1):

$$\hat{y}_{\text{avg}} = \frac{1}{n} (s_1 + \dots + s_n) , \quad (7.5)$$

$$\lambda_1 = |y - \hat{y}_{\text{avg}}| , \quad (7.6)$$

Fast Thinking (Causal Process 2):

$$\hat{y}_{\text{peakEnd}} = \frac{1}{2} (\text{Peak}(s_1, \dots, s_n) + s_n) , \quad (7.7)$$

$$\lambda_2 = |y - \hat{y}_{\text{peakEnd}}| , \quad (7.8)$$

where λ_i indicates the alignment of the actual sentiment y with the Causal Process i , and $\text{Peak}(\cdot)$ selects the sentiment with the strongest intensity by its distance from the neutral sentiment 3, which is the middle point among the sentiment range 1–5, i.e., $\text{Peak}(s_1, \dots, s_n) := s_{\arg\max_i |s_i - 3|}$.

Here, we interpret λ_i as an indicator for each causal process, where a small value (with the best value being zero) implies the alignment with the process i . We show two examples in Table 7.2, one aligning well with the Causal Process C1 with a small λ_1 , and another aligning well with the Causal Process C2 with a small λ_2 .

7.3.4 Findings on SA Datasets

Dataset Setup We adopt three commonly used datasets in SA: Yelp (Zhang et al., 2015b), Amazon (Keung et al., 2020), and App Review (Grano et al., 2017). For the Amazon data, we concatenate each review's title with its text. Since the model performance on many

		Random	GPT-2 XL	LLaMa-7B	Alpaca-7B	GPT-3	GPT-3.5	GPT-4
F1	Overall	19.82 \pm 2.07	10.23 \pm 4.12	31.78 \pm 5.32	46.01 \pm 5.35	52.71 \pm 1.73	57.98 \pm 5.11	59.54 \pm 4.69
	C1 Subset	21.36 \pm 2.26	5.80 \pm 3.11	27.30 \pm 4.73	37.77 \pm 7.66	43.96 \pm 2.93	58.64 \pm 1.48	58.62 \pm 2.54
	C2 Subset	20.43 \pm 2.95	16.37 \pm 5.33	37.66 \pm 7.86	55.82 \pm 4.02	65.40 \pm 1.37	59.09 \pm 9.13	62.57 \pm 6.85
Acc	Overall	19.78 \pm 2.07	23.06 \pm 2.10	39.28 \pm 5.07	47.72 \pm 4.19	53.22 \pm 1.35	58.36 \pm 4.13	59.84 \pm 4.17
	C1 Subset	20.61 \pm 2.23	16.18 \pm 1.59	36.55 \pm 4.05	42.14 \pm 5.26	43.61 \pm 2.89	59.89 \pm 1.09	59.62 \pm 1.96
	C2 Subset	18.86 \pm 2.78	30.79 \pm 2.68	42.33 \pm 7.24	53.93 \pm 3.91	63.93 \pm 1.28	56.66 \pm 8.08	60.08 \pm 7.05

Table 7.3: Performance of different models on the five-class classification of Yelp-5. We use five paraphrases for the prompt (in Appendix A.6.1.3), and report the average performance with the standard deviation.

		Random	GPT-2 XL	LLaMa-7B	Alpaca-7B	GPT-3	GPT-3.5	GPT-4
F1	D=C1, P=C1	20.47 \pm 2.47	6.12 \pm 2.77	55.16 \pm 7.16	52.74 \pm 5.04	38.36 \pm 6.66	60.62 \pm 3.24	54.23 \pm 4.17
	D=C1, P=C2	20.26 \pm 2.31	15.98 \pm 3.59	36.22 \pm 8.95	35.10 \pm 5.40	54.44 \pm 1.64	52.85 \pm 6.01	58.58 \pm 3.33
	D=C2, P=C1	22.35 \pm 3.02	31.98 \pm 8.66	56.69 \pm 8.46	54.74 \pm 14.26	74.64 \pm 3.06	78.18 \pm 1.21	72.52 \pm 3.68
	D=C2, P=C2	20.35 \pm 2.18	48.50 \pm 7.66	66.82 \pm 7.90	71.22 \pm 3.99	77.09 \pm 1.33	78.16 \pm 1.81	76.80 \pm 1.36
Acc	D=C1, P=C1	19.60 \pm 2.67	12.96 \pm 2.91	58.23 \pm 5.30	55.81 \pm 4.00	43.16 \pm 6.20	60.39 \pm 3.25	54.06 \pm 3.97
	D=C1, P=C2	19.68 \pm 2.46	22.61 \pm 5.73	43.07 \pm 8.18	41.45 \pm 4.06	56.36 \pm 1.94	53.11 \pm 5.89	58.68 \pm 3.45
	D=C2, P=C1	20.97 \pm 3.19	43.51 \pm 6.30	57.54 \pm 7.21	54.82 \pm 12.59	76.60 \pm 3.23	77.38 \pm 1.43	70.69 \pm 3.93
	D=C2, P=C2	19.03 \pm 2.18	51.59 \pm 5.09	68.16 \pm 9.24	71.83 \pm 4.33	76.70 \pm 1.35	78.92 \pm 1.31	75.38 \pm 1.70

Table 7.4: Performance on Yelp using the *causal prompts* on the two causal subsets. We experiment all combinations of data nature (“D”) and causal prompt type (“P”), and report the average performance across the five paraphrases for each prompt, with the standard deviation.

binary classification datasets is saturated (Poria et al., 2023; Yang et al., 2019), we use the 5-way classification version of the SA datasets when applicable.

Since we need to utilize the emotion arc, we keep only reviews with at least five sentences, after sentence tokenization using the Spacy package (Honnibal and Montani, 2017). We apply this filtering above on the test set of the Yelp dataset, the English test of Amazon, and the unsplit entire dataset of App Review. We report the statistics of remaining samples in Table 7.1.

To obtain the emotion arcs, we calculate the sentiment score of each sentence by the sentiment-analysis pipeline¹ from Huggingface (Wolf et al., 2020).

Causal Discovery For each input sample, we process them as in Table 7.2, namely first obtaining the sentence-level sentiments to form the emotion arc, and then calculating the alignment scores λ_1 and λ_2 for each causal process, respectively. We consider an example as *dominated* by the causal process C_i if the alignment score λ_i is more optimal than the other.

We report the resulting statistics in Table 7.1. For each dataset, we describe their overall statistics, as well as the statistics of data with the underlying causal process of C1, and

¹<https://huggingface.co/distilbert-base-uncased-finetuned-sst-2-english>. We map its output value in 0–1 to the 5-class labels, by converting a value in 0–0.2 to the class label of 1, 0.2–0.4 to label 2, ..., and 0.8–1 to label 5.

that of C2. We can see that Yelp and Amazon have an almost balanced split of C1 and C2, while App Review has 61% C2 data compared to 39% C1 data. See Appendix A.6.2.2 for an additional visualization of the λ_1 - λ_2 distribution across the 1K data points.

7.4 Improving Sentiment Classifiers with Causal Alignment

Using our proposed causal discovery method, we have identified two distinct subsets with their corresponding causal processes C1 and C2. Now, we address the last practical question proposed in Section 7.2.3:

Can causal alignment help us improve SA in the era of LLMs?

Specifically, we take the commonly used approach in the era of LLMs, i.e., prompting pre-trained LLMs for the SA task, and look into how alignment with the underlying causal process could help SA performance. We answer the following three subquestions in this section:

- Q1. Using the standard SA prompt, do models perform differently on C1/C2 data? (Section 7.4.1)
- Q2. Does it help if we make the prompt aware of the underlying causality, i.e., use causal prompts? (Section 7.4.2)
- Q3. When prompted causally, do LLMs really understand the causal processes? (Section 7.4.3)

7.4.1 Q1: Do Models Perform Differently on C1/C2 Data?

Experimental Setup The first question is whether models perform differently on data with the causal nature of C1 or C2. We use the subsets identified by our psychologically-grounded causal discovery, and test a variety of available autoregressive LLMs, including the open-weight GPT-2 (Radford et al., 2019), LLaMa (Touvron et al., 2023), and Alpaca (Taori et al., 2023b); as well as the closed-weight models with OpenAI API, the instruction-tuned GPT-3 (text-davinci-002) (Brown et al., 2020; Ouyang et al., 2022), GPT-3.5 (gpt-3.5-turbo-0613), GPT-4 (gpt-4-0613) (OpenAI, 2023). We also add a random baseline which uniformly samples the label space for each input.

We use the standard prompt formulation for SA in the format of “[Instruction] Review Text: {x}\n Label:”. The experiments are on a set of randomly selected 1K samples from the test set of Yelp-5 (Zhang et al., 2015b), due to the time- and cost-expensive inference of the above LLMs. (E.g., LLaMa/Alpaca takes 96 GPU hours to run.) See more experimental details in Appendix A.6.1.1.

Results We show the performance of the six LLMs in Table 7.3, and report the F1 and accuracy across the five-class classification on Yelp-5. We can see that the existing LLMs perform the best on the subset with the causal process C2, implying that the decision pattern of LLMs is closer to the Fast Thinking system, which takes the peak-end average of the emotion arc.

7.4.2 Q2: Do Causal Prompts Help?

Designing Causal Prompts Inspired by the fact that models perform differently on C1/C2 data, our next question is, will it help if we directly give a hint to the LLMs about the underlying causal graph?

Prompt Design	
C1	As a customer writing a review, I initially <i>composed</i> the following feedback: "[review]" After carefully considering the facts, I selected a star rating from the options of "1", "2", "3", "4", or "5". My final rating was:
C2	As a customer writing a review, I initially <i>selected</i> a star rating from the options "1", "2", "3", "4", and "5", and then provided the following explanations in my review: "[review]"
The review <i>clarifies</i> why I gave a rating of	

Table 7.5: Causally-aware prompts describing the SA task in contexts with the C1 and C2 causal graphs.

To this end, we propose the idea of *causal prompts*, which are prompts that describe the causal story behind the input and output variables. We list our designed prompts for the C1 and C2 stories in Table 7.5.

Results We report the performance for all combinations of the dataset natures and prompt natures in Table 7.4, where we find that the most-performant setting uses Prompt C2 on the data subset with the same causal nature, C2. This alignment leads to the best performance across almost all models by both F1 and accuracy. On the C2 data, we also see that Prompt C2 outperforms the standard SA prompt in Table 7.3 by a substantial margin, such as 32.13 F1 points increase for GPT-2, and 14.23 F1 points increase for GPT-4.

However, although Prompt C2 shows a strong performance, the other causal prompt, i.e., Prompt C1, does not always help the data subset C1 in all cases, from which we raise a further question – how well do LLMs really mechanistically understand our prompts? We explore this question in the next section.

7.4.3 Q3: Can LLMs Correctly Capture the Causal Stories in the Prompts?

Although the proposal of the two causal prompts is intuitive for humans, we still need to inspect whether LLMs are able to understand them correctly.

Method Mechanistically, for a model to solve SA for the causal process C1 correctly, it needs to treat the sentence-level sentiments across all sentences *equally*; and for a model to solve SA for the causal process C2 correctly, it needs to pay *more* attention to the peak and end sentiments on the emotion arc.

Targeting the two mechanisms, we use causal tracing (Meng et al., 2022) to attribute the final sentiment prediction to the source sentences in the input. Briefly, causal tracing uses causal mediation analysis (Pearl, 2001) to quantify the causal contribution of the internal neuron activations of a model to its final prediction (Vig et al., 2020a). We use causal tracing to inspect the causal effects of the hidden states on the model prediction, using

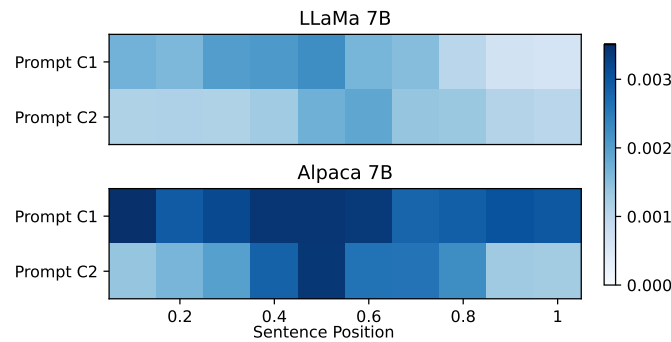


Figure 7.2: Causal attribution in LLaMa-7B and Alpaca-7B, showing how much each sentence contributes to the prediction probability.

the open-weight models, LLaMa and Alpaca. We use the causal effects of the first-layer neurons for each sentence, which we aggregate to obtain the final prediction. See implementation details in Appendix A.6.1.2.

Results We plot the causal attribution results of how much each sentence contributes to the final prediction in Figure 7.2. Here, the ideal behavior of the models is that Prompt C1 should trigger *uniform* attention over all the sentences, which is roughly observed through the more even shades of color of the “Prompt C1” row than the “Prompt C2” row in Figure 7.2 in the row of “Prompt C1”. On the other hand, Prompt C2 should trigger *more* attention to the sentences corresponding to the peak and end sentiments. For this, we see the models have high attention to the middle sentence, as in the “Prompt C2” row in Figure 7.2. Additionally, the average causal effect of the peak sentence on predictions under Prompt C2 is 0.0013 for LLaMa and 0.0029 for Alpaca, quantitatively aligning with our expectation that the peak sentence would show a high contribution. Nonetheless, note that no model sufficiently attends to the end sentence under Prompt C2. This implies that they do not fully grasp the expected contribution pattern of the peak-end rule, missing the significant role of the end sentence.

7.5 Related Work

SA The task of SA aims to identify the sentiment given a piece of text. It has a rich history originating from subjectivity analysis (Wiebe, 1994; Hatzivassiloglou and Wiebe, 2000), and developing rapidly with the availability of large opinionated online data such as reviews with star ratings (Turney, 2002; Nasukawa and Yi, 2003; Zhang et al., 2015b; Keung et al., 2020, *inter alia*). Most literature on SA focuses on building computational models, from using traditional linguistic rules (Hatzivassiloglou and McKeown, 1997; Choi and Cardie, 2008), to the application of machine learning methods, from traditional naive bayes and support vector machines (Pang et al., 2002; Moraes et al., 2013; Tan et al., 2009), to early deep learning models (Socher et al., 2013; Kim, 2014; Xing et al., 2020), and finally entering the era of LLMs (Hoang et al., 2019; Raffel et al., 2020; Yang et al., 2019).

Psychology and Affective Science In the study of emotion, or affect science (Salovey and Mayer, 2004; Barrett, 2006; Feinstein, 2013), previous work finds that not only the

emotion people perceive influences or prime how they communicate in the moment (Barrett, 2006), but language can also influence emotion, which can be observed in functional magnetic neuroimaging (Satpute et al., 2013), and also experiments showing the act of self-reporting the emotion in writing can change the physical reaction to the emotion (Kassam and Mendes, 2013). In his seminal work, Kahneman (2011) uses the two systems of thinking to reveal the mechanisms of how people come up with their sentiment, where fast thinking conforms to the peak-end rule (Kahneman et al., 1993), and slow thinking is more reflective of the overall sentiment.

Cause-Effect Distinction Distinguishing the cause from effect based on observational data is a long-standing and fundamental problem in causality (Hoyer et al., 2008; Zhang and Hyvärinen, 2009; Janzing, 2019). Existing methods to address this problem are based on statistics (Hoyer et al., 2008; Peters et al., 2010; Shajarisales et al., 2015; Mooij et al., 2014), physics (Janzing, 2007; Janzing et al., 2016), information theory (Janzing et al., 2012; Chaves et al., 2014; Mejia et al., 2022), and algorithmic complexity (Janzing and Schölkopf, 2010; Jin et al., 2021c). However, we are the first to look at the rich nature of NLP datasets, and directly approach the difference in the causal and anticausal mechanisms grounded in interdisciplinary insights.

As for our causal prompts, the most similar studies are the non-causally-grounded explorations for prompt tuning, such as by varying the patterns of masked language modeling (Schick and Schütze, 2022) and using the noisy channel method (Min et al., 2022). However, these studies are not aware of the underlying causal processes, thus neglecting the connection of prompts with the causal nature of data, and also the explicit causal story of the sentiment-review relation.

7.6 Conclusion

In conclusion, we have formulated the task of SA into a prediction problem and a causal discovery problem. We first identified the cause-effect relation among existing SA datasets, namely whether the review primes the rating, or the sentimental judgment primed the review writing process. To achieve this causal discovery, we obtain insights from existing psychology studies, namely aligning the above two causal processes with the famous Fast Thinking and Slow Thinking systems, with their distinct qualitative signals. Given the causal understanding of the dataset, we further improve the performance of LLMs on SA using our proposed causal prompts. Our research paves the way for more causally-aware future research in SA.

Limitations and Future Work

This study has several limitations. First, the rapid progression of LLMs makes it challenging to keep up with all newly proposed models and architectures. Since our work covers only a set of recent LLMs at the time of this study, we encourage future research to apply our methods to additional LLMs and other SA datasets.

Although our study is grounded in well-established psychological theories, there remains the possibility that new theories could emerge, necessitating updates to the calculation of the λ values for the two causal processes. However, the causal processes identified in this

work appear plausible, as evidenced by the effectiveness of the causally aligned prompts in improving language model performance.

Regarding the causal graph formulation, we focus on basic bivariate causal graphs, but future work could include more variables, such as confounders, mediators, and colliders.

The nature of this work is to introduce a paradigm shift for SA, and formulate the task differently. Therefore, we see lots of space for future extensions, such as to explore the causal nature of SA in different settings, different languages, and also aspect-based sentiment analysis (Pontiki et al., 2014; Xing et al., 2020; Hua et al., 2023).

Ethical Considerations

Regarding data concerns and user privacy, our study employs several established NLP datasets, and the examples we cite do not include sensitive user information.

Concerning potential stakeholders and misuse, this research primarily introduces a new perspective on the SA task. A possible negative impact concerns the general application of SA, which could be used to analyze user mentality for surveillance or fraudulent purposes. We acknowledge that studies on SA inherently involve these risks, and we firmly oppose the misuse of SA models in such contexts.

Part IV

Causality for Text-Based Computational Social Science

Mining the Cause of Political Decision-Making from Social Media: A Case Study of COVID-19 Policies across the US States

Mining the causes of political decision-making is an active research area in the field of political science. In the past, most studies have focused on long-term policies that are collected over several decades of time, and have primarily relied on surveys as the main source of predictors. However, the recent COVID-19 pandemic has given rise to a new political phenomenon, where political decision-making consists of frequent short-term decisions, all on the same controlled topic—the pandemic. In this paper, we focus on the question of how public opinion influences policy decisions, while controlling for confounders such as COVID-19 case increases or unemployment rates. Using a dataset consisting of Twitter data from the 50 US states, we classify the sentiments toward governors of each state, and conduct controlled studies and comparisons. Based on the compiled samples of sentiments, policies, and confounders, we conduct causal inference to discover trends in political decision-making across different states. Our code and data are publicly available at <https://github.com/zhijing-jin/covid-twitter-and-policy>.

8.1 Introduction

Policy responsiveness is the study of the factors that policies respond to (Stimson et al., 1995). One major direction is that politicians tend to make policies that align with the expectations of their constituents, in order to run successful re-election in the next term (Canes-Wrone et al., 2002).

An overview of existing studies on policy responsiveness reveals several patterns, summarized in Table 8.1. First, most work focuses on the *long-term* setting, where the policies are collected over a span of several decades, e.g., Caughey and Warshaw (2018)’s collection of public opinion surveys and state policymaking data over 1936-2014, and Lax and Phillips (2009)’s collection of public opinion polls and gradual policy changes over 1999-2008. Second, the data sources of existing studies are mostly surveys and polls, which can be time-consuming and expensive to collect (Lax and Phillips, 2012). Third, the resulting data are often of relatively small sizes, for both the number of policies and the number of public opinion.

Different from previous work on long-term policies, our work focuses on the special case of COVID pandemic, during which political leaders make a number of frequent, short-term policies on the same topic: social distancing. Moreover, instead of collecting surveys, we use Twitter to collect public opinion, which is instant, costless, and massive, e.g., tril-

	Previous Work	This Work
Policy Type	Long-term, gradual (over decades)	Short-term (weekly/-monthly)
Policy Sparsity	Less policies on the same topic	Many policies on the same topic across states
Data Source	Surveys	Trillions of tweets
Data Collection	–	NLP & Causality

Table 8.1: Comparison of the characteristics and paradigms of existing work versus our work.

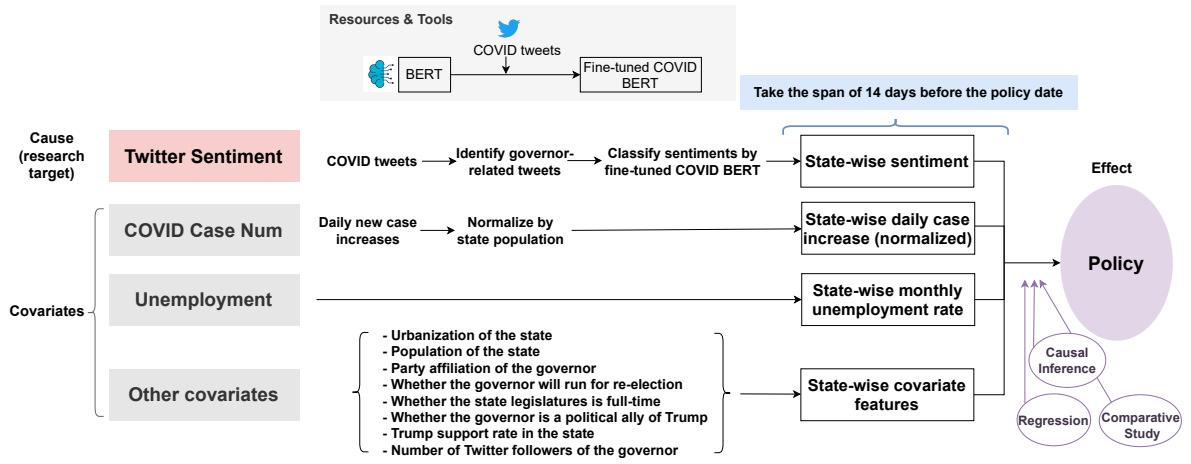


Figure 8.1: The data collection pipeline and architecture of our system to predict the state-wise COVID policies.

lions of data points. We limit our scope to US policies because the 50 states provide abundant policy data, and a good background for both controlled groups and comparative studies.

We present one of the first efforts to address policy responsiveness for short-term policies, namely the causal impact of public Twitter sentiments on political decision-making. This is distinct from existing studies on COVID policies that mostly explore the impact of policies, such as predicting public compliance (Grossman et al., 2020; Allcott et al., 2020; Barrios and Hochberg, 2020; Gadarian et al., 2021; DeFranza et al., 2020). Specifically, since governors have legislative powers through executive orders, we focus our study on each state governor’s decisions and how public opinion towards the governor impacts their decisions. For example, governors that optimize short-term public opinion are more likely to re-open the state even when case numbers are still high.

Our workflow is illustrated in Figure 8.1. We start by collecting 10.4M governor-targeted COVID tweets, which we annotate for sentiment with a BERT-based classifier. Next, we

annotate 838 social distancing policies and collect data on ten potential confounders such as average daily case increases or unemployment rates. Finally, we conduct multiple analyses on the causal effect of Twitter sentiment on COVID policies. For interpretability, we first use a multivariate linear regression to identify correlations of sentiments and policies, in addition to considering all the confounders. We also use do-calculus (Pearl, 1995) to quantify the causal impact of Twitter sentiment on policies. We also conduct cross-state comparisons, cross-time period analysis, and multiple other analyses.

The main contributions of our work are as follows. First, we compile a dataset of public opinion targeted at governors of the 50 US states with 10.4M tweets. Second, we annotate a dataset of 838 COVID policy changes of all 50 states, along with data of ten confounders of each state. Third, we conduct regression analyses and causal analyses on the effect of Twitter sentiment on policies. Finally, we implement additional fine-grained analyses such as cross-state comparisons, cross-time period analysis, and multiple other analyses.

8.2 Related Work

Policy Responsiveness Policy responsiveness (i.e., public opinion—*causes*→policies) is an active research field in political science, where people study how policies respond to different factors (Stimson et al., 1995). Studies show that policy preferences of the state public can be a predictor of future state policies (Caughey and Warshaw, 2018). For example, Lax and Phillips (2009) show that more LGBT tolerance leads to more pro-gay legislation in response. Most policies and public opinion studied in existing literature are often long-term and gradual, taking several decades to observe (Lax and Phillips, 2009, 2012; Caughey and Warshaw, 2018).

Crisis Management Policies Another related topic is crisis management policies, where most studies focus on the reverse causal problem of our study – how crisis management policies impact public opinion (i.e., policies—*causes*→public opinion). A well-known phenomenon is the rally “round the flag” effect, which shows that during a crisis, there will be an increased short-run public support for the political leader (Mueller, 1970, 1973; Baum, 2002), due to patriotism (Mueller, 1970; Parker, 1995), lack of opposing views or criticism (Brody and Shapiro, 1989), and traditional media coverage (Brody, 1991).

To the best of our knowledge, there is not much research on how public opinion influence policies (i.e., public opinion—*causes*→policies) during a crisis. Our work is one of the few to address this direction of causality.

COVID Policies There are several different causal analyses related to COVID-19 policies, although different from our research theme. Existing studies focus on how social distancing policies mitigate COVID spread (i.e., policies—*causes*→pandemic spread) (Kraemer et al., 2020), what features in public attitudes impact the compliance to COVID policies (i.e., public attitudes/ideology—*causes*→policy compliance) (Grossman et al., 2020; Allcott et al., 2020; Barrios and Hochberg, 2020; Gadarian et al., 2021), how policies change the public support of leaders (i.e., policy—*causes*→public support). Bol et al. (2021); Ajzenman et al. (2020), how pandemic characteristics affect Twitter sentiment (Gencoglu and Gruber, 2020), and how political partisanship impacts policies (i.e., partisanship—*causes*→policy designs) (Adolph et al., 2021). However, there is no existing work using public sentiments

	Positive	Neutral	Negative
Percentage	15.8%	36.5%	47.7%
Length	15.51	12.21	16.39
Topics	we, support, thank, great, gov- ernors, covid, action	people, masks, covid, cases, state, today, total	cases, state, covid, close, deaths, people, trump
4- Grams	- great governors responded executive - responded executive action promptly - quickly , support americans	- positive patients nursing homes - governors ordered covid positive - today 's update numbers	- covid patients nursing homes - america 's governors forced - covid patients nursing homes
Example	"I am a small business owner, we kept health insurance for the furloughed staff of my two restaurants, month after month, even while one restau- rant was closed and the other only has limited service. Why? Because I have a conscience. We are in a pandemic."	"Today: @GovInslee 3 pm news conference on WA's coronavirus response. Inslee to be joined by state schools chief. Your daily #covid19 updates via @seattletimes"	"And the politicians that are doing the conditioning are out, maskless, celebrat- ing with their family and friends... @GavinNewsom Glad I never once fell for it. Covid-19 was always just a power-grab for politicians"

Table 8.2: Label distribution (Percentage), average number of words per tweet (Length), topics extracted by LDA topic modeling (Blei et al., 2003), top 4-grams, and examples of positive, neutral, and negative tweets.

(e.g., from social media) to model COVID policies.

Opinion Mining from Social Media Social media, such as Twitter, is a popular source to collect public opinions (Thelwall et al., 2011; Paltoglou and Thelwall, 2012; Pak and Paroubek, 2010; Rosenthal et al., 2015). Arunachalam and Sarkar (2013) suggest that Twitter can be a useful resource for governments to collect public opinion. Existing usage of Twitter for political analyses mostly targets at election result prediction (Beverungen and Kalita, 2011; Mohammad et al., 2015; Tjong Kim Sang and Bos, 2012), and opinion towards political parties (Pla and Hurtado, 2014) and presidents (Marchetti-Bowick and Chambers, 2012). To the best of our knowledge, this work is one of the first to use Twitter sentiment for causal analysis of policies.

8.3 Governor-Targeted Public Opinion

To investigate the causality between public opinion and each state governor's policy decisions, we first describe how we mine public opinion in this Section; we then describe the process we use to collect policies and other confounders in Section 8.4.

We collect governor-targeted public opinion in two steps: (1) retrieve governor-related COVID tweets (Section 8.3.1), and (2) train a sentiment classification model for the COVID tweets and compile sentiments towards governors (Section 8.3.2).

8.3.1 Retrieve Governor-Related COVID Tweets

We use the COVID-related tweet IDs curated by Chen et al. (2020).¹ Chen et al. (2020) identified these tweets by tracking COVID-related keywords and accounts. We provide the list of keywords and accounts they used in Appendix A.7.1.1. We hydrate the tweet IDs to obtain raw tweets using an academic Twitter Developer account. This process took several months to complete, and resulted in a dataset of 1.01TB. The retrieved 1,443,871,617 Tweets span from January 2020 to April 2021.

Since this study focuses on governor’s policy decision-making process, we focus on the public opinion that are more directly related to the governors. Specifically, we focus on tweets that tagged, replied to, or retweeted state governors. We obtain 10,484,084 tweets by this filter. On average, each of the 50 states has about 209K tweets that address the state governor. The rationale of this filter is that the governors and their teams are likely to have directly seen (a portion of) these tweets, since they showed up in governor’s Twitter account.

8.3.2 Classify Sentiments towards Governors

Existing studies on COVID Twitter sentiment analysis (Manguri et al., 2020; Kaur and Sharma, 2020; Vijay et al., 2020; Chakraborty et al., 2020a; Singh et al., 2021a) mostly use TextBlob (Loria, 2018), or some simple supervised models (Machuca et al., 2021; Kaur et al., 2021; Mansoor et al., 2020).

For our study, we use the state-of-the-art BERT model pretrained on COVID tweets by Müller et al. (2020).² We finetune this pretrained COVID BERT on the Twitter sentiment analysis data from SemEval 2017 Task 4 Subtask A (Rosenthal et al., 2017). Given tweets collected from a diverse range of topics on Twitter, the model learns a three-way classification (positive, negative, neutral). In the training set, there are 19,902 samples with positive sentiments, 22,591 samples with neutral sentiments, and 7,840 samples with negative sentiments.

We tokenize the input using the BERT tokenizer provided by the Transformers Python package (Wolf et al., 2020). We add [CLS] and [SEP] tokens at start and end of the input, respectively. The input is first encoded by the pretrained COVID BERT. Then, we use the contextualized vector C of the [CLS] token as the aggregate sentence representation. The model is finetuned on the classification task by training an additional feed-forward layer $\log(\text{softmax}(CW))$ that assigns the softmax probability distribution to each sentiment class.

Prior to training, we preprocess the tweets by deleting the retweet tags, and pseudonymising each tweet by replacing all URLs with a common text token. We also replace all unicode emoticons with textual ASCII representations. During training, we use a batch size of 32 and fine-tune for 5 epochs. We use a dropout of 0.1 for all layers, and the Adam optimizer (Kingma and Ba, 2017) with a learning rate of $1e-5$. Additionally, due to the specific nature of our classification task (i.e., mining opinion towards the governor), we add a post-processing step to classify a tweet as supportive of a governor (i.e., positive) if the tweet retweets a tweet from the governor’s official account.

¹COVID-related Tweet IDs: <https://github.com/echen102/COVID-19-TweetIDs>

²<https://huggingface.co/digitalepidemiologylab/covid-twitter-bert-v2>

Model Performance. We evaluate our model accuracy on two test sets. First, on the test set of SemEval 2017, our finetuned model achieves 79.22% accuracy and 79.29% F1. Second, we also evaluate our model performance on our own test set. Since the features of general tweets provided in SemEval 2017 might differ from COVID-specific tweets, we extracted 500 random tweets from the Twitter data we collected in Section 8.3.1. We asked a native English speaker in the US to annotate the Twitter sentiment with regard to the state governor that the tweet addresses. The annotator has passed a small test batch before annotating the entire test set.

We use the TextBlob classifier as our baseline, since it is the most commonly used in existing COVID Twitter sentiment analysis literature. On our test set’s three-way classification, the TextBlob baseline has 23.35% accuracy and 16.67% weighted F1. Our finetuned BERT classifier has 60.23% accuracy and 62.31% weighted F1. Detailed scores per class is in Appendix A.7.1.3. When applying the sentiment classifier, we care more about whether the average sentiment over a time period is accurate, so we also turn the test set into groups of tweets each containing 20 random samples. The average mean squared error (MSE) for the average sentiment of each group is 0.03889 for the BERT model, and 0.22749 for the TextBlob model. We apply the finetuned COVID BERT classifier on the governor-related tweets we extracted previously. As listed in Table 8.2, among 10.4M tweets, 15.8% are positive, 36.5% neutral, and 47.7% negative.³

We use Latent Dirichlet Allocation (LDA) topic modeling (Blei et al., 2003) to extract key topics of each category. Typical topic words in positive tweets include “we,” “support,” “thank,” “great,” and “governors,” while negative tweets tend to mention more about “america’s governors forced ...” and support Trump, perhaps Trump’s tweets on “liberation.”

8.4 Collection of Policies and Confounders

We focus on state-wide social distancing policies, and collect 838 social distancing policies from 50 states over the period January 2020 – April 2021 (described in Section 8.4.1).

Since we want to focus on the causal effect of public sentiment on policy, we must control for possible confounding factors. In particular, case numbers and unemployment rates are potentially the most important confounders, the collection of which is introduced in Section 8.4.2. In addition, we also collect eight other potential confounders suggested by political science experts (described in Section 8.4.3). The collection process is illustrated in Figure 8.1.

8.4.1 Social Distancing Policy Annotation

We annotate the social distancing policies related to COVID for each of the 50 states in the US. For each state, the annotators are asked to go through the entire list of COVID-related executive orders from January 2020 to April 2021. In cases where the states do not use executive orders for COVID regulations, we also consider proclamations and state guidance on social distancing.

The policies are rated on a scale of 0 (loosest) - 5 (strictest). We provide guidance as to

³Note that label imbalance is commonly observed on Twitter data (Guerra et al., 2014).

the level of strictness that each number indicates, as detailed in Appendix A.7.1.2. Four annotators are asked to conduct the ratings. Since the annotation is very tedious, taking up to 3 hours per state, we do not conduct double annotations. Instead, given our original annotations (for which we score each policy based on its official legal document in PDF), we did a quick second pass by confirming that our scores roughly match the succinct 1~2-sentence textual summary of each policy provided by the Johns Hopkins Coronavirus Resource Center.⁴

8.4.2 Key Confounders: State-Level Case Numbers and Unemployment Rates

We collect COVID daily new confirmed case numbers from the open-source COVID database⁵ curated by the Kaiser Family Foundation. For a fair comparison across states, we normalize the case numbers by the population of the state. We retrieve the seasonally adjusted data of monthly unemployment rates for each state from the U.S. Bureau of Labor Statistics.⁶

8.4.3 Additional Confounders

For additional confounders, we collect both state data as well as governor features.

State Features. For state features, we collect the population⁷ and urbanization rate from US 2010 Census (Census Bureau, 2012).⁸ In addition, we also collect the last US presidential election returns of each state.⁹ Note that it is necessary to use pre-policy data, so we collect the presidential election returns from 2016 but not from 2020. For the presidential election returns, we obtain the percentage of votes for Donald Trump to indicate Trump's support rate.

Governor Features. For each governor, we collect their party affiliation, whether the governor will run for the next gubernatorial election,¹⁰ and whether the state legislatures are full-time or not, collected from National Conference of State Legislatures.¹¹ In addition, we also annotate whether the governor is a political ally of Trump or not. We conduct the annotation based on the background and past news reports of each governor. For corner cases, we quote additional evidence in our annotation, e.g., for republican governors who do not support Trump, and democratic governors who support Trump. We also collect the number of Twitter followers for each governor, since it might be correlated with how much attention the governor pays to the twitter reactions.

Table 8.3 lists the statistics of the confounder data we collected.

⁴Social distancing policy summaries: <https://coronavirus.jhu.edu/data/state-timeline>

⁵COVID case number data: <https://github.com/KFFData/COVID-19-Data>

⁶Monthly unemployment data: <https://www.bls.gov/web/laus/ststdsadata.zip>

⁷Population data: <https://www.census.gov/programs-surveys/decennial-census/data/tables.2010.html>

⁸Urbanization data: <https://www.icip.iastate.edu/tables/population/urban-pct-states>.

⁹Presidential election return data: <https://www.nytimes.com/elections/2016/results/president>

¹⁰For simplicity, we collect the pre-COVID data at the time point of January 2020, and do not consider the change of governorships in two states in early 2021.

¹¹<https://www.ncsl.org/>

Numerical Features			
	Mean (\pm std)	Min	Max
Daily Case Increase (%)	0.02 (\pm 0.02)	0.0	0.45
Unemployment Rate (%)	5.51 (\pm 3.25)	2.0	29.5
Urbanization (%)	73.58 (\pm 14.56)	38.7	95
Population (M)	12.94 (\pm 45.68)	0.57	325.38
Trump's Support Rate (%)	48.29 (\pm 11.93)	4	68
# Twitter Followers (K)	237 (\pm 458)	7	2596
Binary Features			
	Yes	No	
Gov Is Republican	26	24	
Will Run for Re-election	39	11	
Full-Time Legislatures	10	40	
Trump's Political Ally	22	28	

Table 8.3: Statistics of the ten confounders collected for policy prediction task.

8.5 Mining Decisive Factors of COVID Policies

Since we are interested in discovering the key factors that changes the decisions of policy-makers, we focus on the change of policies (e.g., changing from complete close down to reopening K-12 schools) rather than absolute values of the policy strictness. For each policy in state s on date t , we calculate the change Δpolicy as the difference of this policy from the previous policy that was issued.

Since sentiment may change rapidly and many policies are updated frequently during COVID, for each policy change Δpolicy , we focus on the average sentiment over the time span $(t - \Delta t, t)$ from Δt days prior to the policy date t . Here, we set $\Delta t = 14$ since many epidemiology reports are based on 14-day statistics, e.g., the 14-day notification rate.

When building the policy prediction model, we also need to account for confounders. For the confounders, most are static over time for a given state, except for the daily case increases and the unemployment rates that change over time, for which we take the average over the 14-day time span.

Based on the data above, we seek to answer the following questions: (Q1) What variables are indicative of policy changes?, and (Q2) What causal impact does sentiment have on the policies?

8.5.1 Q1: What Variables Are Indicative of Policy Changes?

To aim for interpretability, we choose a multivariate linear regression as our model, which is commonly used in political science literature on COVID policies (Grossman et al., 2020;

Allcott et al., 2020; Barrios and Hochberg, 2020; Gadarian et al., 2021). Specifically, we model the policy change Δpolicy as a function of all variables, including our main focus – Twitter sentiments – and all the confounders, which form in total 11 variables.¹²

Sentiment, Case Numbers, Unemployment Are Important The first experiment is to compare how well different combinations of input variables fit the policy change. We use mean squared error (MSE) as the measure of model capability.

When taking into consideration all variables, the model has an MSE score of 0.368. As a further step, we test whether a smaller number of inputs can achieve similar results. We find that when only taking three variables as inputs, the MSE is 0.369, which is 0.001 from the model taking in all variables. Among all combinations of three variables, the proposed three key variables, sentiment, case numbers, and unemployment rates, achieve the best performance of 0.369.

Note that it is reasonable that with rational decision-making, politicians consider the case numbers and unemployment rates when making COVID policies. The focus of this study is to show the *additional* effect of sentiment, the role of which is not explicitly pointed out in previous COVID policy research.

The Role of Non-Sentiment Variables First, given the presence of the sentiment variable in the model, we test the additional effect of non-sentiment variables. As shown in Table 8.4, case numbers and unemployment rate both lead to non-trivial improvement of the models, and unemployment is more important.

Additional Non-Sentiment Variables	MSE (↓)
Sentiment-only	0.618
+ Case	0.532
+ Unemp	0.407
+ Case, Unemp	0.369
+ Case, Unemp, Others	0.368

Table 8.4: The MSE of models taking as input the additional non-sentiment variables, such as case increases (Case), unemployment (Unemp), and other confounders (Others).

The Role of Sentiment Second, we look into the role of sentiment. We take the optimal 11-variable, 3-variable, and 2-variable models, and conduct ablation studies to inspect how much does sentiment contribute exclusively in Table 8.5.

We show that for each model, sentiment has a crucial impact of more than 0.032 on the model performance. Note that in linear regression, we do not need to explicitly disentangle the correlations within sentiments and other confounders – in Table 8.5, the effect of sentiment is demonstrated in addition to fitting all other variables that may contain correlations.

¹²For each input variable, we first normalize by adjusting mean to zero and standard deviation to 1.

Model	MSE (\downarrow)
11-Variable model	0.368
–Senti	Deterioration of 0.032
3-Variable model	0.369
–Senti	Deterioration of 0.032
2-Variable model	0.407
–Senti	Deterioration of 0.034

Table 8.5: Ablation study of sentiment for the optimal 11-, 3-, 2-variable models. Note that the 11-variable model is the full model taking in all variables.

8.5.2 Q2: What Causal Impact Does Sentiment Have on the Policies?

In the previous section, we investigated the most indicative variables of policies. The experiments indicate how important each variable is to the regression target, i.e., how well they serve as a predictor, although such *correlation* does not necessarily capture *causation*. In this section, we are interested in the causal impact of sentiment on policies, and we use causal inference methods to quantify the impact.

Formulation by Do-Calculus Formally, we are interested in the effect of a cause X (i.e., Twitter sentiment) on the outcome Y (i.e., policy change) in the presence of the confounder Z (i.e., case numbers, unemployment, etc.), as shown in Figure 8.2.

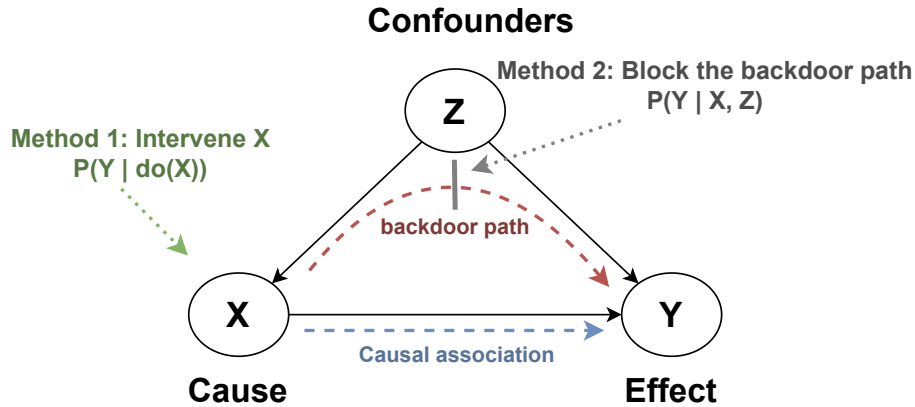


Figure 8.2: Backdoor Adjustment.

To formulate the causal impact, Pearl (1995) defines a language for causality called do-calculus, by which the causal impact of X on Y is formulated as the interventional distribution:

$$P(Y|\text{do}(X)), \quad (8.1)$$

where $\text{do}(X)$ refers to an intervention on the cause X .

Note that the interventional distribution $P(Y|\text{do}(X))$ may be different from the observational distribution $P(Y|X)$ in the presence of the confounder Z . Specifically, in the above Figure 8.2, there are two ways how X correlates with Y . The first is the causal path $X \rightarrow Y$, and the second is the backdoor path $X \leftarrow Z \rightarrow Y$.

There are two ways to account for the backdoor path: Method 1 needs to intervene on X , e.g., create a counterfactual situation where all confounders are the same but the Twitter sentiment can be set to negative vs. positive. In our study of Twitter opinion on COVID policies, this is not a feasible experiment to conduct, due to the fundamental problem of causal inference (Rubin, 1974; Holland, 1986) (namely, for each sample i , we are usually only able to observe one value of X but not both). The other method, backdoor adjustment, circumvents the problem, which will be introduced in the following.

Backdoor Adjustment The key challenge in the above causal inference is that we need to account for the confounder Z . Backdoor adjustment (Pearl, 1995) presents an approach to estimate the causal impact of X on Y by using only *observational* data. Basically, we need to block all backdoor paths by conditioning on nodes that can break the unwanted connections between X and Y . Moreover, these nodes should not contain any descendants of X . In our case, we condition on the confounder Z , and turn the interventional distribution into the observational distribution:

$$P(Y|\text{do}(X)) = \sum_Z P(Y|X, Z)P(Z) . \quad (8.2)$$

The causal impact of X (i.e., positive or negative sentiment) on Y (i.e., policy change) be-

$$\begin{aligned} \beta &= \mathbb{E}[Y|\text{do}(X = 1)] - \mathbb{E}[Y|\text{do}(X = -1)] \\ \text{comes} \quad &= \sum_Z (\mathbb{E}[Y|X = 1, Z] - \mathbb{E}[Y|X = -1, Z])P(Z) \\ &= \mathbb{E}_Z[\mathbb{E}[Y|X = 1, Z] - \mathbb{E}[Y|X = -1, Z]] . \end{aligned} \quad (8.3)$$

Results We apply Eq. (8.3) to all states using a 10-dim vector Z that encodes all confounders.¹³ Then we rank the states by β values, which represents the causal impact of sentiment on the state policies.

Top 5 States with Large β		Top 5 States with Small $ \beta $	
State	β Value	State	β Value
Colorado	4.292	Arizona	0.053
Massachusetts	1.157	West Virginia	0.030
Florida	1.124	Pennsylvania	0.023
Texas	1.095	Nebraska	-0.001
South Dakota	1.057	Alabama	-0.065

Table 8.6: Top five states with the largest β values, and the β values that are closest to zero.

In Table 8.6, we show the top five states with highest β values, and five states with β values that are the closest to zero. The higher the β value, there exists more alignment between people’s sentiment and the state policy strictness in the state.

There are some associations between our results and real-world patterns. For instance, among the top five states in Table 8.6, Colorado’s high β value reflects its Democratic

¹³Due to length restrictions, please refer to the arXiv version of our paper for additional implementation details of the backdoor adjustment.

governor's large net favorable rating compared to the Republican politicians.¹⁴ Massachusetts also has a high governor approval rate, and most people support the COVID policies. The three Republican states, South Dakota, Texas, and Florida, also have high β , but they are in a different scenario. The loose policies in all these states are in line with general sentiment across the states to refuse restrictions.

8.6 Fine-Grained Analyses

8.6.1 Early-Stage vs. Late-Stage Decisions

Since the COVID pandemic is an unprecedented situation, it is likely that in early stages of the pandemic, politicians tend to rely on their pre-judgements, and as time goes on, they form a better understanding of the situation and adjust their reaction towards the public opinion. We compare the causal impact of sentiment on policies in the first three months of the outbreak (i.e., from March to June 1, 2020) and afterwards (i.e., from June 1, 2020 to now). Table 8.7 shows that the states with the most changes in β are Montana, Washington, Georgia, Tennessee, and Indiana.

State	Change in β before and after June 1
Montana	+9.39
Washington	+4.03
Georgia	+3.15
Tennessee	+2.94
Indiana	+2.53

Table 8.7: Top 5 states with the most change in the causal impact of sentiment on policies from March to June 1, 2020, versus from June 1, 2020 to April, 2021.

8.6.2 Cross-State Comparison

For cross-state comparison, we identify states that are similar in terms of the confounders, and then compare how different policies are a result of different public sentiments. For simplicity, we consider the two most important confounders, case numbers and unemployment rates. We evaluate the similarity matching on the two time series across different states by the dynamic time warping algorithm (Berndt and Clifford, 1994), and extract state pairs that are the most similar in terms of the confounders.

In Figure 8.3, we show an example pair of states, Mississippi (MS) and Georgia (GA), which have highly similar case numbers and unemployment rates at most time steps. Note that we use the New York (NY) state to show in contrast how the above pair is different from another unrelated state.

In the comparative study of MS and GA, they can be considered as counterfactuals for each other. In their policy curves, the policy strictness in MS responds to the COVID case numbers (e.g., the policies are stricter on the rising slope of case numbers), but the policies in GA remain loose even during the rising trends in July – August 2020, and November 2020 – January 2021. We look into the sentiment differences across the two states: For example, during November 2020 – January 2021, GA experienced a very low average

¹⁴For example, see this poll result by Colorado Poll reported by Denver Post.

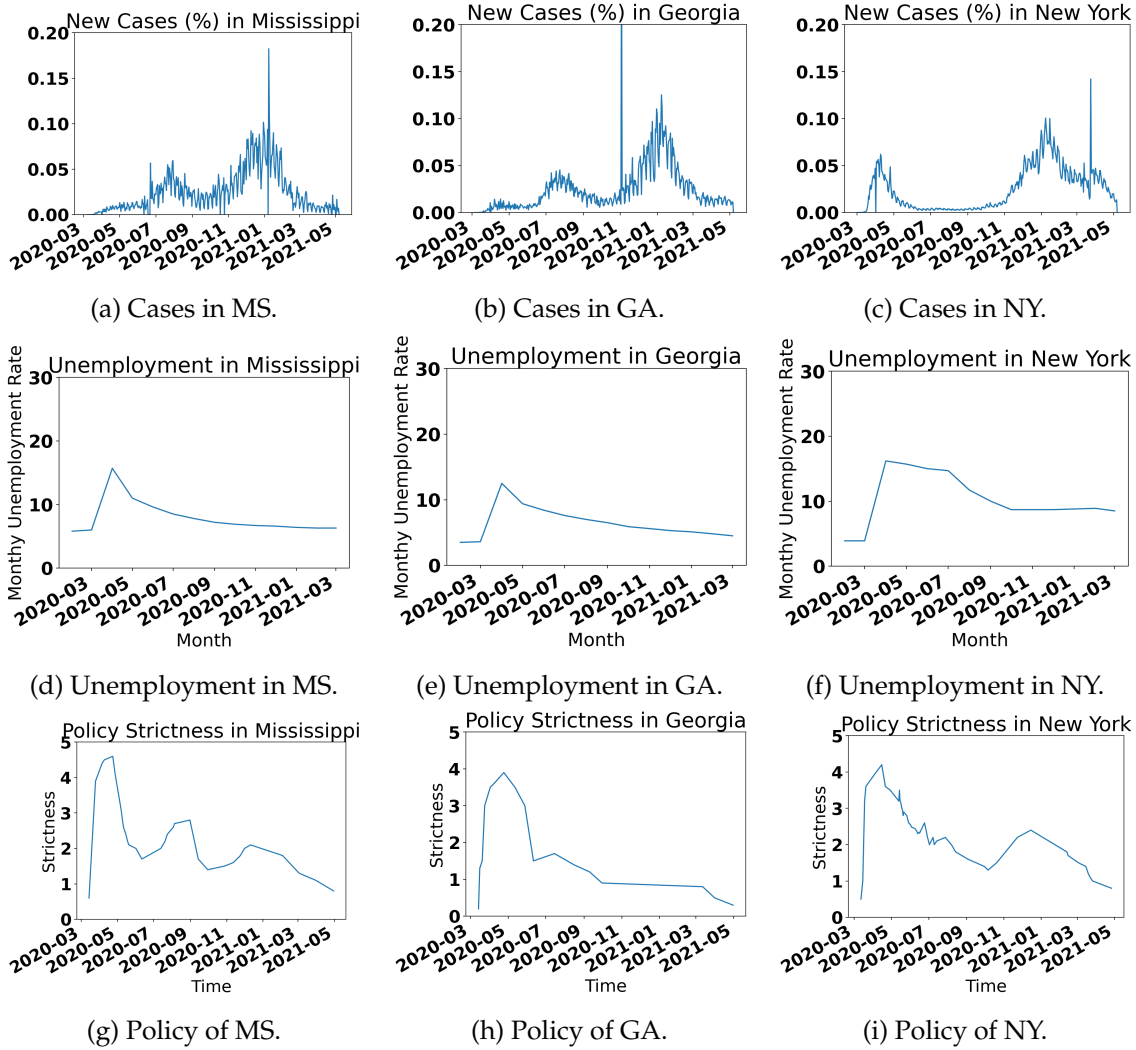


Figure 8.3: Comparative study of states. MS and GA is a pair of states with the most similar confounders, and NY is an irrelevant state to contrast how different MS and GA are from other states. Note that unemployment data is only available until March 2021.

sentiment of -0.58 in the $[-1, 1]$ scale, whereas MS experienced a milder sentiment of -0.04. By the controlled comparison, the more negative sentiment is the potential cause for looser policies in GA.

8.7 Additional Discussions

Fine-Grained Opinions behind the Sentiments. To further interpret why positive tweets usually lead to stricter social distancing policies (and negative tweets lead to looser policies), we look into the correlation of Twitter sentiment and the user's opinion towards social distancing policies. Note that usually it is not easy to directly get an unsupervised intent classifier on COVID specific tweets. Hence, we ask the annotators to classify the opinion on social distancing for the 500 tweets in our test set as supportive, against, and not related to social distancing. Among the tweets about social distancing with positive sentiment, 95.13% support social distancing. Among the tweets about social distancing

with negative sentiment, 69.38% are against social distancing and ask for the reopening of the state.

Additional Analyses. We put our additional analyses in Appendix A.7.2, including correlation across all variables, and alternative causal analysis models such as difference-in-differences (Abadie, 2005), and continuous-valued propensity score matching (Hirano and Imbens, 2004; Bia and Mattei, 2008).

Limitations. There are several limitations of this study. For example, a common limitation of many causal inference settings is the uncertainty of hidden confounders. In our study, we list all the variables that we believe should be considered, but future studies can investigate the effect of other confounders.

Another limitation is the accuracy of the Twitter sentiment classifier. Since the Twitter sentiment during COVID is very task-specific, modeling the sentiments can be very challenging. For example, our model often misclassifies “increased positive cases” as a positive sentiment. Another challenge is that some tweets refer to a url. These cases are difficult to deal with, and might be worth more detailed analyses in future studies.

In the data setting, one limitation is that for causal inference, modeling the whole time series is extremely challenging, so we empirically take the 14-day time span, which is a commonly used time span for many other COVID measures.

Future Work. This work is the first work to use NLP and causal inference to address policy responsiveness, and we explicitly measure the alignment of government policies and people’s voice. This signal can be very important for the government and decision-makers.

In future work, a similar approach can be used together with other variables (e.g., economic growth, participation in health/vaccination campaigns, well-being) to determine to which extent such people-government alignment relates to societal outcomes.

8.8 Conclusion

In this paper, we conducted multi-faceted analyses on the causal impact of Twitter sentiment on COVID policies in the 50 US states. To enable our study, we compile a large dataset of over 10 million governor-targeted COVID tweets, we annotate 838 state-level policies, and we collect data ten potential confounders such as daily COVID cases and unemployment rates. We use a multivariate linear regression and do-calculus to quantify both the correlation of Twitter sentiment as well as its causal impact on policies, in the presence of other confounders. To our knowledge, this is one of the first studies to utilize massive social media data on crisis policy responsiveness, and lays the foundation for future work at the intersection of NLP and policy analyses.

Ethical Considerations

Use of Data For the data used in this study, the COVID-related tweets are a subset of the existing dataset provided by Chen et al. (2020). Following the data regulations of Twitter,

we will not publicize the raw tweet text. If necessary, we can provide the list of tweet IDs to future researchers. For the policy strictness we annotated, we will open-source it since it is public information that can benefit societies affected by the pandemic, and has no privacy or ethical issues. For other confounding variables, the data are also public information.

Potential Stakeholders This research can be used for policy-makers or political science researchers. The research on causality between public opinion and political decision-making helps make policies more interpretable. One potential caveat is that there might be parties that maliciously manipulate sentiments on Twitter to affect politicians. A mitigation method is to control the flow of misinformation, terrorism and violent extremism on social media. The ideal use of the study is to reflect the process how a democracy system surveys the opinion from people, and makes policies that best balances people's long-term and short-term interests.

CausalCite: A Causal Formulation of Paper Citations

Citation count of a paper is a commonly used proxy for evaluating the significance of a paper in the scientific community. Yet citation measures are widely criticized for failing to accurately reflect the true impact of a paper. Thus, we propose CAUSALCITE, a new way to measure the significance of a paper by assessing the causal impact of the paper on its follow-up papers. CAUSALCITE is based on a novel causal inference method, TEXTMATCH, which adapts the traditional matching framework to high-dimensional text embeddings. TEXTMATCH encodes each paper using text embeddings from large language models (LLMs), extracts similar samples by cosine similarity, and synthesizes a counterfactual sample as the weighted average of similar papers according to their similarity values. We demonstrate the effectiveness of CAUSALCITE on various criteria, such as high correlation with paper impact as reported by scientific experts on a previous dataset of 1K papers, (test-of-time) awards for past papers, and its stability across various subfields of AI. We also provide a set of findings that can serve as suggested ways for future researchers to use our metric for a better understanding of the quality of a paper. Our code is available at <https://github.com/causalNLP/causal-cite>.

9.1 Introduction

Recent years have seen explosive growth in the number of scientific publications, making it increasingly challenging for scientists to navigate the vast landscape of scientific literature. Therefore, identifying a good paper has become a crucial challenge for the scientific community, not only for technical research purposes, but also for making decisions, such as funding allocation (Carlsson, 2009), research evaluation (Moed, 2006), recruitment (Gary Holden and Barker, 2005), and university ranking and evaluation (Piro and Sivertsen, 2016).

A traditional approach to recognize paper quality is peer review, a mechanism that requires large efforts, and yet has inherent randomness and flaws (Cortes and Lawrence, 2021; Rogers et al., 2023; Shah, 2022; Prechelt et al., 2018; Resnik et al., 2008). Moreover, the number of papers after peer review is still overwhelmingly large for researchers to read, leaving the challenge of identifying truly impactful research unaddressed. Another commonly used metric is citations. However, this metric faces criticism for biases, such as a preference for survey, toolkit, and dataset papers (Zhu et al., 2015; Valenzuela et al., 2015). Together with altmetrics (Wilsdon, 2016), which incorporates social media attention to a paper, both metrics also bias towards papers from major publishing countries (Rungta et al., 2022; Gomez et al., 2022), with extensive publicity and promotion, and authored by established figures.

To provide a more equitable assessment of paper quality, we employ the causal inference

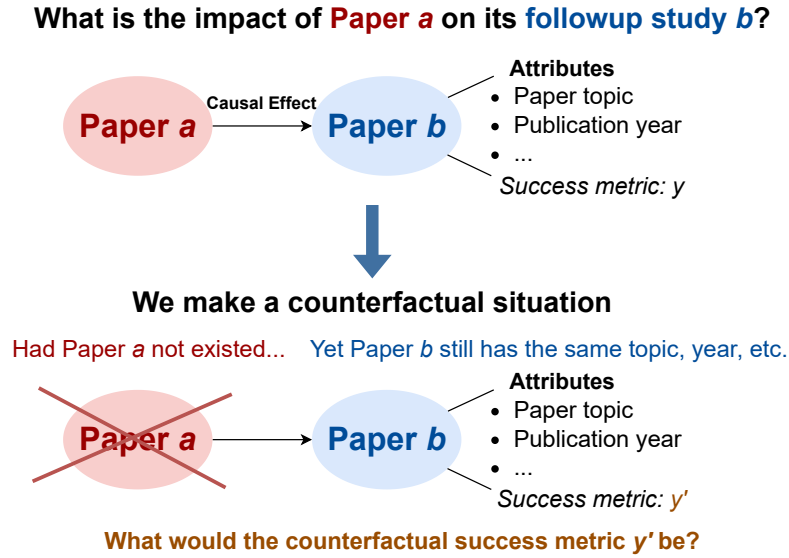


Figure 9.1: An overview of our research question.

framework (Hernán and Robins, 2010) to quantify a paper’s impact by how much of the academic success in the follow-up papers should be *causally attributed* to this paper. We introduce CAUSALCITE, an enhanced citation based metric that poses the following *counterfactual* question (also shown in Figure 9.1): “*had this paper never been published, what would have happened to its follow-up studies?*” To compute the causal attribution of each follow-up paper, we contrast its citations (the treatment group) with citations of papers that address a similar topic, but are not built on the paper of interest (the control group).

Traditionally, this problem is solved by using the matching method (Rosenbaum and Rubin, 1983) in causal inference, which discretizes the value of the confounder variable, and compares the treatment and control groups with regard to each discretized value of the confounder variable. However, this approach does not apply when the confounder variable is high-dimensional, e.g., text data, such as the content of the paper. Thus, we improve the matching method to adapt for textual confounders, by marrying recent advancement of large language models (LLMs) with traditional causal inference. Specifically, we propose TEXTMATCH, which uses LLMs to encode an academic paper as a high-dimensional text embedding to represent the confounders, and then, instead of iterating over discretized values of the confounder, we match each paper in the treatment group with papers from the control group with high cosine similarity by the text embeddings.

TEXTMATCH makes contributions in three different aspects: (1) it relaxes the previous constraint that the confounder variable should be binned into a limited set of intervals, and makes the matching method applicable for high-dimensional continuous variable type for the confounder; (2) since there are millions of papers, we enable efficient matching via a matching-and-reranking approach, first using information retrieval (IR) (Manning et al., 2008) to extract a small set of candidates, and then applying semantic textual similarity (STS) (Majumder et al., 2016; Chandrasekaran and Mago, 2022) for fine-grained reranking; and (3) we enable a more stable causal effect estimation by leveraging all the close matches to synthesize the *counterfactual citation score* by a weighted average according to the similarity scores of the matched papers.

CAUSALCITE quantifies scientific impact via a causal lens, offering an alternative understanding of a paper’s impact within the academic community. To test its effectiveness, we conduct extensive experiments using the Semantic Scholar corpus (Lo et al., 2020; Kinney et al., 2023), comprising of 206M papers and 2.4B citation links. We empirically validate CAUSALCITE by showing higher predictive accuracy of paper impact (as judged by scientific experts on a past dataset of 1K papers (Zhu et al., 2015)) compared to citations and other previous impact assessment metrics. We further show a stronger correlation of the metric with the test-of-time (ToT) paper awards. We find that, unlike citation counts, our metric exhibits a greater balance across various research domains in AI, e.g., general AI, NLP, and computer vision (CV). While citation numbers for papers in these domains vary significantly – for example, while an average CV paper has many more citations than an average NLP paper, CAUSALCITE scores papers across AI sub-fields more similarly.

After demonstrating the desirable properties of our metric, we also present several case studies of its applications. Our findings reveal that the quality of conference best papers is noisier on average than that of ToT papers (Section 9.5.1). We then showcase and present CAUSALCITE for several well-known papers (Section 9.5.3) and utilize CAUSALCITE to identify high-quality papers that are less recognized by citation counts (Section 9.5.4).

In conclusion, our contributions are as follows:

1. We introduce CAUSALCITE, a counterfactual causal effect-based formulation for paper citations.
2. We develop TEXTMATCH, a new method that leverages LLMs and causal inference to estimate the counterfactual causal effect of a paper.
3. We conduct comprehensive analyses, including various performance evaluations and present new findings using our metric.

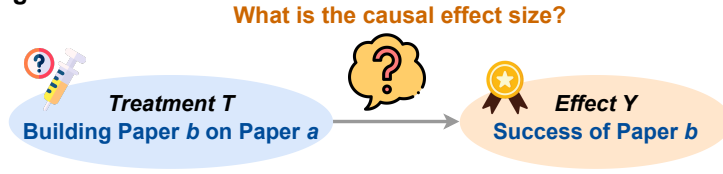
9.2 Problem Formulation

Our problem formulation involves a citation graph and a causal graph. We use lower-case letters for specific papers and uppercase for an arbitrary paper treated as a random variable.

Citation Graph In the citation graph $\mathbb{G} := (\mathbb{P}, \mathbb{L})$, \mathbb{P} is a set of papers, and each edge $\ell_{i,j} \in \mathbb{L}$ indicates that an earlier paper p_i influences (i.e., is cited by) a follow-up paper p_j . To obtain the citation graph, we use the Semantic Scholar Academic Graph dataset (Kinney et al., 2023) with 206M papers and 2.4B citation edges.

Causal Graph. The causal graph, shown in Figure 9.2, highlights the contribution of a paper a to a follow-up paper b . We use a binary variable T to indicate if a influences b and an effect variable Y to represent the success of b . We use \log_{10} of citation counts to quantify Y , although other transformations can also be used. We introduce two sets of variables in this causal graph: (i) The set of confounders, which are the common causes of T and Y . For instance, the research area of b impacts both the likelihood of a paper citing a and its own citation count. (ii) Descendants of the treatment, comprising mediators (e.g., paper a influencing the quality of paper b and subsequently influencing its citations) and

Target:



We use the causal graph to identify the correct variables to control for:

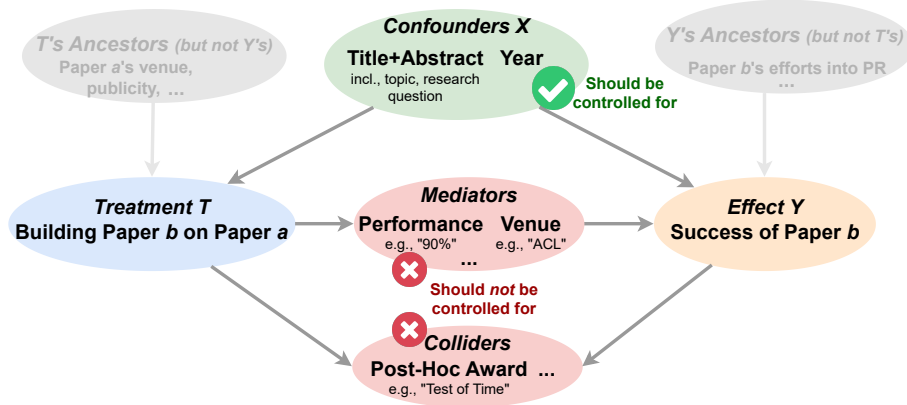


Figure 9.2: The causal graph of our study.

colliders (e.g., both the influence from *a* and the citations of *b* influencing later awards received by *b*).

9.2.1 CausalCite Indices

In this section, we introduce various indices that measure the causal impact of a paper.

Two-Paper Interaction: Pairwise Causal Impact (PCI). To examine the causal impact of a paper *a* on a follow-up paper *b*, we define the pairwise causal impact $\text{PCI}(a, b)$ by unit-level causal effect:

$$\text{PCI}(a, b) := y^{t=1} - y^{t=0}, \quad (9.1)$$

where we compare the outcomes *Y* of the paper *b* had it been influenced by paper *a* or not, denoted as the actual $y^{t=1}$ and the counterfactual $y^{t=0}$, respectively. Note that the counterfactual $y^{t=0}$ can never be observed, but only estimated by statistical methods, as we will discuss in Section 9.3.2.

Single-Paper Quality Metrics: Total Causal Impact (TCI) and Average Causal Impact (ACI). Let S denote the set of all follow-up studies of paper *a*. We define total causal impact $\text{TCI}(a)$ as the sum of the pairwise causal impact index $\text{PCI}(a, b)$ across all $b \in S$. That is,

$$\text{TCI}(a) := \sum_{b \in S} \text{PCI}(a, b). \quad (9.2)$$

This definition provides an aggregated measure of a paper's influence across all its follow-up papers.

As the causal inference literature is usually interested in the average treatment effect, we further define the average causal impact (ACI) index as the average per paper PCI:

$$\text{ACI}(a) := \frac{\text{TCI}(a)}{|S|} = \frac{1}{|S|} \sum_{b \in S} (y^{t=1} - y^{t=0}) . \quad (9.3)$$

We note that $\text{ACI}(a)$ is equal to the average treatment effect on the treated (ATT) of paper a (Pearl, 2009b).

9.3 The TextMatch Method

As illustrated in Figure 9.1, the objective of our study is to quantify the causal effect of the treatment T (i.e., whether paper b is built on paper a) on the effect Y (i.e., the outcome of paper b). To approach this, we envision a counterfactual scenario: what if paper a had never been published, yet certain key characteristics of paper b remain unchanged? The critical question then becomes: which key characteristics of paper b should be *controlled* for in this hypothetical situation?

9.3.1 What Does Causal Inference Tell Us about What Variables to Control for, and What Not?

In causal inference, selecting the appropriate variables for control is a delicate and crucial process that affects the accuracy of the analysis. Pearl’s seminal work on causality guides us in differentiating between various types of variables (Pearl, 2009b).

Firstly, we must control for *confounders* – variables that influence both the treatment and the outcome. Confounders can create spurious correlations; if not controlled, they can lead us to mistakenly attribute the effect of these external factors to the treatment itself. For example, in assessing the impact of one paper on another, if both papers are in a trending research area, the apparent influence might be due to the popularity of the topic rather than the papers’ content.

However, not all variables warrant control. Mediators and colliders should be explicitly avoided in control. Mediators are part of the causal pathway between the treatment and outcome. By controlling them, we would block the very effect we are trying to measure. Colliders, affected by both the treatment and the outcome, can introduce bias when controlled. Controlling a collider can inadvertently create associations that do not naturally exist. In general, this also includes not controlling for the descendants of the treatment, as it could obscure the direct impact we intend to study.

Lastly, variables that do not share a causal path with both the treatment and outcome, known as *unshared ancestors*, are less critical in our analysis. They do not contribute to or confound the causal relationship we are exploring, and thus, controlling for them does not add value to our causal understanding.

9.3.2 Can Existing Causal Inference Methods Handle This Control?

Several causal inference methods have been proposed to address the problem of estimating treatment effects while controlling for confounders. Next, we will discuss the workings and limitations of three classical methods.

Randomized Control Trials (RCTs) Assumes Intervenableity. The ideal way to obtain causal effects is through randomized control trials (RCTs). For example, when testing a drug, we randomly split all patients into two groups, the control group and the treatment group, where the random splitting ensures the same distribution of the confounders across the two groups such as gender and age. However, RCTs are usually not easily achievable, in some cases too expensive (e.g., tracking hundreds of people’s daily lives for 50 years), and in other cases unethical (e.g., forcing a random person to smoke), or infeasible (e.g., getting a time machine to change a past event in history).

For our research question on a paper’s impact, utilizing RCTs is impractical as it is infeasible to randomly divide researchers into two groups, instructing one group to base their research on a specific paper a while the other group does not, and then observe the citation count of their papers years later.

Ratio Matching Iterates over Discretized Confounder Values. In the absence of RCTs, matching is as an alternate method for determining causal effects from observational data. In this case, we can let the treatment assignment happen naturally, such as taking the naturally existing set of papers and running causal inference by adjusting for the variables that block all paths. Given a set of naturally observed papers, one of the most commonly used causal inference methods is ratio matching (Rosenbaum and Rubin, 1983), whose basic idea is to iterate over all possible values x of the adjustment variables X and obtain the difference between the treatment group \mathcal{T} and control group \mathcal{C} :

$$\widehat{ACI}(a) = \sum_x P(x) \left(\frac{1}{|\mathcal{T}_x|} \sum_{i \in \mathcal{T}_x} y_i - \frac{1}{|\mathcal{C}_x|} \sum_{j \in \mathcal{C}_x} y_j \right), \quad (9.4)$$

where for each value x , we extract all the units corresponding to this value in the treatment and control sets, compute the average of the effect variable Y for each set, and obtain the difference.

While ratio matching is practical when there is a small set of values for the adjustment variables to sum over, its applicability dwindles with high-dimensional variables like text embeddings in our context. This scenario may generate numerous intervals to sum over, presenting numerical challenges and potential breaches of the positivity assumption.

One-to-One Matching Is Susceptible to Variance. To handle high-dimensional adjustment variables, one possible way is to avoid pre-defining all their possible intervals, but, instead, iterating over each unit in the treatment group to match for its closest control unit (e.g., McGue et al., 2010; Sato et al., 2022). Consider a given follow-up paper b , and a set of candidate control papers C , where each paper c_i has a citation count y_i , and vector representation t_i of the confounders (e.g., research topic). One-to-one matching estimates PCI as

$$\begin{aligned} \widehat{PCI}(a, b) &= y_b - y_{\arg\max_{c_i \in C} m_i} \\ &= y_b - y_{\arg\max_{c_i \in C} \text{sim}(t_b, t_i)}, \end{aligned} \quad (9.5)$$

where we approximate the counterfactual sample by the paper $c_i \in C$ which is the most similar to paper b by the matching score m_i , which is obtained by the cosine similarity

sim of the confounder vectors. A limitation of the one-to-one matching method is that it might induce large instability in the result, as only taking one paper with similar contents may have a large variance in citations when the matched paper slightly differs.

9.3.3 How Do We Extending Causal Inference to Text Variables?

9.3.3.1 Theoretical Formulation of TextMatch: Stabilizing Text Matching by Synthesis

To fill in the aforementioned gap in the existing matching methods, we propose **TEXTMATCH**, which mitigates the instability issue of one-to-one matching by replacing it with a convex combination of a set of matched samples to form a synthetic counterfactual sample. Specifically, we identify a set of papers $c_i \in C$ with high matching scores m_i to the paper b , and synthesize the counterfactual sample by an interpolation of them:

$$\widehat{PCI}(a, b) = y_b - \sum_{c_i \in C} w_i y_i = y_b - \sum_{c_i \in C} \frac{m_i}{\sum_{c_i \in C} m_i} y_i, \quad (9.6)$$

where the weight w_i of each paper c_i is proportional to the matching score m_i and normalized.

The contributions of our method are as follows: (1) we adapt the traditional matching methods from low-dimensional covariates to any high-dimensional variables such as text embeddings; (2) different from the ratio matching, we do not stratify the covariates, but synthesize a counterfactual sample for each observed treated units; (3) due to this iteration over each treated unit instead of taking the population-level statistics, we closely control for exogenous variables for the ATT estimation, which circumvents that need for the structural causal models; (4) we further stabilize the estimand by a convex combination of a set of similar papers. Note that the contribution of Eq. (9.6) might seem to bear similarity with synthetic control (Abadie and Gardeazabal, 2003; Abadie et al., 2010), but they are fundamentally different, in that synthetic control runs on time series, and fit for the weights w_i by linear regression between the time series of the treated unit and a set of time series from the control units, using each time step's values in the regression loss function.

9.3.3.2 Overall Algorithm

To operationalize our theoretical formulation above, we introduce our overall algorithm in Algorithm 1. We briefly give an overview of the the algorithm with more details to be elaborated in later sections. We use the weighted average of the matched samples following our **TEXTMATCH** method in Eq. (9.6) through lines 25 to 34. In our experiments, we use the interpolation of up to top 10 matched papers. We encourage future work to explore other hyperparameter settings too. Given the PCI estimation, the main spirit of the **GETACIANDTCI(a)** function is to average or sum over all the follow-up studies of paper a , following the theoretical formulation in Eqs. (9.2) and (9.3) and implemented in our algorithm through lines 7 to 12.

9.3.3.3 Key Challenges and Mitigation Methods

We address several technical challenges below.

Algorithm 1 Get causal impact indices ACI and TCI

```

1: Input: Paper  $a$ .
2: procedure GETACIANDTCI( $a$ )
3:    $D \leftarrow \text{GetDesc}(a)$  ▷ Get descendants by DFS
4:    $B \leftarrow \text{GetChildren}(a)$ 
5:    $B' \leftarrow \text{SampleSubset}(B)$  ▷ See Section 9.3.3.3.4
6:    $C \leftarrow \text{EntireSet} \setminus \{D \cup \{a\}\}$  ▷ Get non-descendants
7:    $ACI \leftarrow 0$ 
8:   for each  $b_i$  in  $B'$  do
9:      $I_i \leftarrow \text{GETPCI}(a, b_i, C)$ 
10:     $ACI \leftarrow ACI + \frac{1}{|B'|} \cdot I_i$ 
11:   end for
12:    $TCI \leftarrow ACI \cdot |B|$ 
13:   return  $ACI$  and  $TCI$ 
14: end procedure

15: procedure GETPCI( $a, b, C$ )
16:    $C_{\text{sameYear}} \leftarrow \text{FilterByYear}(C, b_{\text{year}})$ 
17:   for each  $p_i$  in  $C_{\text{sameYear}} \cup \{b\}$  do
18:      $t_i \leftarrow \text{RemoveMediator}(\text{TitleAbstract}_i)$ 
19:   end for
20:    $C_{\text{coarse}} \leftarrow \text{BM25}(b, C_{\text{sameYear}}, \text{topk} = 100)$ 
21:   for each  $c_i$  in  $C_{\text{coarse}}$  do
22:      $m_i \leftarrow \text{Sim}(t_b, t_i)$ 
23:   end for
24:    $C_{\text{top10}} \leftarrow \text{argmax}_{10\_m}(C_{\text{coarse}})$ 

25:    $M \leftarrow 0$ 
26:   for each  $c_i$  in  $C_{\text{top10}}$  do ▷ For the normalization later
27:      $M \leftarrow M + m_i$ 
28:   end for
29:    $\hat{y}^{t=0} \leftarrow 0$ 
30:   for each  $c_i$  in  $C_{\text{top10}}$  do
31:      $w_i \leftarrow \frac{m_i}{M}$ 
32:      $\hat{y}^{t=0} \leftarrow \hat{y}^{t=0} + w_i \cdot y_i$  ▷ Apply Eq. (9.6)
33:   end for
34:   return  $y_b - \hat{y}^{t=0}$ 
35: end procedure

```

9.3.3.3.1 Confounders of Various Types

First, as we mentioned in the causal graph in Figure 9.2, the confounder set consists of a text variable (title and abstract concatenated together) and an ordinal variable (publication year). Therefore, the similarity operation Sim between two papers should be customized. For our specific use case, we first filter by the publication year in line 16, as it is not fair to compare the citations of papers published in different years. Then, we apply the cosine similarity method paper embeddings as in line 22. As a general solution, we recommend to separate hard logical constraints, and soft matching preferences, where the hard constraints should be imposed to filter the data first, and then all the rest of the variables can be concatenated to apply the similarity metric on.

9.3.3.3.2 Excluding the Mediators from Confounders

Another key challenge to highlight is that the text variable we use for the confounder might accidentally include some mediator information. For example, the quality or per-

formance of a paper could be expressed in the abstract, such as “we achieved 90% accuracy.” Therefore, we conduct a specific preprocessing procedure before feeding the text variable to the similarity function. For the RemoveMediator function in line 18, we exclude all numerical expressions such as percentage numbers, as well as descriptions such as “state-of-the-art.” For generalizability, the essence of this step is an entanglement action to separate the confounder variable (in this case, the research content) and all the descendants of the treatment variable (in this case, mentions of the performance). For more complicated cases in future work, we recommend a separate disentanglement model to be applied here.

9.3.3.3 Efficient Matching-and-Reranking Method

Since we use one of the largest available paper databases, the Semantic Scholar dataset (Kinney et al., 2023) containing 206M papers, we need to optimize our algorithm for large-scale paper matching. For example, after we filter by the publication year, the number of candidate papers C_{sameYear} could be up to 8.8M. In order to conduct text matching across millions of papers, we use a *matching-and-reranking* approach, by combining two NLP tasks, information retrieval (IR) (Manning et al., 2008) and semantic textual similarity (STS) (Majumder et al., 2016; Chandrasekaran and Mago, 2022).

Specifically, we first run large-scale matching to obtain 100 candidate papers (line 20) using the common IR method, BM25 (Robertson and Zaragoza, 2009). Briefly, BM25 is a bag-of-words retrieval function that uses term frequencies and document lengths to estimate relevancy between two text documents. Deploying this method, we can find a set of candidate papers for, for example, two million papers, at a speed 250x faster than the text embedding cosine similarity matching. Then, we conduct a fine-grained reranking using cosine similarity (lines 21 to 23). In the cosine similarity matching process, we use the MPNet model (Song et al., 2020) to encode the text of each paper c_i into an embedding t_i , with which we get the matching score m_i according to Eq. (9.5) in line 22, and the normalized weight w_i by Eq. (9.6) in line 31.

9.3.3.4 Numerical Estimation

Given the large number of papers, it is numerically challenging to aggregate the TCI from individual PCIs, because the number of follow-up papers for a study can be up to tens of thousands, such as the 57,200 citations by 2023 for the ImageNet paper (Deng et al., 2009). To avoid extensively running PCI for all follow-up papers, we propose a new numerical estimation method using a carefully designed random paper subset.

A naive way to achieve this aggregation is Monte Carlo (MC) sampling. However, unfortunately, MC sampling requires very large sample sizes when it comes to estimating long-tailed distributions, which is the usual case of citations. Since citations are more likely to be concentrated in the head part of the distribution, we cannot afford the computational budget for huge sample sizes that cover the tails of the distribution. Instead, we propose a novel numerical estimation method for sampling the follow-up papers, inspired by importance sampling (Singh, 2014; Kloek and van Dijk, 1976).

Our numerical estimation method works as follows: First, we propose the formulation that the relation between ACI and TCI is an integral over all possible paper b 's. Then, we formulated the above sampling problem as integral estimation or area-under-the-curve

estimation. We draw inspiration from Simpson’s method, which estimates integrals by binning the input variable into small intervals. Analogously, although we cannot run through all PCIs, we use citations as a proxy, bin the large set of follow-up papers according to their citations into n equally-sized intervals, and perform random sampling over each bin, which we then sum over. In this way, we make sure that our samples come from all parts of the long-tailed distribution and are a more accurate numerical estimate for the actual TCI.

9.4 Performance Evaluation

The contribution of a paper is inherently multi-dimensional, making it infeasible to encapsulate its richness fully through a scalar. Yet the demand for a single, comprehensible metric for research impact persists, fueling the continued use of traditional citations despite their known limitations. In this section, we show how our new metrics significantly improve upon traditional citations by providing quantitative evaluations comparing the effectiveness of citations, Semantic Scholar’s highly influential (SSHI) citations (Valenzuela et al., 2015), and our CAUSALCITE metric.

9.4.1 Experimental Setup

Dataset We use the Semantic Scholar dataset (Lo et al., 2020; Kinney et al., 2023)¹ which includes a corpus of 206M scientific papers, and a citation graph of 2.4B+ citation edges. For each paper, we obtain the title and abstract for the matching process. We list some more details of the dataset in Appendix A.8.2, such as the number of papers reaching 8M per year after 2012.

Selecting the Text Encoder When projecting the text into the vector space, we need a text encoder with a strong representation power for scientific publications, and is sensitive towards two-paper similarity comparisons regarding their abstracts containing key information such as the research topics. For the representation power for scientific publications, instead of general-domain models such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019), we consider LLM variants² pretrained on large-scale scientific text, such as SciBERT (Beltagy et al., 2019), SPECTER (Cohan et al., 2020), and MPNet (Song et al., 2020).

To check the quality of two-paper similarity measures, we conduct a small-scale empirical study comparing human-ranked paper similarity and model-identified semantic similarity in Appendix A.8.1.3, according to which MPNet outperforms the other two models.

Implementation Details We deploy the *all-mpnet-base-v2* checkpoint of the MPNet using the *transformers* Python package (Wolf et al., 2020), and set the batch size to be 32. For the set of matched papers, we consider papers with cosine similarity scores higher than 0.81, which we optimize empirically on 100 random paper pairs. We take the top ten most similar papers above the threshold. In special cases where there is no matched paper above the threshold, it means that no other paper works on the same idea as Paper b, and

¹<https://api.semanticscholar.org/api-docs/datasets>

²Note that we follow the standard notion by Yang et al. (2023) to refer to BERT and its variants as LLMs.

we make the counterfactual citation number to be zero, which also reflects the quality of Paper b as its novelty is high.

To enable efficient operations on the large-scale citation graph, we use the Dask framework,³ which optimizes for data processing and distributed computing. We optimize our program to take around 100GB RAM, and on average 25 minutes for each $\text{PCI}(a, b)$ after matching against up to millions of candidates. More implementation details are in Appendix A.8.1.1. For the estimation of TCI, we empirically select the sample size to be 40, which is a balance between the computational time and performance, as found in Appendix A.8.1.2.

9.4.2 Author-Identified Paper Impact

In this experiment, we follow the evaluation setup in Valenzuela et al. (2015) to use an annotated dataset (Zhu et al., 2015) comprised of 1,037 papers, annotated according to whether they serve as significant prior work for a given follow-up study. Although paper quality evaluation can be tricky, this dataset was cleverly annotated by first collecting a set of follow-up studies and letting one of the authors of each paper go through the references they cite and select the ones that significantly impact their work. In other words, for a given paper b, each reference a is annotated as whether a has significantly impacted b or not.

Table 9.1 reports the accuracy of our CAUSALCITE metric, together with two existing citation metrics: citations, and SSHI citations (Valenzuela et al., 2015). See the detailed derivation of the accuracy scores in Appendix A.8.3.2. From this table, we can see that our CAUSALCITE metric achieves the highest accuracy, 80.29%, which is 5 points higher than SSHI, and 9 points higher than the traditional citations.

9.4.3 Test-of-Time Paper Analysis

The test-of-time (ToT) paper award is a prestigious honor bestowed upon papers that have made substantial and enduring impacts in their field. In this section, we collect a dataset of 792 papers, including 72 ToT papers, and a control group of 10 randomly selected non-ToT papers from the same conference and year as each ToT paper. To collect this ToT paper dataset, we look into ten leading AI conferences spanning general AI (NeurIPS, ICLR, ICML, and AAAI), NLP (ACL, EMNLP, and NAACL), and CV (CVPR, ECCV, and ICCV), for which we go through each of their websites to identify all available ToT papers.⁴

In Table 9.2, we show the correlations of various metrics with the ToT awards. In this table, CAUSALCITE achieves the highest correlation of 0.639, which is +30.14% better than that of citations. Furthermore, we visualize the correspondence of our metric and ToT, and observe a substantial difference between the CAUSALCITE distributions of ToT vs. non-ToT papers in Figure 9.3. We also show three examples of ToT papers in Figure 9.4, where the ToT papers differ from the non-ToT papers by one or two orders of magnitude.

Metric	Accuracy
Citations	71.33
SSHI Citations	75.25
CAUSALCITE	80.29

Table 9.1: Accuracy of all three citation metrics.

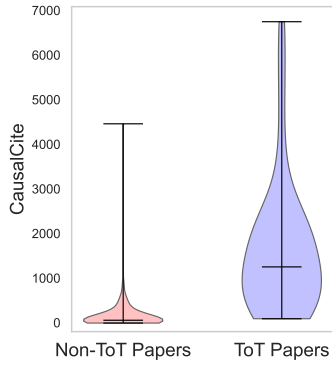


Figure 9.3: Distributions of ToT (mean: 142) and non-ToT papers (mean: 1,623).

Metric	Corr. Coef.
Citations	0.491
SSHI Citations	0.317
TCI	0.640

Table 9.2: Correlation coefficients of each metric and ToT paper award by Point Biserial Correlation (Tate, 1954).

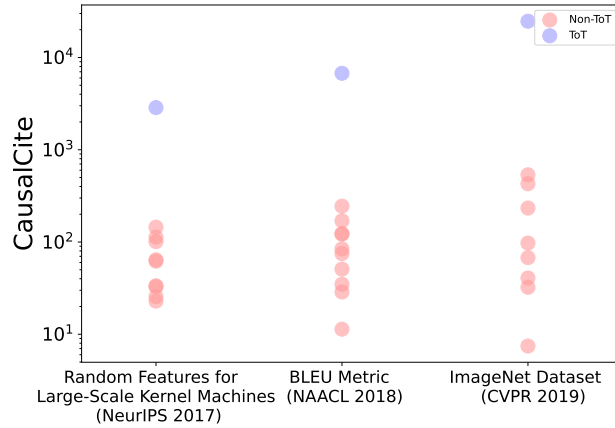


Figure 9.4: The CAUSALCITE values of three example ToT papers from general AI, NLP, and CV.

Research Area	ACI	Citations	SSHI
General AI (n=16)	0.748	2,024	267
CV (n=36)	0.734	7,238	1,088
NLP (n=20)	0.763	1,785	461

Table 9.3: The average of each metric by research area on our collected set of 72 ToT papers.

9.4.4 Topic Invariance of CausalCite

A well-known issue with citations is their inconsistency across different fields. What might be considered a large number of citations in one field might be seen as average in another. In contrast, we show that our ACI index does not suffer from this issue. We show this using our ToT dataset, where we control for the quality of the papers to be ToT but vary the domain by the three fields: general AI, CV, and NLP. We observe in Table 9.3 that even though some domains have significantly more citations (for instance, CV ToT papers have, on average, 4.05 times more citations than NLP), the ACI remains consistent across various fields.

9.5 Findings

Having demonstrated the effectiveness of our metrics, we now explore some open-ended questions: (1) Do best papers have high causal impact? (Section 9.5.1) (2) How does the

³<https://dask.org/>

⁴We get this list by selecting the top conferences on Google Scholar using the h5-Index ranking in each of the above domains: general AI (link), CV (link), and NLP (link).

CAUSALCITE value distribute across papers? (Section 9.5.2) (3) What is the impact of some famous papers evaluated by CAUSALCITE? (Section 9.5.3) (4) Can we use this metric to correct for citations? (Section 9.5.4).

9.5.1 Do Best Papers Have High Causal Impact?

Selecting best paper awards is an arguably much harder task than ToT papers, as it is difficult to predict of the impact of a paper when it is just newly published. Therefore, we are interested in the actual causal impact of best papers. Similar to our study on ToT papers, we collect a dataset of 444 papers including 74 best papers and a control set of random 5 non-best papers from the same conference in the same year, using the same set of the top ten leading AI conferences. We find that the correlation of the CAUSALCITE metric with best papers is 0.348, which is very low compared to the 0.639 correlation with the ToT papers. This shows that the best papers do not necessarily have a high causal impact. One interpretation can be that the best paper evaluation is a forecasting task, which is much more challenging than the retrospective task of ToT paper selection.

9.5.2 What Is the Nature of the CausalCite Distribution?

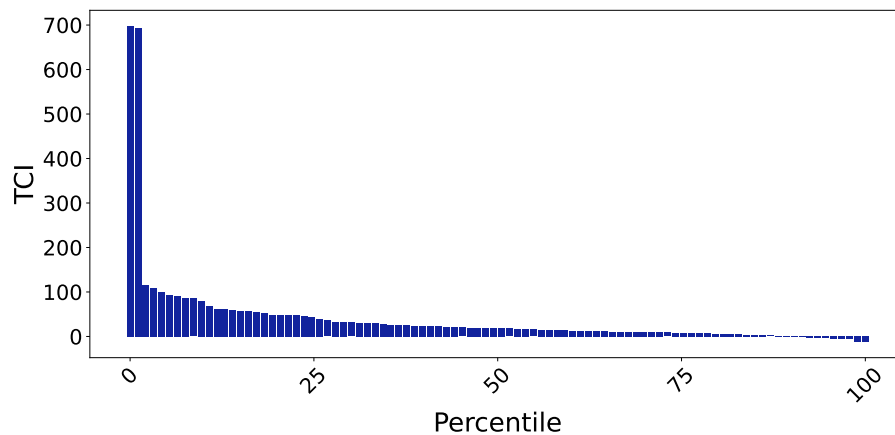


Figure 9.5: The distribution of TCI values by percentile of 100 random papers, which shows a long tail indicating that high impact is concentrated in a relatively small portion of papers.

We explore how the CAUSALCITE scores are distributed across papers in general. We plot Figure 9.5 using a random set of 100 papers from the Semantic Scholar dataset, which is a reasonably large size given the computation budget mentioned in Section 9.4.1. From this plot, we can see a power law distribution with a long tail, echoing with the common belief that the paper impact follows the power law, with high impact concentrated in a relatively small portion of papers.

9.5.3 Selected Paper Case Study

In addition to the shape of the overall distribution, we also look at our metric's correspondence to some selected papers shown in Table 9.4. For example, we know that the Transformer paper (Vaswani et al., 2017) is a more foundational work than its follow-up work BERT (Devlin et al., 2019), and BERT is more foundational than its later variant,

Paper Name	TCI	Citations	ACI
Transformers	52,507	68,064	0.771
BERT	40,675	59,486	0.683
RoBERTa	6,932	14,434	0.480

Table 9.4: Case study of some selected NLP papers.

RoBERTa (Liu et al., 2019). This monotonic trend is confirmed in their TCI and ACI values too. Again, this is a preliminary case study, and we welcome future work to cover more papers.

9.5.4 Discovering Quality Papers beyond Citations

Another important contribution of our metric is that it can help discover papers that are traditionally overlooked by citations. To achieve the discovery, we formulate the problem as outlier detection, where we first use a linear projection to handle the trivial alignment of citations and CAUSALCITE, and then analyze the outliers using the interquartile range (IQR) method (Smiti, 2020). See the exact calculation in Appendix A.8.3.1. We show the three subsets of papers in Table 9.5, where the two outlier categories, the overcited and undercited papers, correspond to the false positive and false negative oversight by citations, respectively. An additional note is that, when we look into some characteristics of the three categories, we find that the citation frequency in result section, i.e., the percentage of times they are cited in results section compared to all the citations, correlates with these categories. Specifically, we find that the undercited papers tend to have more of their citations concentrated in the results section, which usually indicates that this paper constitutes an important baseline for a follow-up study, while the overcited papers tend to be cited out of the results section, which tends to imply a less significant citation.

Paper Category	Result Citations	Residual
Overcited Papers (7.04%)	1.26	-1.792
Aligned Papers (91.20%)	1.51	0.118
Undercited Papers (1.76%)	1.90	1.047

Table 9.5: We use our CAUSALCITE metric to discover outlier papers that are overlooked by citations. For each paper category, we include their portion relative to the entire population, the percentage of citations occurred in the result section (Result Citations), and average residual value by linear regression.

9.6 Related Work

The quantification of scientific impact has a rich history and continuously evolves with technology. Bibliometric analysis has been largely influenced by early methods that relied on citation counts (Garfield et al., 1964; Garfield, 1972, 1964). Hou (2017) investigate the evolution of citation analysis, employing reference publication year spectroscopy (RPYS) to trace its historical development in scientometrics. Donthu et al. (2021) provide practical guidelines for conducting bibliometric analysis, focusing on robust methodologies to analyze scientific data and identify emerging research trends.

Indices such as the h-index, introduced by Hirsch (2005), are established tools for mea-

asuring research impact. The more recent Relative Citation Ratio (RCR), developed by Hutchins et al. (2016), provides a field-normalized alternative to traditional metrics. Valenzuela et al. (2015) introduced SSHI, an approach to identify meaningful citations in scholarly literature. However, these metrics are not without limitations. As Wróblewska (2021) discussed, conventional citation-based metrics often fail to capture the multidimensional nature of research impact. In this context, Elmore (2018) discussed the Altmetric Attention Score, which evaluates the broader societal and online impact of research.

With the increasing availability of large datasets and the advent of digital technologies, new opportunities for bibliometric analysis have emerged. Iqbal et al. (2021) highlighted the role of NLP and machine learning in enhancing in-text citation analysis. Similarly, Umer et al. (2021) explored the use of textual features and SMOTE resampling techniques in scientific paper citation analysis. Jebari et al. (2021) analyzed citation context to detect research topic evolution, showcasing data analysis for scientific discourse. Chang et al. (2023) explored augmenting citations in scientific papers with historical context, offering a novel perspective on citation analysis. Manghi et al. (2021) introduced scientific knowledge graphs, an innovative method for evaluating research impact. Bittmann et al. (2021) explored statistical matching in bibliometrics, discussing its utility and challenges in post-matching analysis. The use of AI in bibliometric analysis is highlighted in research by Chubb et al. (2022) and the systematic review of AI in information systems by Collins et al. (2021). Network analysis approaches, as discussed by Chakraborty et al. (2020b) in the context of patent citations and by Dawson et al. (2014) in learning analytics, further illustrate the diverse applications of advanced methodologies in understanding citation patterns.

9.7 Conclusion

In this study, we propose CAUSALCITE, a novel causal formulation for paper citations. Our method combines traditional causal inference methods with the recent advancement of NLP in LLMs to provide a new causal outlook on paper impact by answering the causal question: “Had this paper never been published, what would be the impact on this paper’s current follow-up studies?”. With extensive experiments and analyses using expert ratings and test-of-time papers as criteria for impact, our new CAUSALCITE metric demonstrates clear improvements over the traditional citation metrics. Finally, we use this metric to investigate several open-ended questions like “Do best papers have high causal impact?”, conduct a case study of famous papers, and suggest future usage of our metric for discovering good papers less recognized by citations for the scientific community.

Limitations and Future Work

There are several limitations for our work. For example, as mentioned previously, our metric has a high computational budget. Future work can explore more efficient optimization methods. Also, we model the content of the paper by its title and abstract, it could also be possible for future work to benefit from modeling the full text, given appropriate license permissions.

As for another limitation, our study is based on data provided by the Semantic Scholar corpus. This corpora has certain properties such as being more comprehensive with com-

puter science papers, but less so in other disciplines. Its citation data also has a delay compared to Google Scholar, so for the newest papers, the citation score may not be accurate, making it more difficult to calculate our metric.

Additionally, our study provides a general framework for causal inference given a causal graph that involves text. It is totally possible that for a more fine-grained problem, the causal graph will change, in which case, we undersuggest future researchers to derive the new backdoor adjustment set, and then adjust the algorithm accordingly. An example of such a variable could be the author information, which might also be a confounder.

Finally, since quality evaluation of a paper is a multi-faceted task, theoretically, a single number can never give more than a rough approximation, because it collapses multiple dimensions into one and loses information. Our argument in this paper is just to show that our formulation is theoretically more accurate than the citation formulation. We take one step further, instead of solving the quality evaluation problem which is much more nuanced. Some intrinsic problems in citations that we can also not solve (because our metrics still rely on using citations, just contrasting them in the right away) include (1) if a paper is newly published, with zero citations, there is no way to obtain a positive causal index, and (2) we do not solve the fair attribution problem when multiple authors share credit of a paper, as our metric is not sensitive towards authors.

Ethical Considerations

Data Collection and Privacy The data used in this work are all from Open Source Semantic Scholar data, with no user privacy concerns. The potential use of this work is for finding papers that are unique and innovative but do not get enough citations due to lack of popularity or awareness of the field. This metric can act as an aid when deciding impact of papers, but we do not suggest its usage without expert involvement. Through this work, we are not trying to demean or criticize anyone's work we only intend to find more papers that have made a valuable contribution to the field.

CS-Centric Perspective The authors of this paper work in Computer Science (mostly Machine Learning) hence a lot of analysis done on the quality of papers that required sanity checks are done on ML papers. The conferences selected for doing the ToT evaluation were also CS Top conferences, hence they might have induced some biases. The metric in general has been created generically and should be applicable to other domains as well, the Author Identified Most Influential Papers study is also done on a generalized dataset, but we encourage readers in other disciplines to try out the metric on papers from their field.

Conclusion & Future Work

This dissertation explored causal methods for NLP. We first started with the question of whether LLMs can do causal reasoning. In Part I, we build formal causal reasoning benchmarks covering two key skills, causal discovery (Chapter 2) and causal effect reasoning (Chapter 3), which existing models struggle to address. We develop CAUSALCoT to enhance model performance, inspired by a combination of symbolic and LLM-driven steps. A natural question next is to interpret how models make their decisions, which leads to our exploration in Part II to use causal methods for intrinsic (Chapter 4) and behavioral interpretability analysis (Chapter 5).

Apart from causal inference to improve performance and to interpret the models, we also look into the causality among the variables in the dataset in Part III. We find that whether the variables hold a causal or an anticausal relation has large impact on their performance in different settings (Chapter 6). We further extend this exploration to cases where the variable causal relations are not evident *a priori*, but discovered using interdisciplinary insights (Chapter 7). Based on the variable relations we also suggest causal prompts to make use of those relations to improve LLMs' decisions.

Lastly, we also demonstrate social applications using causal methods and NLP together in Part IV. Our first study utilizes NLP for sentiment classification on social media text, and perform causal effect estimation between sentiment and social policies (Chapter 8). Our second study looks into causal analysis for paper citations, which is a more technically challenging case to adjust for textual confounders, where we use LLMs to encode the text to high-dimensional vector space, and perform a text-based matching algorithm (Chapter 9).

Building on the insights gained from all the above research, we outline a few open challenges and future directions:

The first direction to highlight is to build a more comprehensive framework of causal reasoning. My previous work distinguishes two types of causal understanding in LLMs: knowledge-based causal understanding (Cui et al., 2024), and knowledge-independent formal reasoning (Jin et al., 2023a, 2024). Looking ahead, there is a large need to systematize reasoning as an interplay of both. To achieve it, one way is to develop more sophisticated models and training strategies that enhance the causal inference abilities of LLMs. This includes exploring novel architectures, fine-tuning methods, and datasets that better capture causal relationships. Another way is to advance tool-augmented LLMs connecting different subskills into an overall pipeline.

Since reasoning does not necessarily need to be constrained to passively processing the text, it is also a promising direction to explore cause of reasoning in interactive systems, such as conversational agents and decision support tools. Research in this area can focus on integrating causal inference with real-time interactions and feedback mechanisms. In addition to the interaction setup, some other modalities of data can also be included as

inputs for causal inference, which can generalize text-only reasoning to image, video, audio, structured information, and many others.

In parallel to improving model performance, we can also advance the use of causal methods for interpretable, robust, and fair NLP models. Interpretability, at its essence, is a causal discovery problem, where we aim to understand what elements or mechanisms in the model contribute to the final prediction. The key technical challenges are to scale up causal inference to a large number (e.g., trillions) of neurons in the model, and to obtain human-understandable high-level interpretation. For robustness and fairness, my existing work uses SCMs to formulate them as causal alignment, between the desired decision-making mechanism, and the model-learned mechanisms. Further work can formalize more evaluation pipelines in causal terms, and apply our framework for robustness and fairness assessments.

Finally, with all the technological advancements in LLMs, we are in an era to see the rise of cross-disciplinary applications to connect Causal NLP with fields such as health-care, economics, and education. We look forward to collaborative research that leverages domain-specific knowledge with Causal NLP methods, leading to impactful advancements in these areas.

If these efforts above are successful, they could lead to the development of more intelligent and robust NLP systems capable of sophisticated causal reasoning. Such advancements have the potential to transform numerous fields. I believe that the continuous improvement and application of causal methods for LLMs will be a critical driver of progress in AI, inspiring future research and opening new avenues for interdisciplinary collaboration.

Appendix

A.1 Additional Materials for Chapter 2

A.1.1 Implementation Details

When finetuning on our data, for GPT-based models, we use the default settings of the OpenAI finetuning API; and for BERT-based models, we use the `transformers` library (Wolf et al., 2020) and train the models on a server with an NVIDIA Tesla A100 GPU with 40G of memory. To fit for the GPU memory, we set the batch size to be 8. We use the validation set to tune the learning rate, which takes value in $\{2e-6, 5e-6, 1e-5, 2e-5, 5e-5\}$; dropout rate, which takes value in $\{0, 0.1, 0.2, 0.3\}$; and weight decay, which takes value in $\{1e-4, 1e-5\}$. We train the models until convergence, which is usually around ten epochs.

Prompts When querying the autoregressive LLMs, we formulate the prompt as follows:

Question: [premise]

Can we deduct the following: [hypothesis]? Just answer "Yes" or "No."

Answer:

A.1.2 Generating Natural Stories

To generate the natural stories based on our symbolic expressions, we utilize the state-of-the-art LLM, GPT-4, which is very good at story generation. We design detailed instructions in the prompt, and generate around 200 stories in our case study. We show two examples stories in Table A.1, and the report the overall statistics in Table A.2.

For more information, the exact prompt we use is "*Here is a causal inference rule: [symbolic form] Please provide a real-world example instantiating this phenomenon. Format it also as "Premise:", "Hypothesis:", and "Relation between the promise and hypothesis:."*"

A.1.3 Templates and Paraphrases

We use the verbalization templates in Table A.3 to compose the hypotheses for all six causal relations.

A.1.4 Change Log for the Dataset Version Update

De-Duplication Strategy As mentioned in Section 2.3.7 in the main paper, our original dataset (v1.0) has duplication due to symmetric relations and verbalizations. We introduce in Table A.4 several reasons for why duplicated hypotheses exist in our original data. One typical reason is symmetric relations such as `Is-Parent(A, B)` and `Is-Child(B,`

A), and, similarly, the paraphrased version of Is-Ancestor(A, B) and Is-Descendent(B, A). Another typical reason is the semantic equivalence in the verbalization templates, which applies to the Has-Collider and Has-Confounder relations. For example, the verbalized texts of Has-Collider(A, B) and Collider(B, A) are “There exists at least one collider (i.e., common effect) of {A and B, B and A},” respectively, which are semantically-equivalent paraphrases of each other, so we randomly keep one out of the two.

Resulting Dataset Statistics after De-Duplication Since the reason for duplication in the first place is due to symmetry in the causal relation, or verbalization, the resulting new data, CLADDER v2.0, is exactly a half of the original data. As we reported previously in Table 2.3 of Section 2.3.7, the total number of samples cuts down to half, while the label distribution and all other properties are the same. To compose each split, we apply the same de-duplication method for the test, train, and development sets. We notice that some duplicates are across the splits, so we prioritize keeping the test and training sets untouched (to minimally affect the experimental results), and then reduce the development set by removing the cross-split duplicates, namely:

- test_2.0 = deduplicate(test_1.0)
- train_2.0 = deduplicate(train_1.0)
- dev_2.0 = deduplicate(dev_1.0) \ {test_2.0, train_2.0}

We expect minimal or almost no change to the experimental results. In case of the slight possibility that this change in the development set might affect the model selection in the training process, future work can feel free to re-train the models and update the exact performance number.

A.1.5 Spurious Correlation Analysis

The inspirations of our two robustness tests (paraphrasing and variable refactorization) come from our data analysis. We check for spurious correlations in the data by reporting in Table A.5 the point-wise mutual information (PMI) between the label and any n-gram with no more than four tokens. In addition, we also report the difference of the PMI with the two labels in the |Diff| column of Table A.5, and report the top 10 n-grams.

The design spirit for our robustness test is that if the models’ correct judgment relies on exploiting these spurious correlations, then such reliance will be broken in our perturbations.

We can see that some spurious correlations are rooted in the framing of the hypothesis, such as “a cause (for)”, and “a direct (one)” (which we use the paraphrasing task to break), and others are connected to the variable names, such as “for D (but)” and “for E (but)” (which we use the variable refactorization to break).

A.1.6 Fine-Grained Error Analysis

In addition to the fine-grained analysis by causal relation type in Table 2.6a for fine-tuned models, we also report such error analysis for non-finetuned models in Table A.6.

These results are particularly revealing, showing how off-the-shelf models perform in recognizing specific relations. Specifically, GPT-3.5 cannot recognize ancestor relations,

whereas GPT-4 fails at all direct causation recognition with parents and children. And RoBERTa MNLI only did collider relation relatively correctly. Note that, when the F1 score is zero, the accuracy number is a result of always predicting the negative class of that relation.

A.1.7 LLM Performance Optimization

Since our experiments in Section 2.4.2 are based on plain, zero-shot prompts, we explore whether better prompting strategies could improve the performance. We enhance the query prompt by incorporating several strategies: (1) Utilizing a system prompt that specifies the model’s expertise (“You are a highly intelligent question-answering bot with profound knowledge of causal inference.”); (2) Including a pair of few-shot examples, one positive and one negative; (3) Implementing chain-of-thought prompting with “Let’s think step by step.” to encourage the language model to generate step-by-step reasoning. In Table A.7, we present the evaluation results on the relatively affordable model, GPT-3.5, where the optimized prompt leads to a 4-point improvement in F1 over the original performance. However, we can see that despite the deployment of all three strategies, the model continues to struggle with this challenging task.

A.2 Additional Materials for Chapter 3

A.2.1 Supplementary for Dataset Generation

A.2.1.1 List of References for Causal Inference

When collecting the causal graphs, query types, and commonsensical stories for our dataset, we took our examples from the following books (sorted by year):

1. Causality (Pearl, 2009b)
2. Causal inference in statistics: A Primer (Glymour et al., 2016)
3. Elements of Causal Inference (Peters et al., 2017)
4. The Book of Why (Pearl and Mackenzie, 2018)
5. Introduction to Causal Inference (Neal, 2020)

And the following papers:

1. Causes and Explanations: A Structural-Model Approach. Part I: Causes (Halpern and Pearl, 2005a)
2. Causes and Explanations: A Structural-Model Approach. Part II: Explanations (Halpern and Pearl, 2005b)
3. Causality and Counterfactuals in the Situation Calculus (Hopkins and Pearl, 2007)
4. Causal inference in statistics: An overview (Pearl, 2009a)

A.2.1.2 Formulation of the Query Types

Here, we introduce all the query types included in our dataset.

Rung-1 Queries: Marginal and Conditional Probabilities. For marginal probabilities, we ask questions about the overall distribution of a variable. For conditional probabilities, we ask whether conditioning on one variable increases or decreases the likelihood of another variable. For the explaining away questions, we condition on a collider node and ask how that affects the correlation between the two parents.

Rung-2 Queries: ATE and Adjustment Set. For ATE questions, we ask whether the treatment ($X = 1$) increases or decreases the likelihood of the effect variable $Y = y$. For adjustment set questions, we ask whether a set of variables should be adjusted for when estimating the causal effect between treatment and effect. By adjusting, we aim to blocked the non-causal paths from the treatments to effect, and hence eliminate spurious correlation. For example, to query whether the set gender is an adjustment set for the effect of a treatment on recovery, we ask *"To estimate the effect of the treatment on recovery, should we directly look at how the treatment correlates with recovery, or should we look at gender-specific correlation?"* In the collider bias questions, similarly to the explaining away questions, we condition on a collider variable and ask about how an intervention on one of the parents (treatment X) affects the other parent (outcome Y). However since by construction X and Y do not have common causes, the answer to this question is always "no".

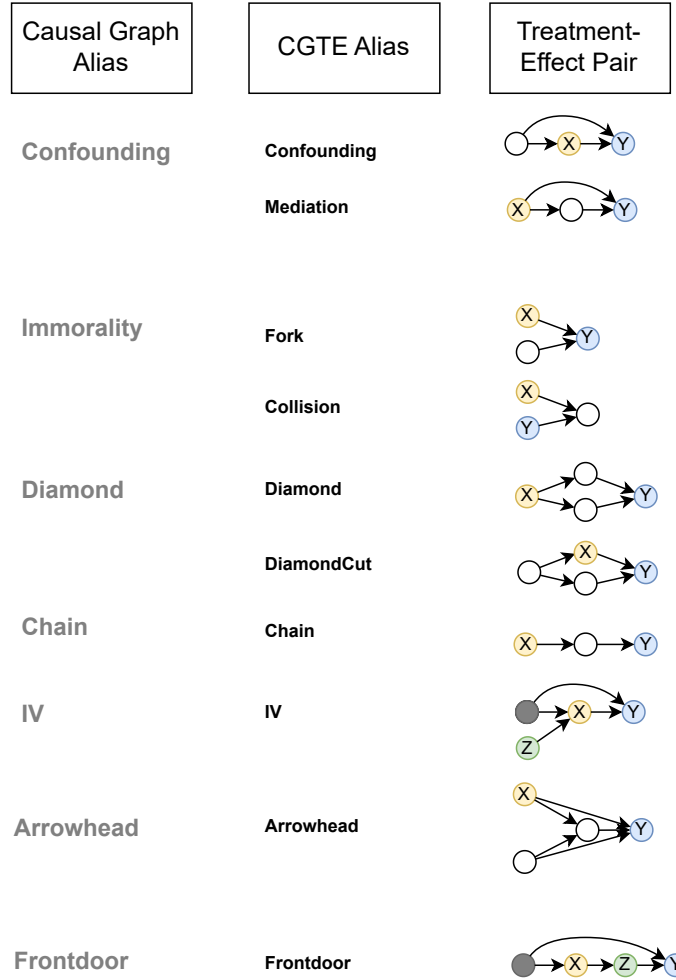


Figure A.1: List of all ten causal graphs with treatment-effect pairs (CGTEs). We omit CGTEs that trivially resemble existing ones.

Rung-3 Queries: Counterfactual Probability, ATT, NDE, and NIE. For counterfactual probability, we ask about what would have been the likelihood of $Y = y$, if the treatment variable X had been x , given sufficient evidence e such that the query is identifiable. For ATT, we ask how the likelihood of $Y = y$ would change for those who received treatment ($X = 1$) if there had been no treatment ($X = 0$). For NDE, we ask whether the $X = 1$ directly increases or decreases the likelihood of the $Y = y$, not through any mediators. For NIE, we ask whether the treatment (setting $X = 1$) increases or decreases the likelihood of $Y = y$ through mediators, not directly.

A.2.1.3 Collection of Causal Graphs

We include all the ten causal graphs with treatment-effect pairs (CGTEs) in Figure A.1.

Note that one causal graph can have several different CGTEs, such as the confounding structure, which has three CGTEs: confounding, mediation, and collision in the triangle form. To generate all the causal graphs and CGTEs here, we iterate all commonly used ones within four nodes in the CI books, and omit CGTEs whose solution by CI methods

trivially resembles existing ones.

A.2.1.4 Data Coverage

Starting from the full set of 12 distinct causal graphs and 10 query types, there are a few combinations that must be omitted as the ground truth answer would be trivial or ill-defined. For example, in the “Immorality” graph, the treatment “X” and outcome “Y” are by construction statistically independent, so their correlation is necessarily 0. Similarly, there are several graphs where certain causal queries are ill-defined or don’t make sense to ask. Specifically:

1. For the Natural Direct Effect, we only include questions on the “IV”, “Arrowhead”, “Confounding”, “Mediation” and “DiamondCut” graphs.
2. For the Natural Indirect Effect, we only include questions on the “Mediation”, “Front-door”, “Arrowhead”, “Diamond” and “Chain” graphs.
3. For the Collider Bias and Explaining Away effect, we only include questions on the “Collision” graph.
4. For the Average Treatment Effect, we include questions on all graphs except “Collision”.
5. For the (deterministic) Counterfactuals, we include questions on all graphs except “Collision”.
6. For the Average Treatment Effect on the Treated (ATT), we include questions on all graphs except “Collision” and “IV”.

The “balanced” benchmark (main benchmark in v1.5), containing 10,112 questions split between all stories, graphs, query types, and commonsensicalness, is balanced such that there are roughly the same number of questions for each distinct story-graph-query combination (ranging from 50-100 per combination) across the different variants: commonsense, anticommonsense, and nonsense. Furthermore, we balance the distribution of correct answers so that there are the same number of “yes”s and “no”s.

The “aggregate” variant (main benchmark in v1.0) contains 10,560 questions and is primarily balanced across all stories. However since the number of stories for each variant (commonsense, anticommonsense, and nonsense) varies significantly, the results in an unbalanced benchmark in terms of sensicalness.

A.2.1.5 Query Form and Text Templates

We provide in Table A.8 the text templates we use for each query type.

A.2.1.6 Nonsensical Stories

To come up with a collection of nonsensical variable names, we use GPT-4 to generate some meaningless words. Specifically, we use the prompt: “Create 100 non-existent words that are short, i.e., within 5-characters.”, with temperature=0 with the OpenAI interface. The collection of nonsensical words we later use as variable names are as follows: ziblo, truq, fyze, glimx, jorv, wexi, snov, yupt, kraz, qixy, vubr, chiz, pliv, moxa, fygo, rukz, tasp, xevo, jyke, wibl, zorf, quzy, nyrp, gwex, smez, vyzt, hupx, cwoj, lirf, ovka, pexu, yigz, twaz, kwox, zuph, fraq, jyxo, swoy, uvzi, nekl, gyzp, rixq, vwem, xyfu, blyz,

qwip, zeku, tijv, yomx, hwaz, czix, plof, muvy, fyqo, rujz, tasb, xevi, jyka, wibm, zorx, quzw, nyro, gwet, smeus, vyta, hupz, cwoi, lirr, ovki, pexy, yigw, twac, kwoz, zupj, fraq, jyxi, swoq, uvzo, nekm, gyzl, rixw, vwen, xyfo, blyx, qwiu, zeky, tijw, yomz, hwax, czir, ploz, muvq, fyqi, rujx, tasn, xevu, jyko, wibp, zory, and quzt.

A.2.1.7 Anti-Commonsensical Stories

For the anti-commonsensical stories, we randomly do one of the actions:

1. Replace the effect variable Y with an attribute that would not be an effect variable in any of the stories. Such replacement variables include: "lip thickness", "earthquakes", "lactose intolerance", "rainfall", "is allergic to peanuts", "brown eyes", "curly hair", "black hair", "foot size", "freckles"
2. Create an irrelevant treatment variable X that does not play a causal role in any of our commonsensical stories. Such as: "can swim", "is religious", "has a brother", "has visited England", "likes spicy food", "is vegetarian", "speaks english", "drinks coffee", "plays card games", "listens to jazz", "solar eclipse", "has a sister", "full moon"

To transform a commonsensical story into an anti-commonsensical story, we apply one of these replacements sampled uniformly, resulting in stories such as:

- Ability to swim has a direct effect on studying habit and exam score. Studying habit has a direct effect on exam score.
- Gender has a direct effect on department competitiveness and peanut allergy. Department competitiveness has a direct effect on peanut allergy.
- Liking spicy food has a direct effect on relationship status. Appearance has a direct effect on relationship status.
- Playing card games has a direct effect on diabetes and lifespan. Smoking has a direct effect on diabetes and lifespan. Diabetes has a direct effect on lifespan. Smoking is unobserved.

For a full list of the replacements and how the replacements are made, check out the code.

A.2.1.8 Explanation Template

Step ① Extract the causal graph: The causal graph expressed in the context is: " \mathcal{G} ".

Step ② Identify the query type: The query type of the above question is "query_type".

Step ③ Formulate the query to its symbolic form: The formal form of the query is "symbolic_expression".

Step ④ Collect all the available data: The available data are: " \mathbf{d} ".

Step ⑤ Derive the estimand: Based on the graph structure and causal query, the question can be simplified into estimand "est".

Step ⑥ Solve for the estimand: Plug in the available data " \mathbf{d} " into "est".
 $\text{est}(\mathbf{d})$

$\approx \text{float}(a)$

Since the estimate for the estimand is $\text{float}(a)$, the overall answer to the question is $\text{bool}(a)$.

A.2.2 Experimental Details

A.2.2.1 CausalCoT Prompt

Q: [question from the dataset]

Guidance: Address the question by following the steps below:

Step 1) Extract the causal graph: Identify the causal graph that depicts the relationships in the scenario. The diagram should simply consist of edges denoted in "var1 -> var2" format, separated by commas.

Step 2) Determine the query type: Identify the type of query implied by the main question. Choices include "marginal probability", "conditional probability", "explaining away effect", "backdoor adjustment set", "average treatment effect", "collider bias", "normal counterfactual question", "average treatment effect on treated", "natural direct effect" or "natural indirect effect". Your answer should only be a term from the list above, enclosed in quotation marks.

Step 3) Formalize the query: Translate the query into its formal mathematical expression based on its type, utilizing the "do(\cdot)" notation or counterfactual notations as needed.

Step 4) Gather all relevant data: Extract all the available data. Your answer should contain nothing but marginal probabilities and conditional probabilities in the form "P(...)=..." or "P(... | ...)=...", each probability being separated by a semicolon. Stick to the previously mentioned denotations for the variables.

Step 5) Deduce the estimand using causal inference: Given all the information above, deduce the estimand using skills such as do-calculus, counterfactual prediction, and the basics of probabilities. Answer step by step.

Step 6) Calculate the estimand: Insert the relevant data in Step 4 into the estimand, perform basic arithmetic calculations, and derive the final answer. There is an identifiable answer. Answer step by step.

A: [LLM previous response]

Q: Based on all the reasoning above, output one word to answer the initial question with just "Yes" or "No".

A: [LLM final answer]

A.2.3 Additional Technical Background for Preliminaries

A.2.3.1 Graphical Models

We adopt the causal inference framework described in (Pearl, 2009b). A causal graph $G := (V, E)$ consists of a set of k vertices $V : \{V_1, \dots, V_k\}$ and directed edges $E := \{e_{ij}\}$, where the existence of each e_{ij} means that there is a direct causation from V_i to V_j , also denoted as $V_i \rightarrow V_j$. We also introduce some notations to describe the relative positions among the nodes. Following a standard assumption in causality (but see, e.g., (Bongers et al., 2021)), we will assume that G is a direct acyclic graph (DAG), where we denote the *parents* of a node V_i as $\text{PA}(V_i) := \{V_j | e_{ij} \in E\}$. We denote *descendants* $\text{DE}(V_i) := \{V_j | V_j \rightarrow \dots \rightarrow V_i \in E\}$ of a node V_i as all the nodes that have at least one direct path leading to a node. We call a node V_k as a *confounder* (i.e., common cause) of the other two nodes V_i and V_j if $e_{ki}, e_{kj} \in E$; a *collider* (i.e., common effect) if $e_{ik}, e_{jk} \in E$; and a *mediator* if $e_{ik}, e_{kj} \in E$.

Among all the variables in V , we use X and Y to denote two special variables, the treatment and effect, respectively.

A.2.3.2 Illustration of the Three Rungs of the Causal Ladder

In Figure A.2, we illustrate the difference among the three rungs by enumerating what actions are performed on the variables other than target variables X and Y .

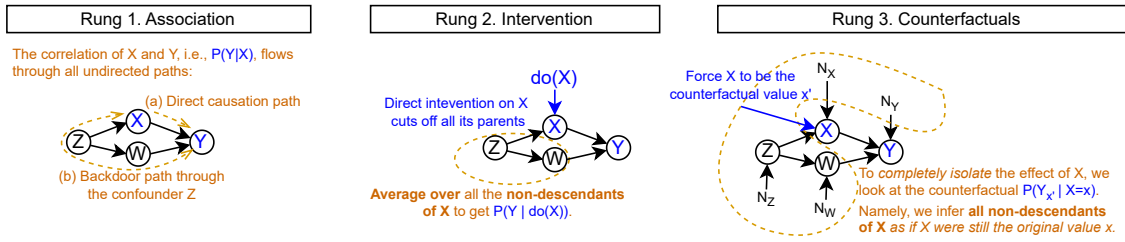


Figure A.2: The Causal Ladder consists of three rungs: association, intervention and counterfactuals. We color in blue the treatment X and effect Y , as well as the actions on X . We color in orange words about how to get the estimand, and we use the orange circle to include all the non-descendants of X .

A.2.3.3 Causal Inference Methods

We introduce do-calculus which can downgrade the Rung-2 queries to Rung-1 quantities when it is applicable, and counterfactual predictions which downgrade the Rung-3 queries.

A.2.3.3.1 Do-Calculus

Do-Operator as a Notation As mentioned in Rung 2, the do-operator is a convenient notation to represent an intervention on a variable. For example, $\text{do}(X = x)$ sets the value of variable X to x .

Three Inference Rules for Climbing the Ladder Do-calculus is a set of rules that allows us to answer higher-rung questions using lower-rung quantities, such as probability dis-

tributions of Rung 1. Given a causal graphical model with and four disjoint sets of variables X , Y , Z , and W , and a joint probability distribution that is Markov and faithful to the graph, do-calculus contains the following three rules:

Rule 1 (Insertion/deletion of observations):

$$P(Y|do(X), Z, W) = P(Y|do(X), W), \quad (A.1)$$

if Y and Z are d-separated by $X \cup W$ in G^* , the graph obtained from \mathcal{G} by removing all arrows pointing into variables in X .

Rule 2 (Action/observation exchange):

$$P(Y|do(X), do(Z), W) = P(Y|do(X), Z, W), \quad (A.2)$$

if Y and Z are d-separated by $X \cup W$ in G^\dagger , the graph obtained from \mathcal{G} by removing all arrows pointing into variables in X and all arrows pointing out of variables in Z .

Rule 3 (Insertion/deletion of actions):

$$P(Y|do(X), do(Z), W) = P(Y|do(X), W), \quad (A.3)$$

if Y and Z are d-separated by $X \cup W$ in G^\ddagger , the graph obtained from \mathcal{G} by first removing all arrows pointing into variables in X (thus creating G^*) and then removing all arrows pointing into variables in Z that are not ancestors of any variable in W in G^* .

These rules are sound and complete (Shpitser and Pearl, 2006b). Namely, iff we have all the terms on the right hand side, then the causal term on the left hand side is identifiable.

Example Application of Do-Calculus Taking the example in Figure 3.2, g_1 maps the query type ATE to its symbolic expression $\mathbb{E}[Y|do(X = 1)] - \mathbb{E}[Y|do(X = 0)]$.

Next, g_2 further simplifies the estimand given the confounding graph, as in the flow chart in the middle of Figure 3.2:

$$ATE := \mathbb{E}[Y|do(X = 1)] - \mathbb{E}[Y|do(X = 0)] \quad (A.4)$$

$$= \sum_z P(Z = z) [\mathbb{E}(Y|X = 1, Z = z) - \mathbb{E}(Y|X = 0, Z = z)], \quad (A.5)$$

which which resolves all the $do(\cdot)$ terms to probability terms. This example shows the famous backdoor adjustment in do-calculus (Pearl, 1995).

A.2.3.3.2 Three Steps for Counterfactual Prediction

Given a SCM M , distribution on the exogenous variables $P(u)$, and evidence e from the model $\langle M, P(u) \rangle$, the probability of the counterfactual "if X had been x then Y would have been y , given we observed e ," denoted $P(Y_x = y|e)$, can be evaluated using the following three steps (Pearl, 2009b):

Abduction: Update the probability distribution $P(u)$ by the evidence e to obtain $P(u|e)$

Action: Modify M by the action $do(X = x)$, i.e. replace X with $X = x$ in the structural equations, to obtain the modified SCM M_x

Prediction: Use the modified model $\langle M_x, P(u|e) \rangle$, to compute the probability of $Y = y$.

A.2.4 Previous Results on CLadder v1.0

A.2.4.1 Dataset Statistics for v1.0

	Total	Rung 1	Rung 2	Rung 3
Size				
# Samples	10,560	3,288	3,288	3,984
Question				
# Sentences/Sample	6.85	6.00	7.00	7.25
# Words/Sample	94.47	76.41	96.84	103.42
# Nodes/Graph	3.54	3.54	3.55	3.54
# Edges/Graph	3.44	3.41	3.43	3.46
Answer				
Positive Class (%)	50	50	50	50
Explanations				
# Sentences/Sample	13.11	12.04	13.76	13.83
# Words/Sample	146.82	141.88	147.88	151.30

Table A.9: Statistics of our CLADDER data v1.0.

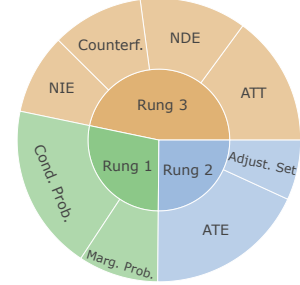


Figure A.3: Distributions of query types in our dataset v1.0.

Our data-generating procedure has the potential to algorithmically generate very large amounts of questions. In practice, we pick a dataset size that is large enough to be representative, and at the same time not too large to be problematic given the expensive inference costs of LLMs. We therefore set our dataset size to be 10K. We report the statistics of our dataset in Table A.9.

The dataset roughly balanced across the query types, graph structures, stories, and ground-truth answers (as seen in Figure A.3). Note that there are some slight adjustments such as more samples for ATE because it allows us to test various techniques, including back-door and front door adjustments. More details on our design choices can be found in Appendix A.2.1.4.

A.2.4.2 Main Results on v1.0

We compare the performance of all models in Table A.10. First, we can see that the causal reasoning task in CLADDER is in general very challenging for all models. And models such as the earlier, non-instruction-tuned GPT-3 and both LLaMa and Alpaca are no better than random performance. With instruction-tuning, models start to show some improvement. And amongst all, our CAUSALCoT achieves the highest performance of 66.64%, which is 2.36 points better than vanilla GPT-4.

Moreover, from the accuracy by empirical alignment level in Table A.10, we can see that the original GPT-4 model performs the best on commonsensical data, but 5.34 points worse on nonsensical data. However, our CAUSALCoT enhances the reasoning ability across all levels, with substantial improvement on anti-commonsensical data and non-sensical data, indicating that CAUSALCoT is particularly beneficial on unseen data.

A.2.4.3 Ablation Study on v1.0

We conduct an ablation study for our multi-step CAUSALCoT. We ablate each of the four subquestions, and observe in Table A.11 that classifying the query type and formalizing it has the most effect on

	Acc.
CAUSALCoT	66.64
w/o Step ①	64.54
w/o Step ②	63.74
w/o Step ③	63.43
w/o Step ④	64.47

Table A.11: Ablation study.

the model’s performance, which might be because that they are the crucial formalization step in order to do the causal inference correctly. Meanwhile, removing Steps ① and ④, which are mostly about parsing the prompt correctly, have the least impact on performance.

A.2.5 More Experiments

A.2.5.1 Details of Our Error Analysis

For Step 2 about the query type prediction, we report the overall F1 classification score, and also F1 by rungs. For the rest of the steps, we manually annotate the correctness of 100 samples of CAUSALCoT. We report the correctness of est by accuracy, and the correctness of the predicted set of available data by taking the F1 with the ground-truth d . For Step 5, we report the accuracy of whether the model simplifies the estimand correctly to est' using causal inference, and also arithmetic correctness (Arith.).

A.2.5.2 ROSCOE Evaluation

We employed the ROSCOE suite of evaluation metrics on step-by-step text reasoning, as introduced by (Golovneva et al., 2022), to automate the evaluation of the outputs from CAUSALCoT on 2,000 randomly sampled questions from our dataset. Differing from conventional metrics, ROSCOE is specifically designed to scrutinize the quality of large language model outputs, focusing on aspects such as semantic consistency, logicity, informativeness, fluency, and factuality, all evaluated within the context of step-by-step reasoning, rather than solely the final response. This allows for a more objective and comprehensive assessment of a model’s output, greatly aiding in the verification of its interpretability. The results of this evaluation can be found in Table A.12 and Figure A.4. We consider the model’s performance as unsatisfying if it falls out of the top quantile, namely receiving a score $s \in [0, 1]$ smaller than 0.25 when the score should be minimized, or greater than 0.75 when it should be maximized.

We can see in the plot that the good-performing aspects are faithfulness to the original question, reasoning alignment with the ground truth, and absence of external hallucinations, which are consistently within the top quantile. This suggests that the model carries out accurate reasoning within the constraints of the fictitious world introduced in each question.

However, there are some performance dips in redundancy, perplexity chain, and missing step metrics. The first two could potentially be attributed to complex elements such as graph notation, while the relatively lower “missing step” score warrants further investigation. Despite these observations, this analysis largely aligns with our qualitative understanding of the models’ good response ability in answering causal questions in our dataset.

	Example 1 (Label=Negative)	Example 2 (Label=Positive)
Symbolic Form	<p><i>Premise:</i> Suppose there is a closed system of 2 variables, A and B. All the statistical relations among these 2 variables are as follows: A correlates with B.</p> <p><i>Hypothesis:</i> A directly affects B.</p> <p><i>Relation between the promise and hypothesis:</i> The premise does not a necessary condition for the hypothesis.</p>	<p><i>Premise:</i> Suppose there is a closed system of 3 variables, A, B and C. All the statistical relations among these 3 variables are as follows: A correlates with C. B correlates with C. However, A is independent of B.</p> <p><i>Hypothesis:</i> A directly affects C.</p> <p><i>Relation between the promise and hypothesis:</i> The premise is a necessary condition for the hypothesis. So if the premise is true, the hypothesis must be true.</p>
Natural Story	<p><i>Premise:</i> Suppose there is a closed system of 2 variables, ice cream sales and swimming pool attendance. All the statistical relations among these 2 variables are as follows: ice cream sales correlate with swimming pool attendance.</p> <p><i>Hypothesis:</i> Ice cream sales directly affect swimming pool attendance.</p> <p><i>Relation between the premise and hypothesis:</i> The premise does not provide a necessary condition for the hypothesis. The correlation between ice cream sales and swimming pool attendance could be due to a third variable, such as hot weather, which increases both ice cream sales and swimming pool attendance. Therefore, it is not necessarily true that ice cream sales directly affect swimming pool attendance.</p>	<p><i>Premise:</i> Let's consider three factors: eating junk food (A), obesity (C), and watching television (B). There is a correlation between eating junk food and obesity, and between watching television and obesity. However, eating junk food and watching television are independent from each other.</p> <p><i>Hypothesis:</i> Eating junk food directly affects obesity.</p> <p><i>Relation between the premise and hypothesis:</i> The premise provides the necessary conditions for the hypothesis. It establishes the independent variables A (eating junk food) and B (watching television) and their correlations with obesity. Given that these are true, it supports the hypothesis that eating junk food directly affects obesity.</p>

Table A.1: Examples of natural stories generated based on the symbolic form in our CLADDER dataset, showing the broad application value of our dataset as the starting point for various verbalizations of the correlation-to-causation inference task.

Test Set Size	102
Dev Set Size	102
# Tokens/Premise	64.88
# Tokens/Hypothesis	13.54
# Tokens/Explanation	64.66
% Positive Labels	1.67

Table A.2: Statistics of our generated natural stories. We report the number of samples in the test and development sets; number of tokens per premise (# Tokens/Premise), hypothesis (# Tokens/Hypothesis), and explanation (# Tokens/Explanation); and percentage of the positive labels (% Positive Labels).

Causal Relation	Hypothesis Template
Is-Parent	{Var i} directly causes {Var j}.
Is-Ancestor	{Var i} causes something else which causes {Var j}.
Is-Child	{Var j} directly causes {Var i}.
Is-Descendant	{Var j} is a cause for {Var i}, but not a direct one.
Has-Collider	There exists at least one collider (i.e., common effect) of {Var i} and {Var j}.
Has-Confounder	There exists at least one confounder (i.e., common cause) of {Var i} and {Var j}.
<i>Paraphrases</i>	
Is-Parent	{Var i} directly affects {Var j}.
Is-Ancestor	{Var i} influences {Var j} through some mediator(s).
Is-Child	{Var j} directly affects {Var i}.
Is-Descendant	{Var j} influences {Var i} through some mediator(s).
Has-Collider	{Var i} and {Var j} together cause some other variable(s).
Has-Confounder	Some variable(s) cause(s) both {Var i} and {Var j}.

Table A.3: Templates and their paraphrases for each causal relation in the hypothesis. We use {Var i} and {Var j} as placeholders for the two variables.

Two Equivalent Forms	Duplication Property	De-Duplication Method
<ul style="list-style-type: none"> { Is-Parent(i, j) Is-Child(j, i) 	Two exact same strings	Keep only one, by forcing $i < j$
<ul style="list-style-type: none"> { Is-Ancestor(i, j) (Original) Is-Descendent(j, i) (Original) 	Two different strings, but semantically equivalent	Randomly sample one out of the two
<ul style="list-style-type: none"> { Is-Ancestor(i, j) (Paraphrased) Is-Descendent(j, i) (Paraphrased) 	Two exact same strings	Keep only one, by forcing $i < j$
<ul style="list-style-type: none"> { Has-Collider(i, j) Has-Collider(j, i) 	Two different strings, but semantically equivalent	Randomly sample one out of the two
<ul style="list-style-type: none"> { Has-Confounder(i, j) Has-Confounder(j, i) 	Two different strings, but semantically equivalent	Randomly sample one out of the two

Table A.4: De-duplication methods for the six causal relation types and their verbalizations.

N-Gram	PMI w/ Non-Ent. Label	PMI w/ Ent. Label	Diff
a cause	1.692209	-1.025611	2.717820
a cause for	1.663640	-0.983790	2.647430
A causes	1.640679	-0.951610	2.592289
A causes something	1.621820	-0.926075	2.547895
a direct	1.606052	-0.905316	2.511369
a direct one	1.592673	-0.888107	2.480781
for D	1.584826	-0.878180	2.463006
for D but	1.583897	-0.877014	2.460911
for E	1.582980	-0.875864	2.458844
for E but	1.582074	-0.874728	2.456802

Table A.5: PMI between the labels and n-grams. The labels include non-entailment (Non-Ent.) and entailment (Ent.). And the n-grams include all with no more than four words. The |Diff| column shows the absolute value of the difference between the PMIs with two labels. We show the top 10 n-grams with the largest differences of their PMIs with the two classes in the |Diff| column.

Selected Models	Relation Type	F1	Precision	Recall	Accuracy
GPT-3.5	All	21.69	17.79	27.78	69.46
GPT-3.5	Is-Parent	8.82	100	4.62	83.47
GPT-3.5	Is-Ancestor	0	0	0	90.67
GPT-3.5	Is-Child	9.84	100	5.17	85.33
GPT-3.5	Is-Descendant	14.29	11.9	17.86	84
GPT-3.5	Has-Collider	34.24	25.51	52.07	35.12
GPT-3.5	Has-Confounder	15.33	8.86	56.76	37.8
GPT-4	All	29.08	20.92	47.66	64.6
GPT-4	Is-Parent	0	0	0	82.67
GPT-4	Is-Ancestor	30.77	31.25	30.3	88
GPT-4	Is-Child	0	0	0	84.53
GPT-4	Is-Descendant	26.98	17.35	60.71	75.47
GPT-4	Has-Collider	44.1	30.18	81.82	32.71
GPT-4	Has-Confounder	20.67	11.53	100	23.86
RoBERTa MNLI	All	22.79	34.73	16.96	82.5
RoBERTa MNLI	Is-Parent	0	0	0	82.67
RoBERTa MNLI	Is-Ancestor	0	0	0	91.2
RoBERTa MNLI	Is-Child	0	0	0	84.53
RoBERTa MNLI	Is-Descendant	0	0	0	92.53
RoBERTa MNLI	Has-Collider	43.45	39.73	47.93	59.52
RoBERTa MNLI	Has-Confounder	0	0	0	84.45

Table A.6: Fine-grained evaluation results for some selected non-fine-tuned models.

	F1	Precision	Recall	Accuracy
GPT-3.5 (plain query; original)	21.69	17.79	27.78	69.46
GPT-3.5 (enhanced query)	25.44	17.29	48.11	52.01

Table A.7: Performance of GPT-3.5 with different queries. We quote the original performance from Table 2.4.

Query Type	Symbolic Expression	Ex-pression	Natural Language Question Template
Rung 1: Association			
Marg. Prob.	$P(Y)$		Is the overall likelihood of $\{v_{\text{noun}}(X = 1)\}$ greater than chance?
Cond. Prob.	$P(Y X)$		Is the chance of $\{v_{\text{noun}}(Y = 1)\}$ larger when observing $\{v_{\text{noun}}(X = 1)\}$?
Rung 2: Intervention			
ATE	$E[Y \text{do}(X = 1)] - E[Y \text{do}(X = 0)]$		Will $\{v_{\text{noun}}(X = 1)\}$ increase the chance of $\{v_{\text{noun}}(Y = 1)\}$?
Adjust. Set	If S opens a backdoor path		To understand how $\{v_{\text{overall}}(X)\}$ affects $\{v_{\text{overall}}(Y = 1)\}$, should we look directly at how $\{v_{\text{overall}}(X)\}$ correlates with $\{v_{\text{overall}}(Y)\}$ in general, or this correlation case by case according to $\{v_{\text{overall}}(S)\}$?
Rung 3: Counterfactuals			
Counterf. Prob.	$P(Y_x = y)$		Can we infer that $\{v_{\text{sent}}(Y = 1)\}$ had it been that $\{v_{\text{cond}}(X = 1)\}$ instead of $X=0$?
ATT	$E[Y_1 - Y_0 X = 1]$		For $\{v_{\text{attr}}(X = 1)\}$, would it be more likely to see $\{v_{\text{noun}}(Y = 1)\}$ $\{v_{\text{cond}}(X = 0)\}$?
NDE	$E[Y_{1,M_0} - Y_{1,M_0}]$	–	If we disregard the mediation effect through $\{v_{\text{overall}}(Y = 1)\}$, would $\{v_{\text{noun}}(X = 1)\}$ still positively affect $\{v_{\text{noun}}(Y = 1)\}$?
NIE	$E[Y_{0,M_1} - Y_{0,M_0}]$	–	Does $\{v_{\text{overall}}(X)\}$ affect $\{v_{\text{overall}}(Y)\}$ through $\{v_{\text{overall}}(\text{OtherVars})\}$?

Table A.8: Example natural language templates for each query type.

	Overall Acc.	Acc. by Rung			Acc. by Empirical Alignment		
		1	2	3	Anti-C.	Nonsens.	Comm.
Random	49.27	50.28	48.40	49.12	49.69	49.01	49.12
LLaMa	45.22	63.33	31.10	41.45	45.31	45.21	45.12
Alpaca	45.54	63.33	31.57	41.91	45.94	45.21	45.49
GPT-3 Non-Instr. (davinci)	47.42	63.88	32.99	44.89	47.0	48.28	46.97
GPT-3 Instr. (text-davinci-001)	57.07	63.95	63.63	48.04	59.12	57.81	54.28
GPT-3 Instr. (text-davinci-002)	56.24	46.03	69.55	55.04	54.75	59.65	54.31
GPT-3 Instr. (text-davinci-003)	62.69	58.0	80.83	54.52	63.93	62.09	62.05
GPT-3.5 (queried in May 2023)	61.71	65.12	69.9	54.11	65.43	55.15	64.55
GPT-4 (queried in May 2023)	64.28	53.94	81.87	63.11	65.75	60.87	66.21
+ CAUSALCoT	66.64	61.67	86.13	58.23	69.32	63.02	67.60

Table A.10: Performance of all models on our CLADDER dataset v1.0. We report the overall accuracy (Acc.), and also fine-grained accuracy by rung and by empirical alignment.

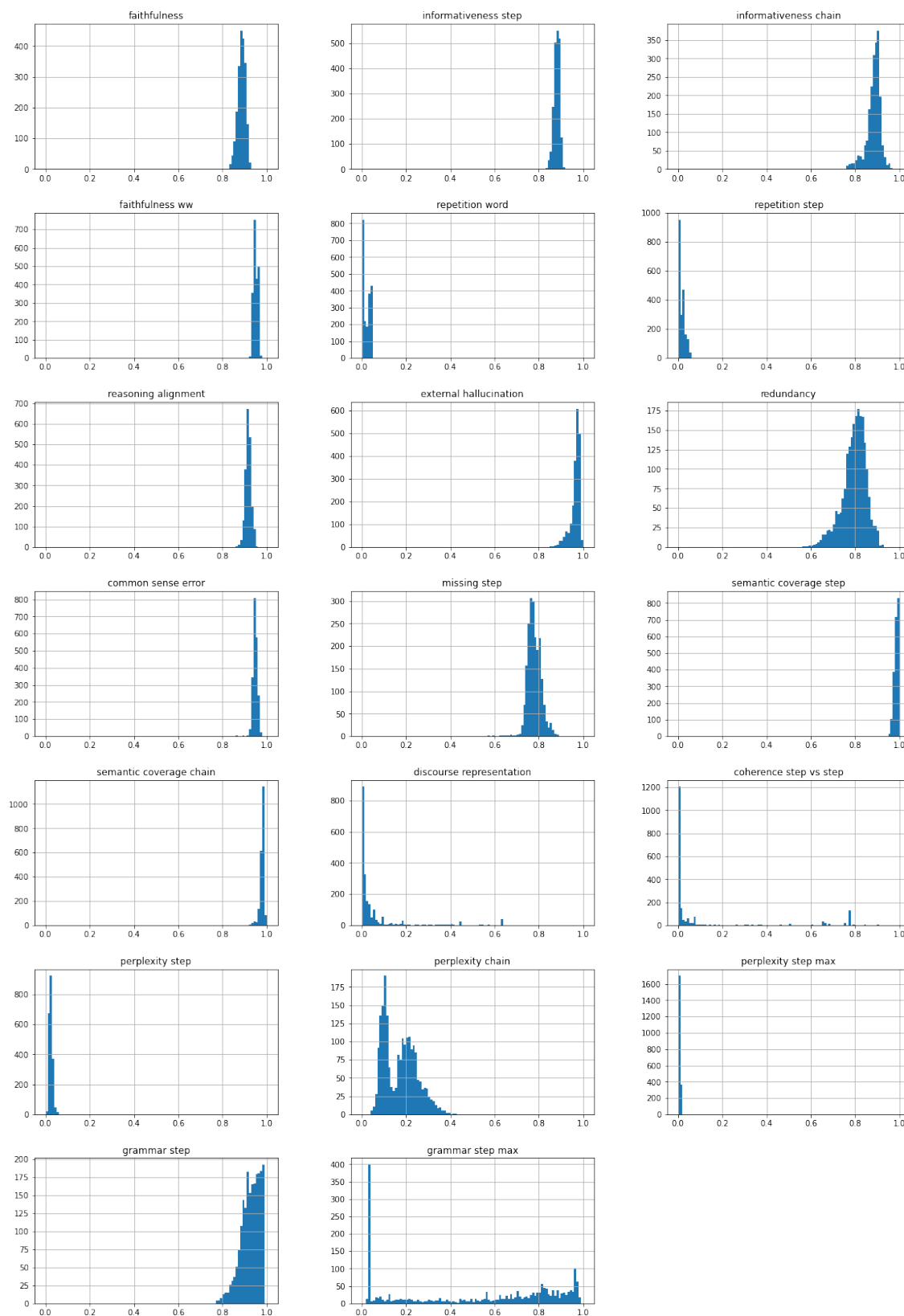


Figure A.4: ROSCOE scores of answers from CAUSALCoT on 2,000 randomly sampled questions from our dataset.

	Mean	Std	Min	25%	50%	75%	Max
Faithfulness	0.89	0.02	0.83	0.88	0.89	0.90	0.93
Informativeness Step	0.88	0.01	0.83	0.87	0.88	0.89	0.92
Informativeness Chain	0.88	0.03	0.76	0.87	0.89	0.90	0.96
Faithfulness Word	0.95	0.01	0.92	0.94	0.95	0.96	0.97
Repetition Word	0.02	0.02	-0.00	0.00	0.02	0.04	0.05
Repetition Step	0.02	0.01	-0.00	0.00	0.01	0.03	0.06
Reasoning Alignment	0.92	0.01	0.86	0.91	0.92	0.93	0.95
External Hallucination	0.97	0.02	0.84	0.96	0.97	0.98	0.99
Redundancy	0.80	0.05	0.56	0.77	0.80	0.83	0.92
Common Sense Error	0.95	0.01	0.86	0.94	0.95	0.96	0.98
Missing Step	0.78	0.03	0.58	0.76	0.78	0.80	0.88
Semantic Coverage Step	0.99	0.01	0.95	0.98	0.99	0.99	1.00
Semantic Coverage Chain	0.98	0.01	0.93	0.98	0.98	0.99	0.99
Discourse Representation	0.06	0.13	0.00	0.01	0.01	0.05	0.67
Coherence Step Vs Step	0.14	0.27	0.00	0.00	0.01	0.07	0.94
Perplexity Step	0.02	0.01	0.00	0.02	0.02	0.03	0.07
Perplexity Chain	0.17	0.07	0.05	0.11	0.17	0.23	0.42
Perplexity Step Max	0.00	0.00	0.00	0.00	0.00	0.01	0.02
Grammar Step	0.93	0.04	0.77	0.90	0.93	0.96	0.99
Grammar Step Max	0.53	0.35	0.02	0.12	0.65	0.85	0.99

Table A.12: Statistics of ROSCOE scores evaluated on answers from CAUSALCoT on 2,000 randomly sampled questions from our dataset.

A.2.6 Comparison with Existing Causality-Related Datasets

We show in Table A.13 the distinction of our work from all existing causality-related datasets that address either the causality-as-knowledge task, or the causality-as-language-comprehension task.

	Question Types				Skill Types		
	Assoc.	Interv.	Counterf.	CI Method	Formalization of Causal Queries	Causal RE	Qualitative Reason- ing
<i>Datasets for Causality as Knowledge (Commonsense Causality)</i>							
COPA (2012)	✗	✓	✗	✗	✗	✗	✗
Event2Mind (2018)	✗	✓	✗	✗	✗	✗	✗
ATOMIC (2019a)	✗	✓	✗	✗	✗	✗	✗
SocialIQA (2019b)	✗	✓	✗	✗	✗	✗	✗
TimeTravel (2019)	✗	✓	✗	✗	✗	✗	✗
Goal-Step (2020a)	✗	✓	✗	✗	✗	✗	✗
Abductive (ART) (2020)	✗	✓	✗	✗	✗	✗	✗
Com2Sense (2021b)	✗	✓	✗	✗	✗	✗	✗
CRASS (2022)	✗	✗	✓	✗	✗	✗	✗
<i>Datasets for Causality as Language Comprehension (Causal Relation Extraction)</i>							
SemEval2021 Task8 (2010)	✗	✗	✗	✗	✗	✓	✗
EventCausality (2011)	✗	✗	✗	✗	✗	✓	✗
Causal-TimeBank (2014)	✗	✗	✗	✗	✗	✓	✗
CaTeRS (2016)	✗	✗	✗	✗	✗	✓	✗
BECauSE (2017)	✗	✗	✗	✗	✗	✓	✗
TellMeWhy (2021)	✗	✗	✗	✗	✗	✓	✗
<i>Datasets for Formal Causal Reasoning</i>							
Corr2Cause (Jin et al., 2024)	✗	✓	✗	✓	✓	✗	✗
CLADDER (Ours)	✓	✓	✓	✓	✓	✓	✓

Table A.13: Comparison of our dataset and existing causal or reasoning datasets. The aim of our dataset is to test the pure reasoning ability of LLMs on causal questions. For each dataset, we first identify whether its question types cover the three rungs: association (Assoc.), intervention (Interv.), and counterfactuals (Counterf.). We also check what skill types the dataset tests: the application of causal inference methods (CI Method), formalization of causal queries, causal relation extraction from the given text (Causal RE), and qualitative reasoning.

A.3 Additional Materials for Chapter 4

A.3.1 Experiments for Pythia-6.9b

This section extends the experimental analysis conducted on GPT-2 to Pythia-6.9b. The goal is to replicate the prior methodology and compare the outcomes across the two different models, thus contributing to a broader understanding of model behaviors under similar conditions.

A.3.1.1 Macroscopic Inspection across Layers and Token Positions

Figure A.5 provides a comparative analysis of the logit values for two specific tokens, labeled as factual and counterfactual, across various positions and layers in Pythia-6.9b.

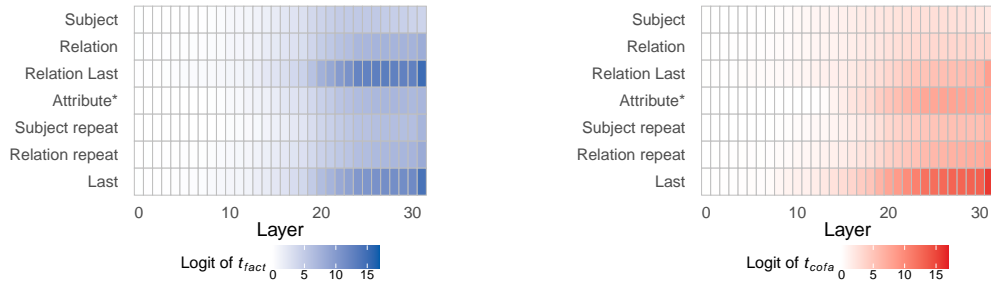


Figure A.5: **Layer-wise Position Analysis of Relevant Tokens in GPT-2-small.** The figure presents the logit values for two pertinent tokens across various positions and layers. The left panel illustrates the logit values for the factual token t_{fact} , while the right panel illustrates the logit values for the counterfactual token t_{cofa} .

A.3.1.2 Intermediate Inspection of Attention and MLP Blocks

This subsection exposes the contributions of Attention and MLP Blocks to the differences in logit values across layers within Pythia-6.9b. Figure A.6 explores how these components influence the computation of logits for two tokens, represented as the difference $\Delta_{\text{cofa}} = \text{Logit}(t_{\text{cofa}}) - \text{Logit}(t_{\text{fact}})$ at the final position of the input. The analysis specifically highlights the distinct effects of these blocks at different stages of the model's operation.

A.3.1.3 Microscopic Inspection of Individual Attention Heads

Figure A.7 quantifies the direct contributions of all attention heads to the difference in logit values, labeled as Δ_{cofa} . It specifically identifies heads that preferentially enhance the logits for t_{fact} (shown in blue) versus those favoring t_{cofa} (depicted in red), offering insights into how attention mechanisms differentially prioritize token attributes.

Figure A.8 presents the attention patterns of the relevant attention heads at the last token position. It shows the consistent pattern of the relevant heads, with a consistent focus on the attribute position.

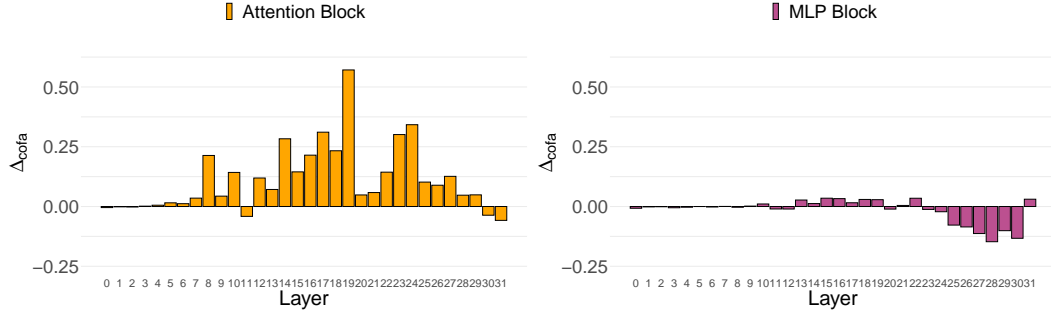


Figure A.6: **Attribution of Logit Differences to Attention and MLP Blocks.** delineates the differential impact of Attention and MLP Blocks on logit values at the terminal input position. The attention mechanism is shown to predominantly influence early layer processing in the left panel, while the right panel details the increased contribution of MLP Blocks to the factual token’s logits in the concluding layers, illustrating the dynamic interplay between these fundamental neural network elements.

A.3.2 Other Experiment for GPT-2

A.3.2.1 Ranks Analysis in the Last Position

We provide additional information in Figure A.9 mapping the logits to ranks of the tokens, and find that the rank of t_{cofa} in the projected logit distribution remains very low: t_{cofa} is among the 20 most likely tokens in the first five layers and between the 20th and the 70th in the last part of the network.

A.3.2.2 Attention Pattern of Relevant Attention Heads

Figure A.10 shows the full attention pattern for the relevant attention heads, as identified in Section 4.6. It is show as the attention pattern is similar between all the relevant attention heads, independently if the heads is favoring t_{fact} or t_{cofa} .

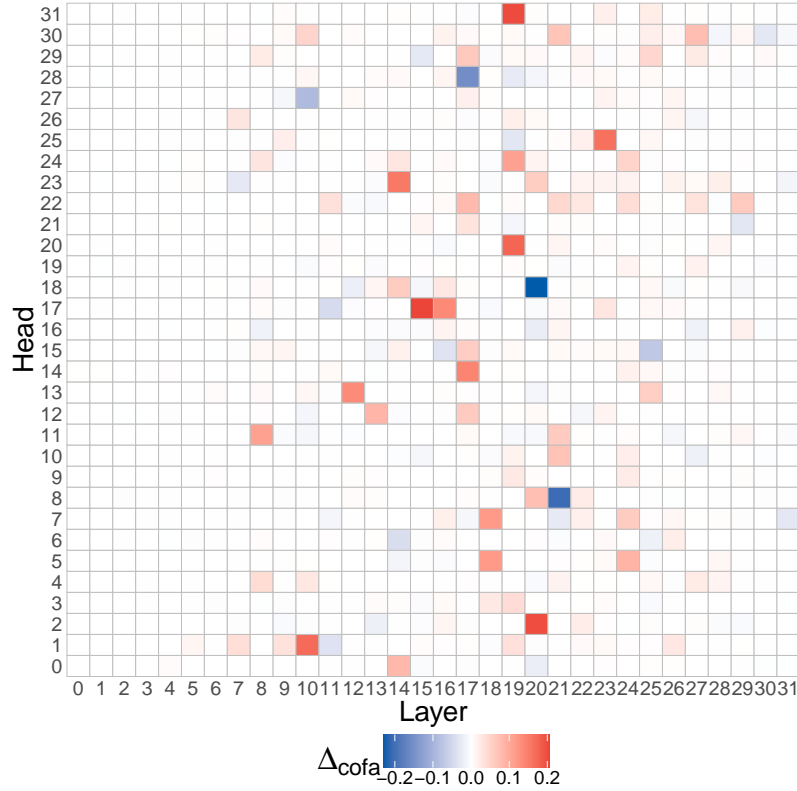


Figure A.7: **Direct Contribution of Attention Heads.** The figure displays the direct contribution of all heads in Pythia-6.9b to the logit difference Δ_{cofa} with heads favoring t_{fact} highlighted in blue and those favoring t_{cofa} in red.

A.4 Additional Materials for Chapter 5

A.4.1 Creation of the Prompts

We consider MWP examples from the union of the three datasets SVAMP, ASDiv-A, and MAWPS. The textual template t of a problem consists of a context (describing a real-world state and/or actions) and a question. In order to obtain suitable prompts for the models, we convert the problems' questions into statements where the result of the problem is expected to be the first token after the prompt. E.g., in the example in Section 5.2, *how many trees will he have?* is converted into *the number of trees that he will have is _*. From the MWP templates of the SVAMP/ASDiv-A/MAWPS collection (we consider all splits), we filter out the templates whose questions do not start with *How many...*, and we use spaCy¹ to identify the subject, the object and the verbs in the sentence. This allows us to convert the last sentence of the template from *The number of... is*. This way, we obtain 437 statement-based MWP templates for two-operand problems and 307 for three-operand problems. We manually checked a subset of the templates to identify possible mistakes in the conversion procedure.

¹<https://spacy.io>

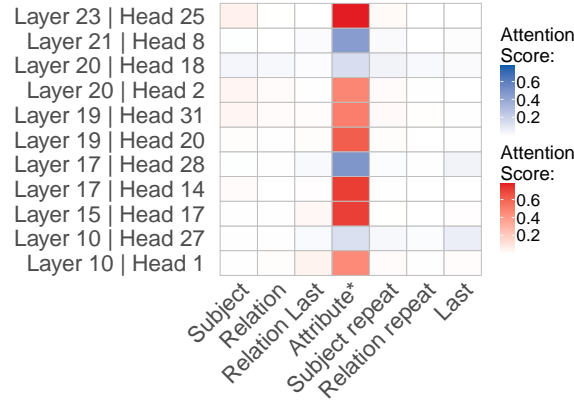


Figure A.8: **Attention Pattern for Relevant Attention Heads.** The panel illustrates the attention patterns of relevant heads for the last position, demonstrating consistent attention to the attribute position by both red and blue heads.

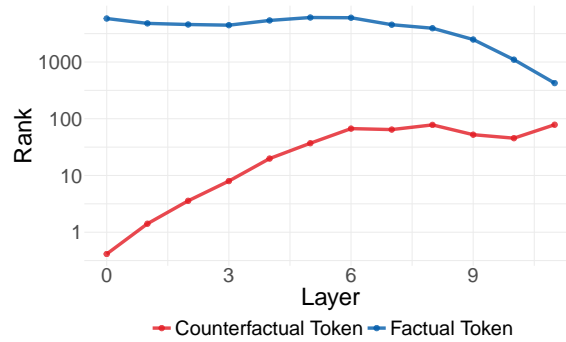


Figure A.9: **Rank of Target Tokens for Attribute Position Across Layers in GPT-2.** This figure depicts the trend where the logit rank for the factual token t_{fact} decreases while the rank for the counterfactual token t_{cofa} increases at the attribute position. In the concluding layers, this pattern is evident as t_{fact} typically secures a lower rank, in contrast to t_{cofa} , which shows an upward trajectory in rank. However, it is important to note that t_{cofa} 's rank consistently remains lower than that of t_{fact} .

A.4.2 Frequently Asked Questions

A.4.2.1 How do the intervention data look like?

In Table A.14 we report examples of MWP pairs representing different types of intervention.

A.4.2.2 What is the accuracy of the evaluated models on the generated problems?

We report the accuracy of the models considered for our main evaluation in terms of accuracy at 1 and accuracy at 10. Results are displayed in Figure A.11. The accuracy of the LLaMA models is 11.1%, 25.7%, 32.8%, and 13.0% respectively for the 7B, 13B, 30B, and Alpaca versions. The accuracy of the GPT-3 Davinci models on the three-operand problems is 2%, 11%, and 15% for the Instruct, Davinci-002, and Davinci-003 versions, respectively.

TCE($N \rightarrow R$)	Ruby has 87 candies. If she shares the candies among 29 friends, the number of candies that each friend gets is	$g = 87/29 = 3$
	Ruby has 35 candies. If she shares the candies among 5 friends, the number of candies that each friend gets is	$g = 35/5 = 7$
DCE($N \rightarrow R$)	The school is composed of 13 buildings each having 10 classrooms. The number of classrooms that the school has is	$g = 10 \times 13 = 130$
	The school is composed of 65 buildings each having 2 classrooms. The number of classrooms that the school has is	$g = 65 \times 2 = 130$
DCE($S \rightarrow R$)	The razorback t-shirt shop ordered 6 cases of t-shirts. If each case contains 17 t-shirts the number of t-shirts that they ordered is	$g = 17 \times 6 = 102$
	The roller coaster at the state fair costs 6 tickets per ride. If 17 friends were going to ride the roller coaster the number of tickets that they would need is	$g = 17 \times 6 = 102$
TCE($T \rightarrow R$)	Sean has 23 whistles. He has 6 more whistles than Charles. The number of whistles that Charles has is	$g = 23 - 6 = 17$
	Jovana filled her bucket with 23 pounds of shells. If she adds 6 more pounds of shell to fill her bucket, the number of pounds that she has is	$g = 23 + 6 = 29$

Table A.14: For each of the causal effects measured (left column), we report a pair of MWPs illustrating the intervention performed (center), along with their respective ground-truth result (left column).

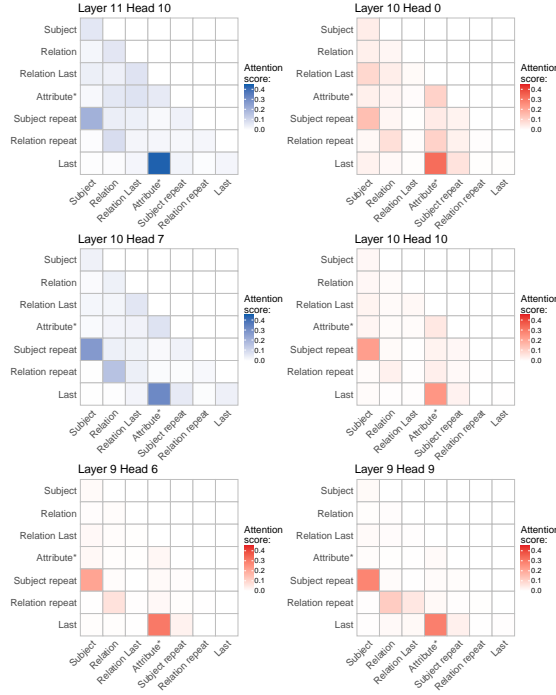


Figure A.10: **Attention Pattern of Significant Heads.** This figure illustrates the comprehensive attention pattern of heads substantially influencing $\Delta_{t_{\text{cofa}}}$. Notably, a similar pattern emerges for both heads favoring t_{cofa} (depicted in red) and those favoring t_{fact} (illustrated in blue), particularly in the attention edge between the attribute and the final position.

A.4.2.3 What is the relation between accuracy and the RCC metric?

We examine the relationship between performance and robustness, computing the Pearson correlation coefficient between accuracy (accuracy@10) and the relative confidence change (RCC) metric. On a per-template basis (500 instances for each template), we found accuracy to be positively correlated with $\text{TCE}(N \text{ on } R)$ and $\text{TCE}(T \text{ on } R)$ (0.24 and 0.49, respectively) and negatively correlated with $\text{DCE}(N \rightarrow R)$ and $\text{DCE}(S \rightarrow R)$ (-0.26 and -0.36, respectively). We see these results as a quantitative validation of the intuition behind our framework: the better the model's performance, the more the model tends to correctly adjust its prediction after a result-altering intervention (higher sensitivity) and to correctly not change its prediction after a result-preserving intervention (higher robustness).

Moreover, we conduct an additional sanity check as in Patel et al. (2021): removing the question from the MWP templates, we observe a sensitivity-robustness degradation to random guessing (i.e., $\text{TCE} \approx \text{DCE}$). This indicates that the measurement of the causal effects within our framework is not affected by patterns in the templates that might have been picked up or memorized by large models.

A.4.3 Computation of Causal Effects for GPT-3

We access GPT-3 through the OpenAI APIs, which allow a user to prompt the model and obtain the probabilities assigned by the model to the k -th most likely vocabulary entries,

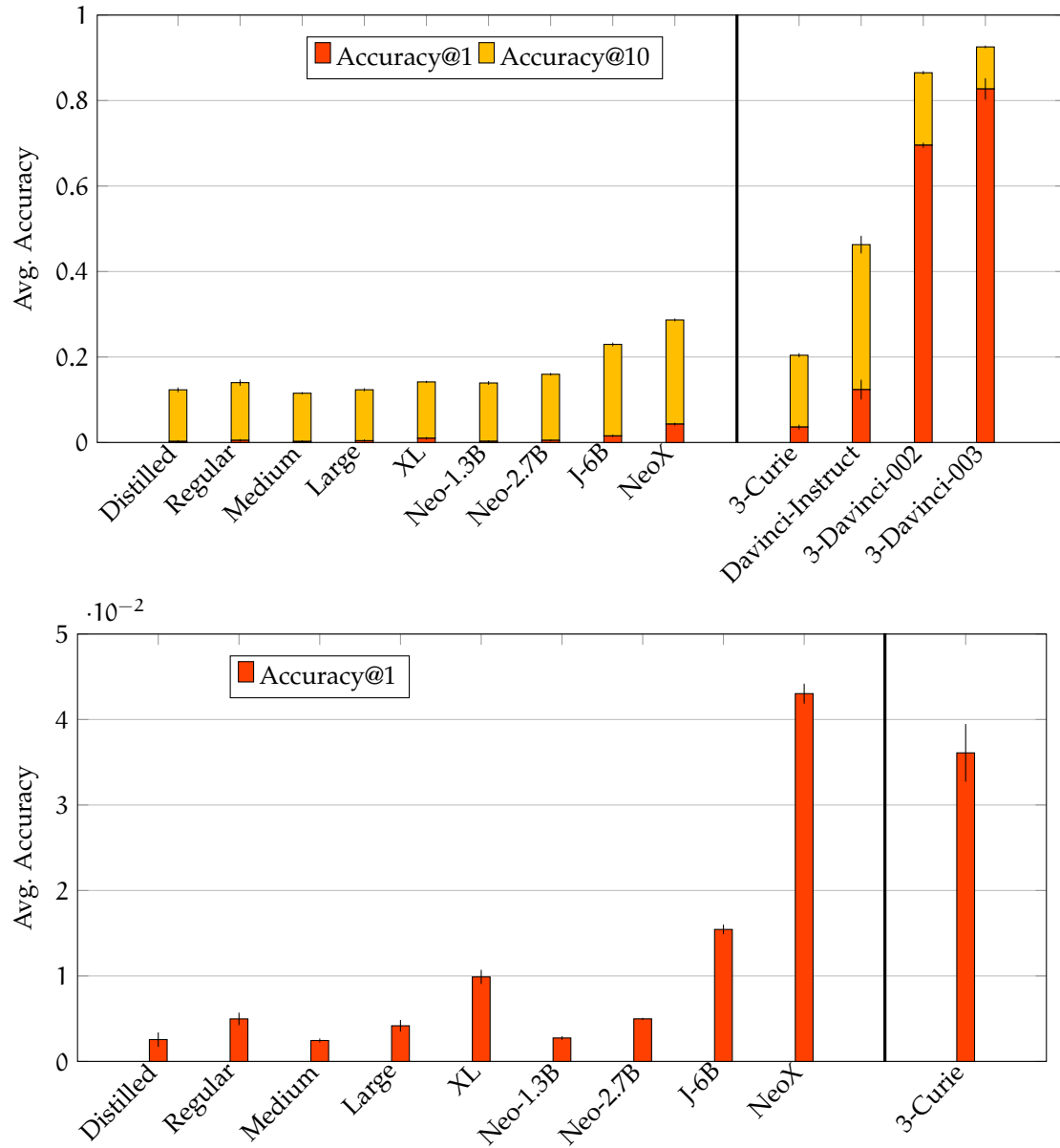


Figure A.11: Average accuracy of the models on the generated instances of MWPs. Results are averaged over two sets consisting of 500 problem instances generated for each template. The lower figure shows a zoomed-in visualization of the accuracy at 1.

for each token generated. To overcome this limitation, we approximate the relative probability change δ_{rcc} as follows, depending on the kind of effect measured.

The limit for k is set by OpenAI to 5. However, for our main set of experiments (i.e., computing the causal effects of N , S , and T) we were granted an increased limit of k to 100. This allowed us to obtain reasonable estimates for the causal effects, as the number of cases in which $P(g)$ is not defined is less than 10% of the number of examples that we consider.

Algorithm 2 Computation of δ_{rcc} for GPT-3

```

1:  $Q = (t, n, g)$ 
2:  $Q' = (t', n', g')$ 
3: if  $P(g)$  is defined then
4:   if  $P'(g)$  is defined then
5:      $\Delta = \frac{P(g) - P'(g)}{P'(g)}$ 
6:   else
7:      $\hat{p}' \leftarrow P'(k\text{-th most likely token})$ 
8:      $\Delta = \frac{P(g) - \hat{p}'}{\hat{p}'}$ 
9:   end if
10: else
11:    $\Delta = 0$ 
12: end if
13: if  $P'(g')$  is defined then
14:   if  $P(g')$  is defined then
15:      $\Delta' = \frac{P'(g') - P(g')}{P(g')}$ 
16:   else
17:      $\hat{p} \leftarrow P(k\text{-th most likely token})$ 
18:      $\Delta' = \frac{P'(g') - \hat{p}}{\hat{p}}$ 
19:   end if
20: else
21:    $\Delta' = 0$ 
22: end if
23:  $\delta_{\text{rcc}} = \frac{1}{2}(\Delta + \Delta')$ 

```

A.4.3.1 TCE(N on R) and TCE(T on R)

In cases when $P(g)$ is defined (i.e. when g appears in the top k token predictions) and $P'(g)$ is not defined, we compute a lower bound on the relative change using the upper bound on $P'(g)$ given by the probability of the k -th most likely token. This gives us a conservative estimate of Δ . For cases in which $P(g)$ is not defined, we cannot say anything about the relative change, and we set $\Delta = 0$. The same applies when swapping P and P' . This procedure is illustrated by Algorithm 2.

A.4.3.2 DCE(N \rightarrow R) and DCE(S \rightarrow R)

In this case, we simply discard the examples for which $P(g)$ is not defined or $P'(g)$ are not defined. In that is not the case, then we compute δ_{rcc} as in Section 5.3.4.

A.4.3.3 Heatmap Illustration

The heatmap for GPT-3 displayed in Figure 5.4 was computed by taking the raw probability score produced by the model over the whole vocabulary, as the limit on the available top predicted tokens makes it impossible to normalize it over the set $\{0, \dots, 300\}$, as done for the other models. The probability was set to 0 when g did not appear in the model's top 5 predictions for the next token after the prompt.

A.4.4 Computing Infrastructure and Inference Details

To run our experiments, we used a single NVIDIA TITAN RTX with 24GB of memory for all the versions of GPT-2 and GPT-Neo. We used a single NVIDIA A100 with 40GB of memory for GPT-J-6B and a single NVIDIA A100 with 80GB of memory for GPT-NeoX and the LLaMA models (two for the 30B version). We accessed GPT-3 using the OpenAI APIs. The longest run (GPT-J) on the four kinds of experiments corresponding to the four kinds of effects measured took ~12 hours, using 500 MWP instances for each of the 437 templates. Due to budget and resource constraints, the experiments on GPT-3, GPT-NeoX, and LLaMA were carried out using 20 examples generated for each template and took ~7 hours. Experiment tracking was carried out using Weights & Biases².

²<http://wandb.ai/>

A.5 Additional Materials for Chapter 6

A.5.1 Meta Study Settings of SSL and DA

For the meta study of SSL, we covered but are not limited to all relevant papers cited by the review on NLP SSL by Søgaard (2013). We went through the leaderboard of many NLP tasks and covered the SSL papers listed on the leaderboards. The papers covered by our meta study are available on our GitHub.

For supervised DA, we searched papers with the keyword domain adaptation and task names from a wide range of tasks that use supervised DA.

Note that for fair comparison, we do not consider papers without a comparable supervised baseline corresponding to the SSL, or a comparable unadapted baseline corresponding to the DA. We do not consider MT DA which tackles the out-of-vocabulary (OOV) problem because $P(E|C)$ may be different for OOV (Habash, 2008; III and Jagarlamudi, 2011).

A.5.2 Experimental Details of Minimum Description Length

We calculate the $MDL(X)$ and $MDL(Y)$ by a language model, and obtain $MDL(X|Y)$ and $MDL(Y|X)$ using translation models. For language model, we use the autoregressive GPT2 (Radford et al., 2019), and for the translation model, we the Marian Neural Machine Translation model (Junczys-Dowmunt et al., 2018) trained on the OPUS Corpus (Tiedemann and Nygaard, 2004). Both these models use the layers from the transformer model (Vaswani et al., 2017). The autoregressive language model consists only of decoder layers, whereas the translation model used six encoder and six decoder layers. Both of these models have roughly the same number of parameters. We used the huggingface implementation (Wolf et al., 2020) of these models for their respective set of languages.

A.6 Additional Materials for Chapter 7

A.6.1 Implementation Details

A.6.1.1 Model Details

Using Closed-Weight Models For the use of GPT model series, we use the OpenAI API,³ with a text generation temperature of 0. We spent around 400 USD across around 20-30K single API calls.

Using Open-Weight Models For reproducibility, we set the generation temperature to 0 for all the models used in our work. For the open-weight models, GPT2-XL, LLaMa-7B and Alpaca-7B, it took around 24 hours on 4 GPUs RTX 2080 to generate their predictions on 1K data points for the 5 paraphrases of the causally-neutral prompt (denoted as C0), and on 500 data points for the 5 paraphrases of the C1 prompt, and 5 paraphrases of the C2 prompt. The causal tracing experiments with LLaMa-7B and Alpaca-7B on 100 data points took around 24 hours each using one GPU V100.

A.6.1.2 Implementation Details for Causal Tracing

We introduce the workings of the causal tracing method (Meng et al., 2022) as follows. First, we compute the hidden states of the residual stream of LLaMa-7B’s layers for two inputs, (1) the original input: the prompt+review, and (2) the corrupted input: prompt + a corrupted version of the review by adding random noise immediately after the token embeddings. Then, we restore one by one the clean state of the residual stream into the corrupted version and measure the effect of the clean state on the probability of the originally predicted token for each token sequence and layer position.

Since this process is highly time-consuming, taking around 12 hours for 50 samples even using the smallest LLaMa model with 7B parameters, we do a case study on the 7B LLaMa and Alpaca using 100 random samples from the 1K test set. For these experiments, we follow the idea of APE (Zhou et al., 2023) to use the best-performing prompts on the 1k test set for C1 and C2.

A.6.1.3 Prompts

A.6.1.3.1 Prompts to Get paraphrases

Since we need to report the average performance across five paraphrases of the same prompt, for each original prompt, we call GPT to generate the four paraphrases.

Below is the prompt that we used for this paraphrase generation process:

You are an expert in prompt engineering for large language models (LLMs). And you are also a native English speaker who writes fluent and grammatically correct text.

Given the following prompt for NLP sentiment analysis, you provide four alternative prompts.

³<https://openai.com/api/>

Original Prompt ##### [Our original prompt]

Alternative Prompt 1

... (Then, we let the model to generate all the way to “Alternative Prompt 4”.)

We queried the GPT-4 model with temperature 0 on June 8, 2023.

A.6.1.3.2 Neutral Prompt

In addition to the standard prompt to query LLMs in the main paper, we show its four paraphrases in Table A.15.

Prompt Design
As a proficient data annotator in natural language processing (NLP), your responsibility is to determine the sentiment of the given review text. Please assign a sentiment value from “1” (very negative) to “5” (very positive). Review Text: “[review]” Sentiment Score:
As a skilled data annotator in the field of natural language processing (NLP), your task is to evaluate the sentiment of the given review text. Please classify the sentiment using a scale from “1” (highly negative) to “5” (highly positive). Review Text: “[review]” Sentiment Rating:
As an expert data annotator for NLP tasks, you are required to assess the sentiment of the provided review text. Kindly rate the sentiment on a scale of “1” (extremely negative) to “5” (extremely positive). Review Text: “[review]” Sentiment Score:
As a proficient data annotator in natural language processing (NLP), your responsibility is to determine the sentiment of the given review text. Please assign a sentiment value from “1” (very negative) to “5” (very positive). Review Text: “[review]” Sentiment Assesment:

Table A.15: Four additional paraphrases of the neutral prompt (C0) generated with GPT-4.

A.6.1.3.3 Causal Prompts

In addition to the standard C1 and C2 prompts in the main paper, we show the four paraphrases for each of them in Tables A.16 and A.17, respectively.

A.6.2 Additional Experimental Results

A.6.2.1 Few-Shot Results

For reproducibility and controllability, we use the zero-shot prompting setting across the experiments in the main paper, to avoid randomness in few-shot prompting according to which examples are selected as the few shots, and the order of the examples.

As a supplementary information in case this is of some readers’ interest, we provide the few-shot prompting results in Tables A.18 and A.19.

Prompt Design
As a customer sharing my experience, I crafted the following review: “[review]” Taking into account the details of my experience, I chose a star rating from the available options of “1”, “2”, “3”, “4”, or “5”. My ultimate rating is:
As a client providing my opinion, I penned down the subsequent evaluation: “[review]” Upon thorough reflection of my encounter, I picked a star rating among the choices of “1”, “2”, “3”, “4”, or “5”. My conclusive rating stands at:
As a patron expressing my thoughts, I drafted the ensuing commentary: “[review]” After meticulously assessing my experience, I opted for a star rating from the range of “1”, “2”, “3”, “4”, or “5”. My definitive rating turned out to be:
As a consumer conveying my perspective, I authored the following assessment: “[review]” By carefully weighing the aspects of my interaction, I determined a star rating from the possibilities of “1”, “2”, “3”, “4”, or “5”. My final verdict on the rating is:

Table A.16: Four additional paraphrases of the causal prompt C1 generated with GPT-4.

A.6.2.2 λ_1 - λ_2 Distribution Plot

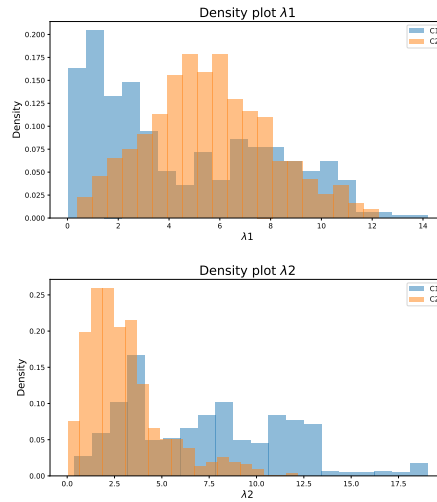


Figure A.12: The λ_1 - λ_2 density plots of C1 (above) and C2 (below).

To provide a clear understanding of the distributions of λ_1 and λ_2 , we include their density plots of the causal processes C1 and C2 in Figure A.12. The mean values of λ_1 and λ_2 for each group are in Table A.20.

Further, we performed the Mann-Whitney U rank test to determine if the underlying distributions of λ_1 and λ_2 for groups C1 and C2 are the same. The results are as follows:

- For λ_1 , the p-value is 8.4572×10^{-71} , leading us to reject the null hypothesis that the two groups come from the same distribution.
- For λ_2 , the p-value is 1.36138×10^{-11} , also leading us to reject the null hypothesis that the distributions are the same.

These statistical results indicate significant differences between the distributions of λ_1 and

Prompt Design
As a customer sharing my experience, I first chose a star rating from the available choices of "1", "2", "3", "4", or "5", and subsequently elaborated on my decision with the following statement: "[review]"
The review elucidates the reasoning behind my assigned rating of
As a client providing my opinion, I initially picked a star rating from the range of "1" to "5", and then proceeded to justify my selection with the following commentary: "[review]"
The review sheds light on the rationale for my given rating of
As a patron expressing my thoughts, I started by selecting a star rating from the scale of "1" to "5", and then offered an explanation for my choice in the following review text: "[review]"
The review expounds on the basis for my designated rating of
As a consumer conveying my perspective, I began by opting for a star rating within the "1" to "5" spectrum, and then detailed my reasoning in the subsequent review passage: "[review]"
The review delineates the grounds for my conferred rating of

Table A.17: Four additional paraphrases of the causal prompt C2 generated with GPT-4.

		Random	GPT-3 Few-Shot
F1	Overall	19.82 ± 2.07	63.35 ± 0.80
	C1 Subset	21.36 ± 2.26	54.44 ± 1.24
	C2 Subset	20.43 ± 2.95	75.65 ± 0.45
Acc	Overall	19.78 ± 2.07	64.14 ± 0.86
	C1 Subset	20.61 ± 2.23	54.22 ± 1.28
	C2 Subset	18.86 ± 2.78	75.18 ± 0.53

Table A.18: Few-shot performance of the standard SA prompts on Yelp-5. We use five paraphrases for the prompt, and report the average performance with the standard deviation.

λ_2 across the causal process groups, which indicate distinct underlying characteristics in the sentiment dynamics of the two groups.

A.6.3 Emotion Arc Clustering

We analyze the emotional arc patterns of Yelp reviews. Reagan et al. (2016) identified 6 basic emotional arc shapes in stories. However, reviews are usually shorter and therefore present fewer variations. We take each sentence of the review and predict its sentiment. Then we divide the review into ten bins and compute an average sentence sentiment for each decile to make reviews with different lengths comparable. Reviews shorter than 10 sentences generate null values for some deciles which we fill with the information of the next decile. In Figure A.14, we illustrate the 4 clusters found, they have the following characteristics:

Positive + Early Rise: This cluster primarily comprises highly positive reviews, where customers express satisfaction and praise for their overall experience. Interestingly, 21.1% of these reviews begin with a negative first sentence, which often indicates initially low expectations or a negative first impression. However, despite the initial negativity, the reviews tend to turn positive as customers elaborate on their positive experiences.

		Random	GPT-3 Few-Shot
F1	Data=C1, Prompt=C1	20.47 ± 2.47	49.18 ± 0.76
	Data=C1, Prompt=C2	20.26 ± 2.31	52.79 ± 2.64
	Data=C2, Prompt=C1	22.35 ± 3.02	80.46 ± 1.29
	Data=C2, Prompt=C2	20.35 ± 2.18	75.88 ± 1.86
Acc	Data=C1, Prompt=C1	19.60 ± 2.67	50.12 ± 0.71
	Data=C1, Prompt=C2	19.68 ± 2.46	53.83 ± 2.46
	Data=C2, Prompt=C1	20.97 ± 3.19	81.21 ± 1.17
	Data=C2, Prompt=C2	19.03 ± 2.18	76.36 ± 1.53

Table A.19: Few-shot performance on Yelp using two different causal prompts on the two causal subsets. We use five paraphrases for each prompt, and report the mean performance with the standard deviation.

	C1	C2
$\mu(\lambda_1)$	4.48	5.62
$\mu(\lambda_2)$	7.31	3.02

Table A.20: Mean values of the lambdas for C1 and C2.

Negative + Early Fall: This cluster mainly consists of predominantly negative reviews. Similarly to the Positive cluster, some reviews (28.14%) start with a sentence with the opposite sentiment, usually indicating high expectations followed by disappointment.

Rise: The main characteristic of this cluster is the positive ending of the review, despite the initial negativity observed in the first half, with an average sentiment of -1.63. An important fraction of the reviews in this cluster (52.49%) start with a positive comment as a summary, but then proceed to highlight the negative aspects of the experience. Despite the initial criticisms, the reviews conclude with positive points, suggesting that the overall experience was still satisfactory.

Fall: In contrast to the previous cluster, the Fall cluster is characterized by a negative ending of the review, despite a generally positive first half with an average sentiment of 2.18. An important proportion (36%) of the reviews in this cluster begin with a negative comment as a summary, but then proceed to describe the positive aspects before eventually highlighting the negative ones. This cluster showcases a shift in sentiment from positive to negative, indicating a decline in satisfaction as the review progresses.

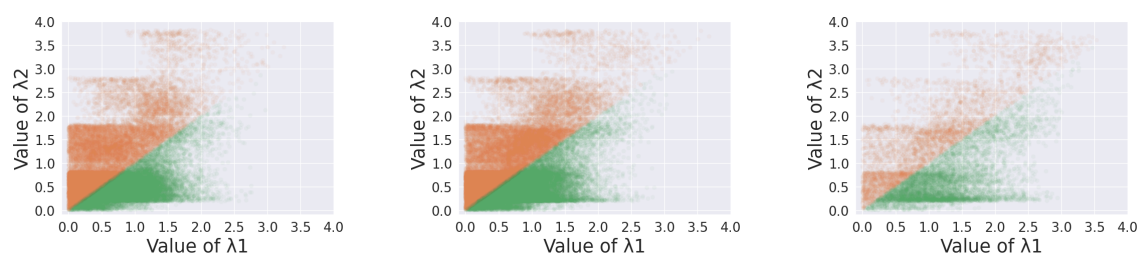


Figure A.13: The λ_1 - λ_2 plot on Yelp-5 (left), Amazon (middle), and App Review (right). We draw the $y = x$ diagonal line, and the orange dots in the upper-left triangle represent the C1-dominant subset, and green dots in the lower-right triangle are the C2-dominant subset.

Positive + Early Rise

Review: Was there last Friday. Seats right in front of the stage. The show was good. The headliner, while a bit long, was good. Fantastic service from our waitresses. Will definitely go back.

Review: This is by far my favorite Panera location in the Pittsburgh area. Friendly, plenty of room to sit, and good quality food & coffee. Panera is a great place to hang out and read the news - they even have free WiFi! Try their toasted sandwiches, especially the chicken bacon dijon.

Negative + Early Fall

Review: Pass on this place, there are better restaurants mere feet away.

The menu here is too large, which is a sure sign none of the food is going to be good. And, its not good. Some of the salads are alright, but its just not good food.

The service is friendly and prompt, but the beer is over priced. They do have a good selection though. This place is open late if you need a bite to eat, but there are so much better options out there.

Review: Wings are overpriced. And the quality of them are bad. They were tough and greasy. The staff are pleasant but then over all experience was too expensive for a sports bar.

Rise

Review: To be honest, I feel that this is one of the most overpriced restaurants in the entire city. The food is average to good, the place is beautiful with outdoor seating, but in my opinion the price is just not worth it. They have a really good happy hour, so I would definitely recommend going to that and maybe trying an appetizer or two.

Review: The first time I came here, I waited in line for 20 minutes. When it was my turn, I realized I left my wallet in the car. It hurt so bad, I didn't come back for a year.

I can walk to this place from my house- which is dangerous because those biscuits are just OH SO DREAMY. I can't describe them. Just get some.

Do I feel guilty about noshing on fabulous Strawberry Napoleons and Jewish Pizza (kind of like a modified, yet TOTALLY delicious fruitcake bar) at 10:15am? Hecks, naw... But they do have quiche and some other breakfast-y items for those who prefer a more traditional approach to your stomach's opening ceremony.

Just go early :) They open at 10 on Saturdays. And bring cash...it's easier that way.

Fall

Review: It's cheap, I'll say that, but otherwise it's bland food served by workers who mostly don't seem to notice they're working, and when they do, only respond snarkily. There are many better vegetarian and vegan options to choose from

Review: I do like my Mad Mex, however predictable and non-authentic it may be. The portion sizes are mammoth and I come away with a satisfied sense of regret. Their beer menu is happily extensive. Charging me \$9 for chips and salsa is a bit of crime, wouldn't ya say though!?! I mean, c'mon! Our service has most times been lacking—a bit rushed and on the inattentive side. Also, why do you require your wait staff to not servestraws/lemons/etc unless asked by cusotmers—weirdness-cut out these odd cost-cutting, anti-service friendly measures please

Table A.21: Example reviews for each emotion arc cluster.

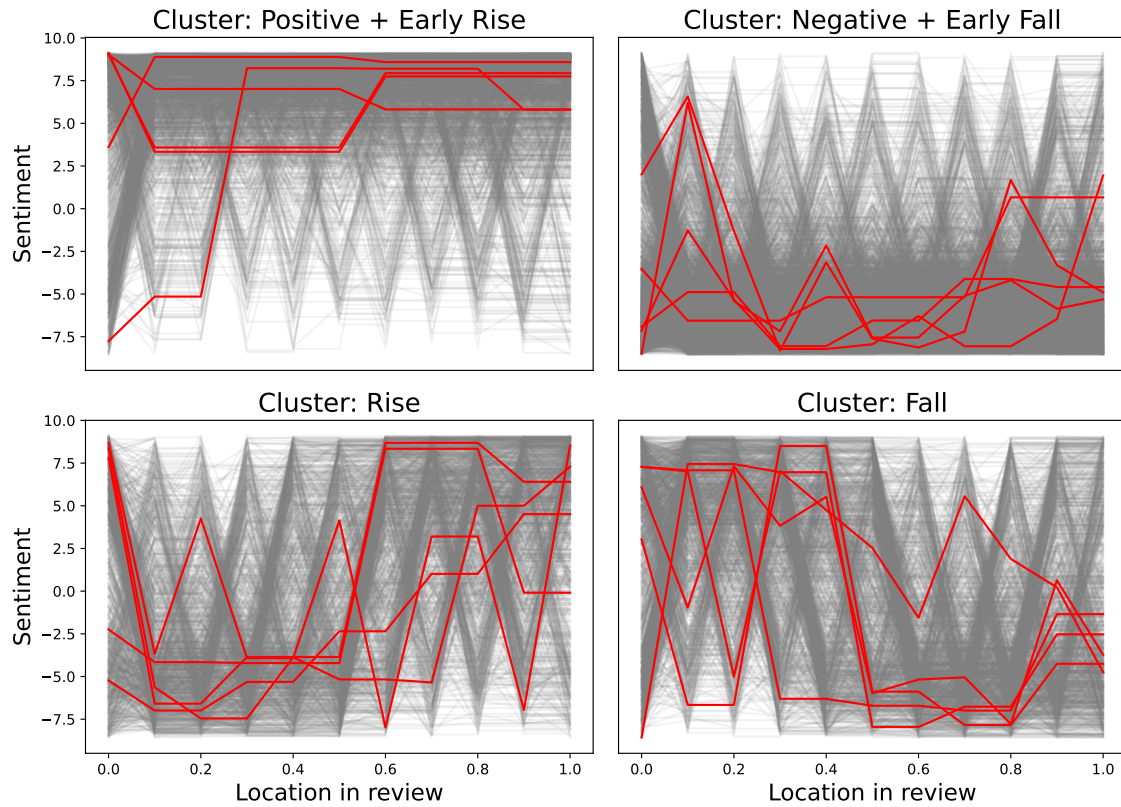


Figure A.14: Four emotion arc clusters.

A.6.4 Additional Interpretability by Shapley Values

We further analyze the effect of each part of the prompts on LLaMa’s predictions. Using 50 reviews, we compute the shapley values of each token. In Figure A.15 we observe that the tokens with the largest shapley values are the ones in the end, which is expected since they are the ones helping to form a grammatically correct sentence. To account for that, we subtracted the average shapley values computed for the other possible start rating answers. In Figure A.16 we show the adjusted shapley values. We observe that the tokens in prompt C1 have a larger effect than the tokens in prompt C2. The words introducing the review have a positive effect on C2 but a negative one on C1. Whereas, the phrase “I chose a star rating” has a negative effect on C2 but a positive one on C1.

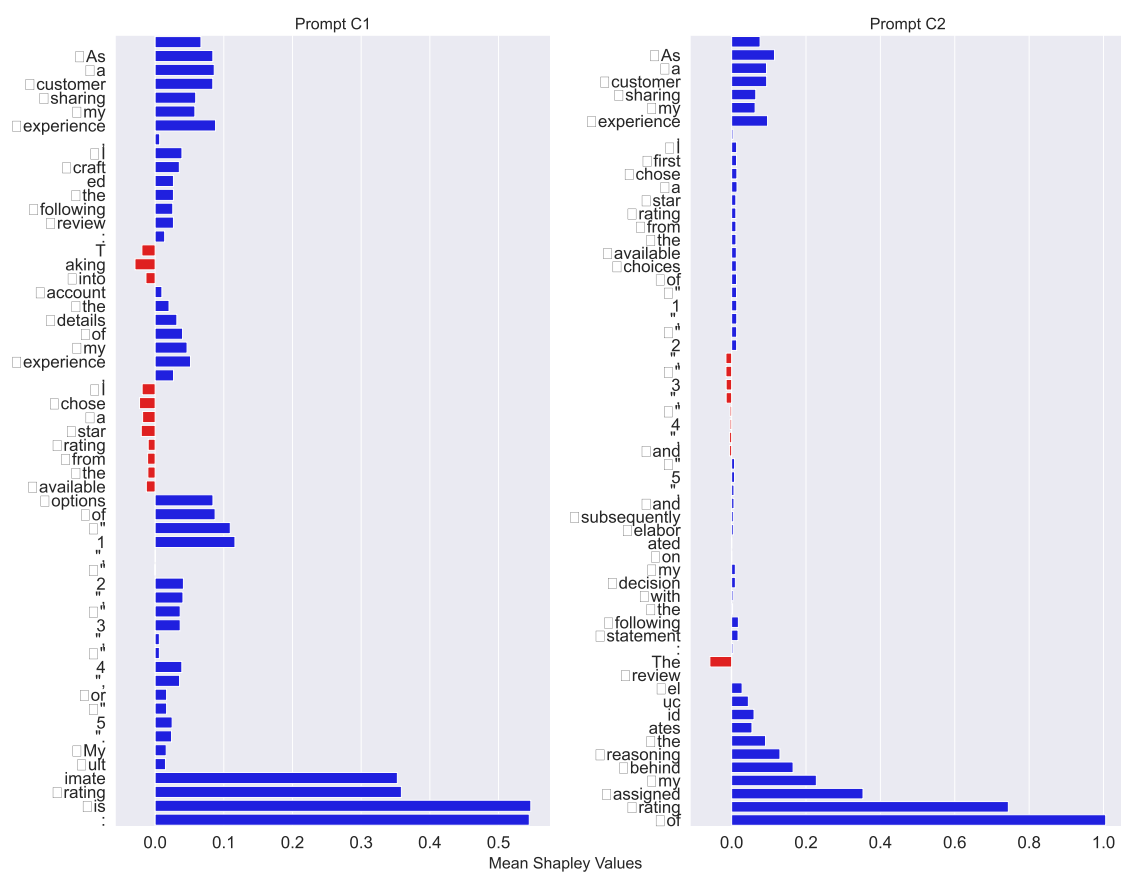


Figure A.15: Shapley values for the two types of prompts.

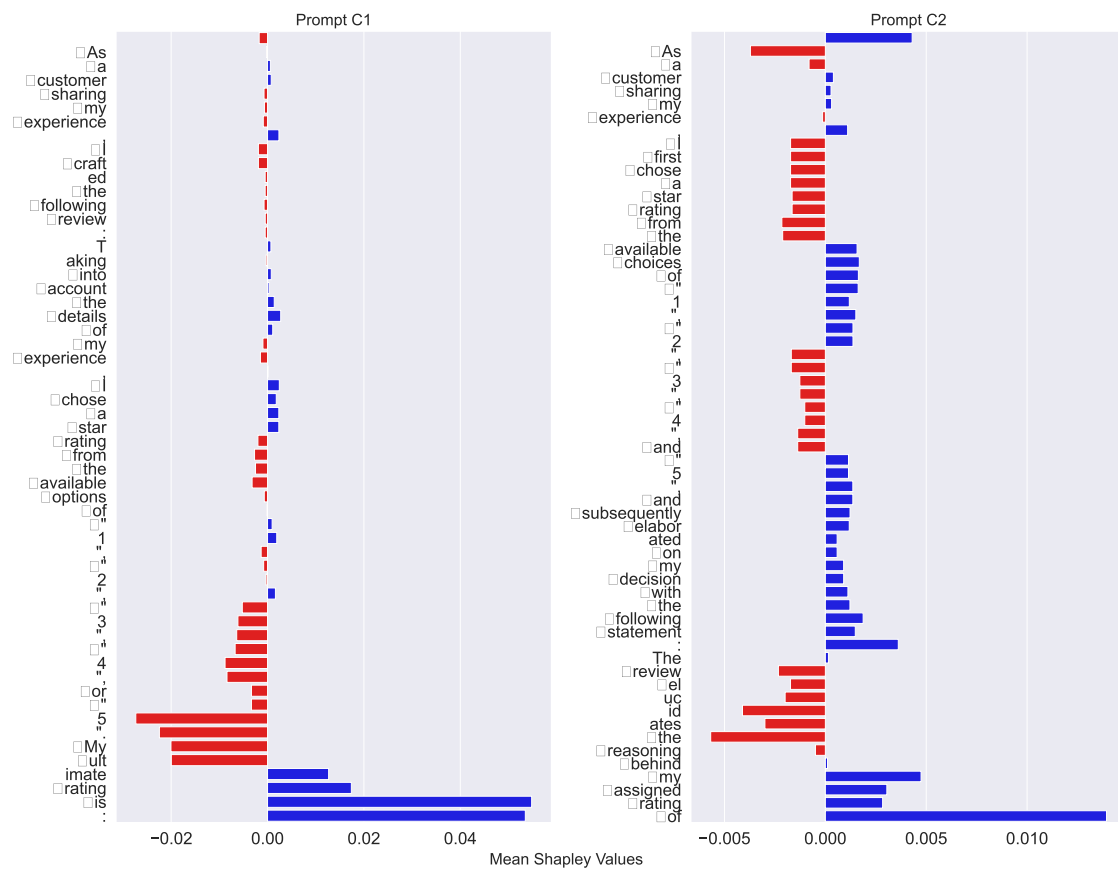


Figure A.16: Adjusted shapley values for the two types of prompts.

A.7 Additional Materials for Chapter 8

A.7.1 Statistics of our Data

A.7.1.1 COVID Twitter Keywords

We list the COVID-related Twitter keywords and accounts tracked by Chen et al. (2020) in Tables A.22 and A.23. They are used to retrieve the 1.01TB raw Twitter data.

Keywords used by Chen et al. (2020)	
14DayQuarantine	covidiot
CDC	epitwitter
COVID	flatten the curve
COVID__19	flattenthecurve
COVID-19	kung flu
China	lock down
Corona	lockdown
Coronavirus	outbreak
Coronials	pandemic
DontBeASpreader	pandemie
DuringMy14DayQuarantine	panic buy
Epidemic	panic buying
GetMePPE	panic shop
InMyQuarantineSurvivalKit	panic shopping
Koronavirus	panic-buy
Kungflu	panic-shop
N95	panicbuy
Ncov	panicbuying
PPEshortage	panicshop
Sinophobia	quarantinelifelife
Social Distancing	quarentinelife
SocialDistancing	saferathome
SocialDistancingNow	sars-cov-2
Wuhan	sflockdown
Wuhancoronavirus	sheltering in place
Wuhanlockdown	shelteringinplace
canceleverything	stay at home
china virus	stay home
chinavirus	stay home challenge
chinese virus	stay safe stay home
chinesevirus	stayathome
corona virus	stayhome
coronakindness	stayhomechallenge
coronapocalypse	staysafestayhome
covid	trump pandemic
covid-19	trumppandemic
covid19	wear a mask
covidiot	wearamask

Table A.22: Keywords used by Chen et al. (2020) to track COVID-related tweets.

Accounts tracked by Chen et al. (2020)	
PneumoniaWuhan	WHO
CoronaVirusInfo	HHSGov
V2019N	NIAIDNews
CDCemergency	DrTedros
CDCgov	

Table A.23: Accounts tracked by Chen et al. (2020) to retrieve COVID-related tweets.

A.7.1.2 Annotation Guidance for Policy Strictness

For each state, the annotators are asked to go to the official website that lists all COVID policies of the state. In most cases, the website lists all executive orders (EOs), proclamations, or other forms of policies issued during 2020 – 2021. Then the annotator is asked to read through the EOs that are related to COVID social distancing policies. For each relevant policy, the annotator is asked to record the start date on which the policy will take effect,⁴ a brief intro of what kind of social distancing policy it is, and a real-valued score in the range of 0 (loosest) to 5 (strictest).

For the scoring criteria, we provide the following guides:

- Score 0: masks are optional, open the schools,, bars, gaming facilities, concert, and almost everything
- Score 1: State of emergency, limit gathering, close K-12
- Score 2: Open 50% capacity for retail business, open religious activities like churches to 50%
- Score 3: Open 25% capacity for retail businesses
- Score 4: Open only business for necessities such as supermarkets, only allow delivery and curbside services, gatherings have to be no more than 10 people
- Score 5: Strict stay at home policy, close every business

A.7.1.3 Accuracy of Twitter Sentiment Classifier

We list the detailed performance report of TextBlob and our COVID BERT in Table A.24, including the overall accuracy, weighted and macro F1 scores, precision and recall for each class, and MSE of the average sentiment of random groups of 20 tweets. Note that since TextBlob predicts a real-valued number in the range of -1 to 1 for the sentiment, we regard [-1, -0.33) as negative, [-0.33, 0.33] as neutral, and (0.33, 1] as positive.

Model	Acc	F1 Score		Positive		Neutral		Negative		MSE on Groups
		Weighted	Macro	P	R	P	R	P	R	
TextBlob	23.35	16.67	19.70	20.34	10.62	20.67	85.19	74.07	6.45	0.43
COVID BERT	60.23	62.31	55.17	51.19	76.11	26.76	35.51	83.68	62.99	0.15

Table A.24: The detailed performance report of the TextBlob baseline, and our COVID BERT model. We report the overall accuracy (Acc), weighted and macro F1 scores, precision (P) and recall (R) for each class, and MSE of the average sentiment of random groups of 20 tweets.

⁴For consistency, we record 0:01am of the first effective date, but not the 11:59pm of the previous day.

A.7.2 Additional Analyses

A.7.2.1 Correlation across All Variables

We can see that, averaging over all 50 states, unemployment correlates the most with policy changes, which is consistent with our analysis in Section 8.5.1. Since different states may have different styles to take sentiment into consideration when making policies, the effect of sentiment on policy changes over all 50 states is relatively mild.

For Twitter sentiment, it correlates largely with case numbers, and urbanization rate of the state.

Interestingly, the case numbers correlate with whether the state governor is a political ally of Trump.

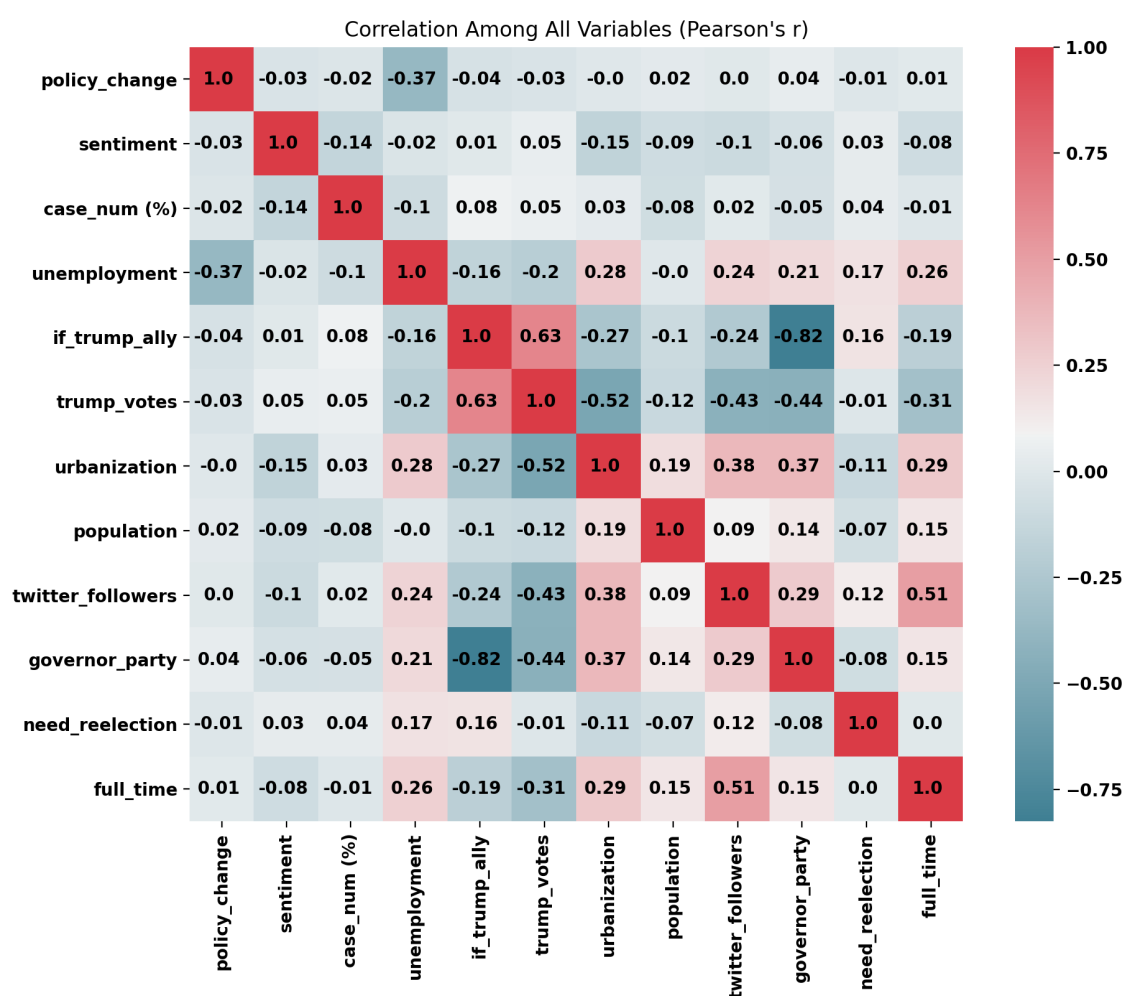


Figure A.17: Correlation across all variables.

A.7.2.2 Alternative Causal Analysis Methods by Potential Outcomes Framework

There are two commonly used frameworks for causal inference, one is the do-calculus we introduced in Section 8.5.2, and the other is the potential outcomes framework (Rubin,

1974, 2005; Imbens and Rubin, 2015). We will introduce two alternative causal inference methods on our problem, using the potential outcomes framework.

Difference-in-Differences One possible limitation of this study is that we treat the data in an i.i.d. way, following most existing studies. An improvement is to treat it as time series. For time series analyses, one commonly used method is the first-difference (FD) estimator, difference in differences (DID) (Abadie, 2005). Specifically, DID takes in the time series data of the cause X , effect Y , and confounders Z , and solves the following regression:

$$\Delta Y = \beta \cdot \Delta X + \Delta Z \quad (\text{A.6})$$

$$Y_t - Y_{t-1} = \beta(X_t - X_{t-1}) + Z_t - Z_{t-1}, \quad (\text{A.7})$$

where t is the time step, and β is the causal effect of X on Y .

After applying DID on all the policies, we obtain β scores for all states, and the top 5 states with largest β are Colorado ($\beta = 0.67$), Kentucky ($\beta = 0.23$), Wyoming ($\beta = 0.22$), Oregon ($\beta = 0.19$), North Carolina ($\beta = 0.17$), Michigan ($\beta = 0.14$), and New York ($\beta = 0.13$).

Continuous-Valued Propensity Score Matching Another commonly used alternative for causal inference is propensity score matching. However, the challenge in our study is that the cause is not categorical, but takes continuous values. To this end, we follow the extension of propensity score matching to continuous treatment (Hirano and Imbens, 2004; Bia and Mattei, 2008). We adopt the stata package of Bia and Mattei (2008) for continuous-valued propensity score matching. The resulting prediction of policies based on Twitter sentiment is a polynomial function with an order of three. As examples, We show the predictions of Texas (TX) and Michigan (MI) in Figure A.18.

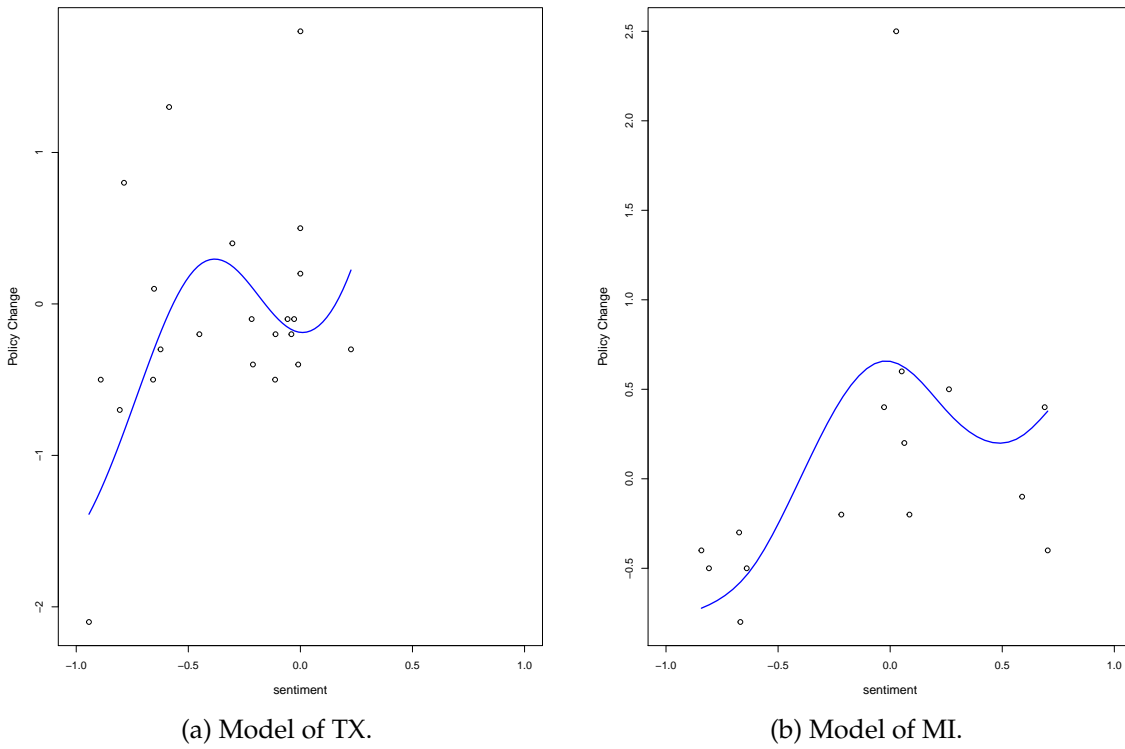


Figure A.18: Causal models by continuous-valued propensity score matching of TX and MI.

A.8 Additional Materials for Chapter 9

A.8.1 Additional Implementation Details

A.8.1.1 Time and Space Complexity Details

For the time cost of running the causal impact indices, each $\text{PCI}(a, b)$ takes around 1,500 seconds, or 25 minutes. Multiplying this by 40 samples per paper a , we spend 16.67 hours to calculate each ACI or TCI for the paper's overall impact. For a fine-grained division into the time cost, the majority of the time is spend on the BM25 indexing (800s) and the sentence embedding cosine similarities calculation (400s). The rest of the time-consuming steps are the BFS search (150-200s every time) to identify descendants and non-descendants of a paper.

For the space complexity, we loaded the 2.4B edges of the citation graph into a parquet gzip format for faster loading, and use Dask's lazy load operation to load it part by part to RAM for better parallelization. The program can fit into different sizes of RAMs by modifying the number of partitions and reducing the number of workers in Dask, at the cost of an increased computation time. On the hard disk, citation graph takes up 19G space, and paper data takes 11G.

A.8.1.2 Numerical Estimation Method: Finding the Sample Size

For our numerical estimation method, we first calculate the ACI on a subset of carefully sampled papers and then aggregate it to TCI. One design choice question is how to decide

the size of this random subset. In our case, we need to balance both the computation time (25 minutes per pairwise paper impact) and the estimation accuracy. To identify the best sample size, we conduct a small-scale study, first obtaining the TCI using our upper-bound budget of $n = 100$ samples and then gradually decreasing the number of samples to see if there is a stable point in the middle which also leads to a result close to that obtained with 100 samples. In Figure A.19, we show the trade-off of the two curves, the error curve and time cost, where we can see $n = 40$ seems to be a good point balancing the two. It is at the elbow of the arrow curve, making it relatively close to the estimation result of $n = 100$, and also in the meantime vastly saving our computational budget, enabling us to run efficient experiments for more analyses.

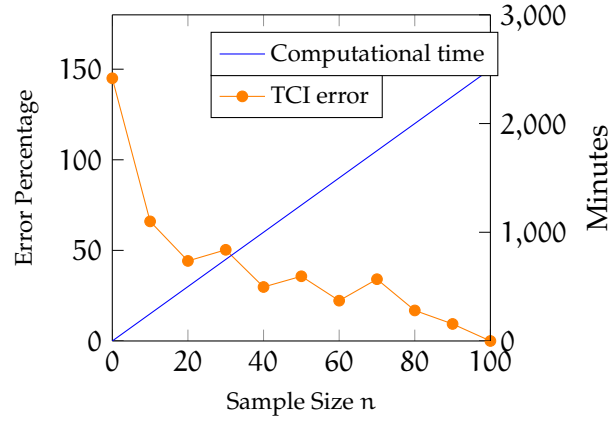


Figure A.19: We show the trade-off of two curves: the error curve (orange), and the time cost curve (blue). For the error curve, we see an elbow point at around $n = 40$, when the error starts to be small. The curve for the computational time is linear, taking 25 minutes for each paper. Balancing the trade-offs, we decided to choose the sample size $n = 40$.

A.8.1.3 Experiment to Select the Best Embedding Method

When selecting the text encoder for our CAUSALCoT method, we compare among the three LLMs pre-trained on scientific papers, SciBERT, MPNet, and SPECTER. Specifically, we conduct a small-scale experiment to see how much the similarities scores based on the embedding of each model align with human annotations. As for the annotation process, we first collect a set of random papers, and for each such paper (which we call a pivot paper), we identify ten papers, from the most similar to the least, with monotonically decreasing similarity. We collect a total of 100 papers consisting of ten such collections, for which we show an example in Table A.25. Then we see how the resulting similarity scores conform to this order by deducting the percentage of papers that are out of place in the ranking.

We find that MPNet correlates the best with human judgments, achieving an accuracy of 82%, which is 10 points better the second best one, SPECTER, which gets 72%, and 18 points better than SciBERT with a score of 64%. It also gives more distinct scores to papers with different levels of similarity. This capability advantage may be attributed to its Siamese network objectives in the training process (Song et al., 2020). We open-sourced our annotated data in the codebase.

Paper Index	Title	SciBERT	SPECTER	MPNet
<i>Pivot Paper: GPT-3 (2020)</i>				
1 (Most similar)	PaLM (2022)	0.9787	0.8689	0.7679
2	GPT-2 (2019)	0.9346	0.9064	0.8196
3	GPT (2018)	0.9488	0.8778	0.7790
4	BERT (2019)	0.9430	0.8321	0.6784
5	Transformers (2017)	0.9202	0.8644	0.6385
6	SciBERT (2019)	0.8396	0.8112	0.5667
7	Latent Diffusion Models (2021)	0.9586	0.7755	0.4567
8	Sentiment Analysis Using DL (2015)	0.7775	0.7298	0.2911
9	Sentiment Analysis Using ML (2014)	0.6462	0.6403	0.2563
10 (Least similar)	New High Energy Accelerator (1952)	0.8033	0.5617	0.0359

Table A.25: An example collection of papers with monotonically decreasing similarity to the pivot paper. As can be seen from the similarities scores produced by the three text embedding methods, MPNet corresponds to the ground truth the most, and also shows clear score distinctions between less similar and more similar papers.

A.8.2 Dataset Overview

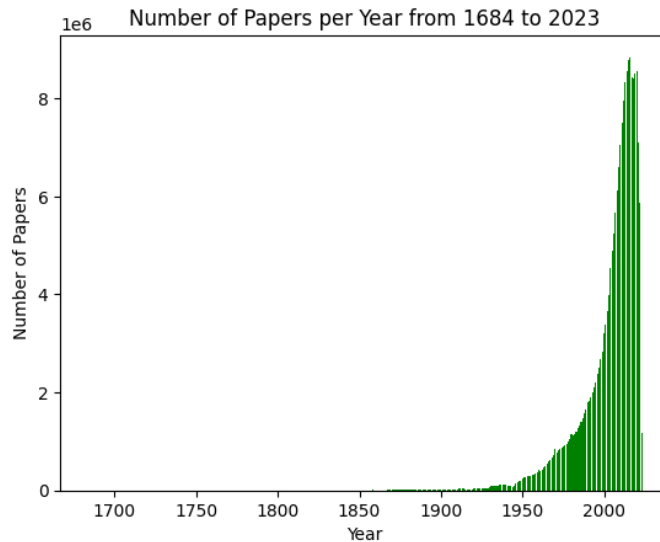


Figure A.20: The number of papers published per year from 1684 to 2023. We can see that in recent years since 2010, there are more than 7 million papers each year.

For the Semantic Scholar dataset (Kinney et al., 2023; Lo et al., 2020), we obtain the set of 206M papers using the “Papers” endpoint to get the Paper Id, Title, Abstract, Year, Citation Count, Influential Citation Count (Valenzuela et al., 2015), and the Reference Count for each paper. The papers come from a variety of fields such as law, computer science, linguistics, chemistry, material science, physics, geology, etc. For the citation network with 2.4B edges, we use the Semantic Scholar Citations API to get each edge of the citation graph in a triplet format of (fromPaper, toPaper, isInfluentialCitations).

In general, the number of publications shows an explosive increase in recent years. Figure A.20 shows the number of papers published per year, which reaches on average 7.5M per year since 2010. Figure A.21 shows the number of references each paper cites, which

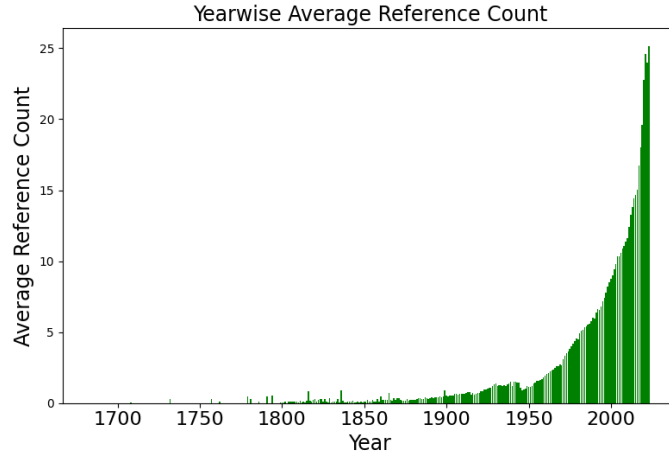


Figure A.21: The year-wise average of the number of references per paper, also with a sharply increasing trend.

also increases from less than five before 1970s, to around 25 in recent years. Both statistics support the need of our paper, which helps distinguish the quality of scientific studies given such massive growths of papers.

A.8.3 Additional Analyses

A.8.3.1 Citation Outlier Analysis

For the outlier detection, we first visualize the scatter plot between our CAUSALCITE and citations. Then, we fit a log-linear regression to learn the line $\log(\text{TCI}) = 1.026 \log(\text{Cit}) - 0.541$, as shown in Figure A.22, with a root mean squared error (RMSE) of 0.6807. After fitting the function, we use the interquartile range (IQR) method (Smiti, 2020), which identify as outliers any samples that are either lower than the first quartile by over 1.5 IQR, or higher than the third quartile by more than 1.5 IQR, where IQR is the difference between the first and third quartile.

We denote as overcited papers the ones that are identified as outliers by the IQR method due to too many citations than what it should have deserved given the CAUSALCITE value. Symmetrically, we denote as undercited papers the ones that are identified as outliers by the IQR method due to too few citations than what it should have deserved given the CAUSALCITE value. And we denote the non-outlier papers as the aligned ones.

A.8.3.2 Additional Information for the Author-Identified Paper Impact Experiment

As mentioned in the main paper, the dataset is annotated by pivoting on each paper b , and going through each of its references a to label whether a has a significant influence on b or not. We show an example of paper b and all its 31 references in Table A.26. We calculate the accuracy of each metric with the spirit that each non-significant paper's impact value should be lower than a significant paper's. Specifically, we go through the score of each non-significant paper, and count its accuracy as 100% if it is lower than all the significant papers', or the more general form $n_{\text{lower}}/|\text{Sig}|$ of conformity, where n_{lower} is the number of significant papers which it is lower than, and $|\text{Sig}|$ is the total number of significant papers. Then we report the overall accuracy for each score by averaging the accuracy numbers on

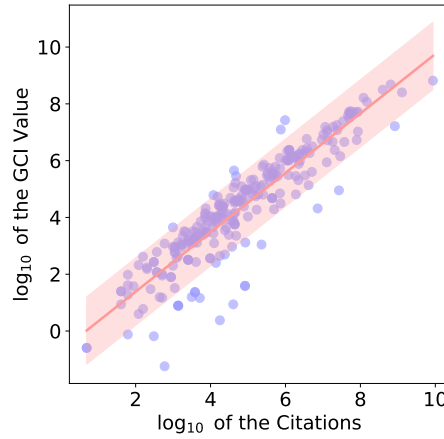


Figure A.22: The scatter plot between our CAUSALCITE and citations, with the fitted function as $\log(\text{TCI}) = 1.026 * \log(\text{Cit}) - 0.541$, and a non-outlier band width of 0.8809.

each non-significant paper. To illustrate the idea better, we show the calculated accuracy numbers for all three metrics on our example batch in Table A.26.

A.8.3.3 Step Curve for PCI Values Given a Fixed Paper b

Apart from the long-tailed curve shape of TCI in Section 9.5.2, we also look into the pairwise paper impacts by PCI. If we fix the paper b , we can see that $\text{PCI}(\cdot, b)$ often has a step curve shape in Figure A.23. The reason behind it lies in the nature of PCI, which is calculated based on the top K papers that are similar in content with paper b , but do not cite paper a . When we go through different references, e.g., from a_1 to a_2 of the same paper b , the semantically matched top K papers could still be largely the same pool, and only change when some papers in the pool need to be swapped when releasing the constraint to be that they can cite a_1 , and adding the constraint that they cannot cite a_2 .

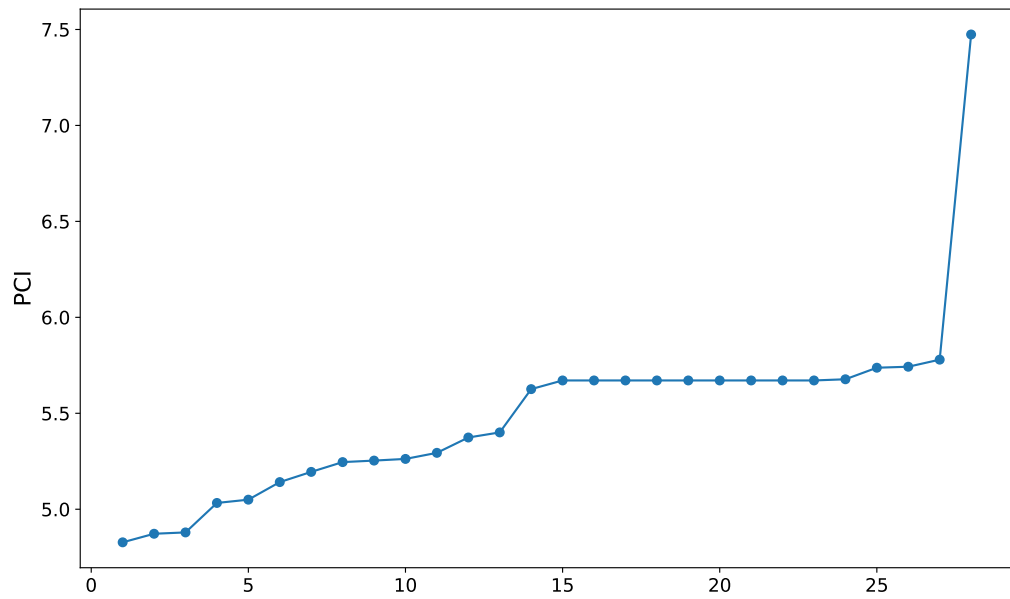


Figure A.23: We take an example paper b, Sentence BERT (Reimers and Gurevych, 2019), and plot its PCI values with all its reference paper a's. We can see clearly that there is a plateau in the curve, showing a step function-like nature.

<i>References of the Paper "Sorting improves word-aligned bitmap indexes"</i>	Label	PCI	Citations	SSHI
- A Quantitative Analysis and Performance Study for Similarity-Search Methods in High-Dimensional Spaces	0	3.519	1777	156
- Optimizing bitmap indices with efficient compression	0	3.519	375	40
- Data Warehouses And Olap: Concepts, Architectures And Solutions	0	3.526	187	11
- Histogram-aware sorting for enhanced word-aligned compression in bitmap indexes	0	3.543	17	1
- CubiST++: Evaluating Ad-Hoc CUBE Queries Using Statistics Trees	0	3.543	5	1
- Improving Performance of Sparse Matrix-Vector Multiplication	0	3.543	114	11
- Binary Gray Codes with Long Bit Runs	0	3.543	53	4
- Analysis of Basic Data Reordering Techniques	0	3.543	16	1
- Tree Based Indexes Versus Bitmap Indexes: A Performance Study	0	3.543	24	0
- Secondary indexing in one dimension: beyond b-trees and bitmap indexes	0	3.543	10	1
- A comparison of five probabilistic view-size estimation techniques in OLAP	0	3.543	24	1
- Compression techniques for fast external sorting	0	3.543	16	0
- A Note on Graph Coloring Extensions and List-Colorings	0	3.543	33	1
- Using Multiset Discrimination to Solve Language Processing Problems Without Hashing	0	3.543	52	2
- Monotone Gray Codes and the Middle Levels Problem	0	3.543	80	5
- The Art in Computer Programming	0	3.543	9242	678
- An Efficient Multi-Component Indexing Embedded Bitmap Compression for Data Reorganization	0	3.543	8	2
- The LitOLAP Project: Data Warehousing with Literature	0	3.543	8	0
- Multi-resolution bitmap indexes for scientific data	0	3.583	96	3
- Notes on design and implementation of compressed bit vectors	0	3.583	81	12
- Compressing Large Boolean Matrices using Reordering Techniques	0	3.595	88	7
- Compressing bitmap indices by data reorganization	1	3.595	53	4
- Model 204 Architecture and Performance	0	3.635	238	10
- On the performance of bitmap indices for high cardinality attributes	1	3.654	196	10
- A performance comparison of bitmap indexes	0	3.655	86	9
- Minimizing I/O Costs of Multi-Dimensional Queries with Bitmap Indices	0	3.692	16	0
- Evaluation Strategies for Bitmap Indices with Binning	0	3.692	69	3
- C-Store: A Column-oriented DBMS	0	3.710	1241	111
- Byte-aligned bitmap compression	0	3.793	209	48
- Bit Transposed Files	0	3.837	84	10
- Space efficient bitmap indexing	0	4.011	96	16

Table A.26: All the reference papers for a given study "Sorting improves word-aligned bitmap indexes." Among all its 31 references, we **boldface** the reference papers that are annotated to be significant influencers. For the three metrics, PCI, citations, and SSHI, we report their impact scores for each reference paper on the given study, where we mark a score **in green** when it conforms to the rule that a non-significant paper's value should be lower than that of a significant paper, and mark a score **in dark green** if it conforms to the rule to have a lower score than one of the significant paper, but violates the rule, i.e., having a higher score than the other significant paper. In this example, our PCI metric has an accuracy score of 79.3%, which is higher than both citations (68.1%), and SSHI (65.0%).

Bibliography

- Alberto Abadie. 2005. Semiparametric difference-in-differences estimators. *The Review of Economic Studies*, 72(1):1–19. 111, 173
- Alberto Abadie, Alexis Diamond, and Jens Hainmueller. 2010. Synthetic control methods for comparative case studies: Estimating the effect of california’s tobacco control program. *Journal of the American statistical Association*, 105(490):493–505. 119
- Alberto Abadie and Javier Gardeazabal. 2003. The economic costs of conflict: A case study of the basque country. *American economic review*, 93(1):113–132. 119
- Steven Abney. 2007. Semisupervised learning for computational linguistics. 82
- Christopher Adolph, Kenya Amano, Bree Bang-Jensen, Nancy Fullman, and John Wilkerson. 2021. Pandemic politics: Timing state-level social distancing responses to covid-19. *Journal of Health Politics, Policy and Law*, 46(2):211–233. 100
- Nicolas Ajzenman, Tiago Cavalcanti, and Daniel Da Mata. 2020. More than words: Leaders’ speech and risky behavior during a pandemic. *Available at SSRN 3582908*. 100
- Guillaume Alain and Yoshua Bengio. 2016. Understanding intermediate layers using linear classifier probes. *ArXiv*, abs/1610.01644. 38
- Tariq Alhindi, Tuhin Chakrabarty, Elena Musi, and Smaranda Muresan. 2022. Multitask instruction-based prompting for fallacy recognition. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8172–8187, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. 18
- Hunt Allcott, Levi Boxell, Jacob Conway, Matthew Gentzkow, Michael Thaler, and David Yang. 2020. Polarization and public health: Partisan differences in social distancing during the coronavirus pandemic. *Journal of Public Economics*, 191:104254. 99, 100, 106
- Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy P. Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul Ronald Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, and et al. 2023. Gemini: A family of highly capable multimodal models. *CoRR*, abs/2312.11805. 36

- Ravi Arunachalam and Sandipan Sarkar. 2013. The new eye of government: Citizen sentiment analysis in social media. In *Proceedings of the IJCNLP 2013 Workshop on Natural Language Processing for Social Media (SocialNLP)*, pages 23–28, Nagoya, Japan. Asian Federation of Natural Language Processing. 101
- Vincent Van Asch and Walter Daelemans. 2016. Predicting the effectiveness of self-training: Application to sentiment classification. *CoRR*, abs/1601.03288. 80
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*. 69
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova Das-Sarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *CoRR*, abs/2204.05862. 29
- David Baidoo-Anu and Leticia Owusu Ansah. 2023. Education in the era of generative artificial intelligence (ai): Understanding the potential benefits of chatgpt in promoting teaching and learning. *Available at SSRN* 4337484. 32
- Elias Bareinboim, Juan D. Correa, Duligur Ibeling, and Thomas Icard. 2022. On pearl’s hierarchy and the foundations of causal inference. In Hector Geffner, Rina Dechter, and Joseph Y. Halpern, editors, *Probabilistic and Causal Inference: The Works of Judea Pearl*, volume 36 of *ACM Books*, pages 507–556. ACM. 21, 22, 23
- Lisa Feldman Barrett. 2006. Solving the emotion paradox: Categorization and the experience of emotion. *Personality and social psychology review*, 10(1):20–46. 83, 93, 94
- John M Barrios and Yael Hochberg. 2020. Risk perception through the lens of politics in the time of the covid-19 pandemic. Technical report, National Bureau of Economic Research. 99, 100, 106
- Matthew A Baum. 2002. The constituent foundations of the rally-round-the-flag phenomenon. *International Studies Quarterly*, 46(2):263–298. 100
- RF Baumeister, E Bratslavsky, M Muraven, and DM Tice. 1998. Ego depletion: is the active self a limited resource? *Journal of Personality and Social Psychology*, 74(5):1252–1265. 87
- Helen Beebee, Christopher Hitchcock, Peter Charles Menzies, and Peter Menzies. 2009. *The Oxford handbook of causation*. Oxford Handbooks. 7
- Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017. What do neural machine translation models learn about morphology? In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 861–872, Vancouver, Canada. Association for Computational Linguistics. 38

- Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. 2023. Eliciting latent predictions from transformers with the tuned lens. *CoRR*, abs/2303.08112. 40, 48
- Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. SciBERT: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3615–3620, Hong Kong, China. Association for Computational Linguistics. 122, 176
- Taylor Berg-Kirkpatrick and Daniel Spokoyny. 2020. An empirical investigation of contextualized number prediction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4754–4764, Online. Association for Computational Linguistics. 62
- Donald J Berndt and James Clifford. 1994. Using dynamic time warping to find patterns in time series. In *KDD workshop*, volume 10, pages 359–370. Seattle, WA, USA:. 109
- Steven Bethard, William Corvey, Sara Klingenstein, and James H. Martin. 2008. Building a corpus of temporal-causal structure. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco. European Language Resources Association (ELRA). 32
- Gary Beverungen and Jugal Kalita. 2011. Evaluating methods for summarizing twitter posts. *Proceedings of the 5th AAAI ICWSM*. 101
- Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen-tau Yih, and Yejin Choi. 2020. Abductive commonsense reasoning. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net. 8, 150
- Michela Bia and Alessandra Mattei. 2008. A stata package for the estimation of the dose-response function through adjustment for the generalized propensity score. *The Stata Journal*, 8(3):354–373. 111, 173
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 2397–2430. PMLR. 41
- Felix Bittmann, Alexander Tekles, and Lutz Bornmann. 2021. Applied usage and performance of statistical matching in bibliometrics: The comparison of milestone and regular papers with multiple measurements of disruptiveness as an empirical example. *Quantitative Science Studies*, 2(4):1246–1270. 127
- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. 2022. Gpt-neox-20b: An open-source autoregressive language model. *CoRR*, abs/2204.06745. 58

- Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella Biderman. 2021. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow. 58
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022. 101, 103
- Léonard Blier and Yann Ollivier. 2018. The description length of deep learning models. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc. 75, 81
- Avrim Blum and Tom M. Mitchell. 1998. Combining labeled and unlabeled data with co-training. In *Proceedings of the Eleventh Annual Conference on Computational Learning Theory, COLT 1998, Madison, Wisconsin, USA, July 24-26, 1998*, pages 92–100. ACM. 78
- Nikolay Bogoychev and Rico Sennrich. 2019. Domain, translationese and noise in synthetic data for neural machine translation. *CoRR*, abs/1911.03362. 69
- Damien Bol, Marco Giani, André Blais, and Peter John Loewen. 2021. The effect of covid-19 lockdowns on political support: Some good news for democracy? *European Journal of Political Research*, 60(2):497–505. 100
- Stephan Bongers, Patrick Forré, Jonas Peters, and Joris M Mooij. 2021. Foundations of structural causal models with cycles and latent variables. *The Annals of Statistics*, 49(5):2885–2915. 140
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. COMET: Commonsense transformers for automatic knowledge graph construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4762–4779, Florence, Italy. Association for Computational Linguistics. 81
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics. 18
- Elizabeth M. Brannon. 2005. The independence of language and mathematical reasoning. *Proceedings of the National Academy of Sciences*, 102(9):3177–3178. 53
- Richard Brody. 1991. *Assessing the president: The media, elite opinion, and public support*. Stanford University Press. 100
- Richard A Brody and Catherine R Shapiro. 1989. A reconsideration of the rally phenomenon in public opinion. *Political behavior annual*, 2:77–102. 100
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. 1, 14, 20, 29, 30, 32, 36, 50, 58, 91, 176

- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott M. Lundberg, Harsha Nori, Hamid Palangi, Marco Túlio Ribeiro, and Yi Zhang. 2023. Sparks of artificial general intelligence: Early experiments with GPT-4. *CoRR*, abs/2303.12712. 29, 32
- Kailash Budhathoki and Jilles Vreeken. 2017. MDL for causal inference on discrete data. In *2017 IEEE International Conference on Data Mining, ICDM 2017, New Orleans, LA, USA, November 18-21, 2017*, pages 751–756. IEEE Computer Society. 81
- Isabel Cachola, Kyle Lo, Arman Cohan, and Daniel Weld. 2020. TLDR: Extreme summarization of scientific documents. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4766–4777, Online. Association for Computational Linguistics. 80
- Brandice Canes-Wrone, David W Brady, and John F Cogan. 2002. Out of step, out of office: Electoral accountability and house members’ voting. *American Political Science Review*, pages 127–140. 98
- Angela Cao, Gregor Williamson, and Jinho D. Choi. 2022. A cognitive approach to annotating causal constructions in a cross-genre corpus. In *Proceedings of the 16th Linguistic Annotation Workshop (LAW-XVI) within LREC2022*, pages 151–159, Marseille, France. European Language Resources Association. 32
- David Card. 1999. The causal effect of education on earnings. *Handbook of labor economics*, 3:1801–1863. 20, 33
- Håkan Carlsson. 2009. Allocation of research funds using bibliometric indicators—asset and challenge to swedish higher education sector. 113
- Devin Caughey and Christopher Warshaw. 2018. Policy preferences and policy change: Dynamic responsiveness in the american states, 1936–2014. *American Political Science Review*, 112(2):249–266. 98, 100
- US Census Bureau. 2012. *United States Summary, 2010: Population and housing unit counts*. US Department of Commerce, Economics and Statistics Administration, U.S. CENSUS BUREAU. 104
- Koyel Chakraborty, Surbhi Bhatia, Siddhartha Bhattacharyya, Jan Platos, Rajib Bag, and Aboul Ella Hassanien. 2020a. Sentiment analysis of COVID-19 tweets by deep learning classifiers - A study to show how popularity is affecting accuracy in social media. *Appl. Soft Comput.*, 97(Part):106754. 102
- Manajit Chakraborty, Maksym Byshkin, and Fabio Crestani. 2020b. Patent citation network analysis: A perspective from descriptive statistics and ergms. *Plos one*, 15(12):e0241797. 127
- Dhivya Chandrasekaran and Vijay Mago. 2022. Evolution of semantic similarity - A survey. *ACM Comput. Surv.*, 54(2):41:1–41:37. 114, 121
- Joseph Chee Chang, Amy X Zhang, Jonathan Bragg, Andrew Head, Kyle Lo, Doug Downey, and Daniel S Weld. 2023. Citesee: Augmenting citations in scientific papers with persistent and personalized historical context. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–15. 127

- Rafael Chaves, Lukas Luft, Thiago O. Maciel, David Gross, Dominik Janzing, and Bernhard Schölkopf. 2014. Inferring latent structures via information inequalities. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence, UAI 2014, Quebec City, Quebec, Canada, July 23-27, 2014*, pages 112–121. AUAI Press. 94
- Emily Chen, Kristina Lerman, and Emilio Ferrara. 2020. Tracking social media discourse about the covid-19 pandemic: Development of a public coronavirus twitter data set. *JMIR Public Health and Surveillance*, 6(2):e19273. 102, 111, 170, 171
- Emily Cheng, Diego Doimo, Corentin Kervadec, Iuri Macocco, Jade Yu, Alessandro Laio, and Marco Baroni. 2024. Emergence of a high-dimensional abstraction phase in language transformers. *arXiv preprint arXiv:2405.15471*. 38
- Raj Chetty, John N Friedman, Nathaniel Hilger, Emmanuel Saez, Diane Whitmore Schanzenbach, and Danny Yagan. 2011. How does your kindergarten classroom affect your earnings? Evidence from project star. *The Quarterly journal of economics*, 126(4):1593–1660. 20
- David Maxwell Chickering. 2002. Optimal structure identification with greedy search. *J. Mach. Learn. Res.*, 3:507–554. 10
- Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1724–1734. ACL. 69
- Gihyeon Choi, Shinhyeok Oh, and Harksoo Kim. 2020. Improving document-level sentiment classification using importance of sentences. *Entropy*, 22(12):1336. 83
- Yejin Choi and Claire Cardie. 2008. Learning with compositional semantics as structural inference for subsentential sentiment analysis. In *2008 Conference on Empirical Methods in Natural Language Processing, EMNLP 2008, Proceedings of the Conference, 25-27 October 2008, Honolulu, Hawaii, USA, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 793–801. ACL. 93
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*. 50, 62, 176
- Jennifer Chubb, Peter Cowling, and Darren Reed. 2022. Speeding up to keep up: exploring the use of ai in the research process. *AI & society*, 37(4):1439–1457. 127
- Kevin Clark, Minh-Thang Luong, Christopher D Manning, and Quoc V Le. 2018. Semi-supervised sequence modeling with cross-view training. *arXiv preprint arXiv:1809.08370*. 79
- Stephen Clark, James Curran, and Miles Osborne. 2003. Bootstrapping POS-taggers using unlabelled data. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pages 49–55. 81

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. 50, 62
- Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S. Weld. 2020. SPECTER: Document-level Representation Learning using Citation-informed Transformers. In *ACL*. 122
- Christopher Collins, Denis Dennehy, Kieran Conboy, and Patrick Mikalef. 2021. Artificial intelligence in information systems research: A systematic literature review and research agenda. *International Journal of Information Management*, 60:102383. 127
- Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. 2023. Towards automated circuit discovery for mechanistic interpretability. *CoRR*, abs/2304.14997. 41
- Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single \mathbb{R}^d vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics. 38
- Corinna Cortes and Neil D. Lawrence. 2021. Inconsistency in conference peer review: Revisiting the 2014 neurips experiment. *CoRR*, abs/2109.09774. 113
- Ernest D. Courant, Milton Stanley Livingston, and Hartland S. Snyder. 1952. The strong-focusing synchrotron—a new high energy accelerator. *Physical Review*, 88:1190–1196. 176
- Robert G Cowell, Philip Dawid, Steffen L Lauritzen, and David J Spiegelhalter. 2007. *Probabilistic networks and expert systems: Exact computational methods for Bayesian networks*. Springer Science & Business Media. 23
- Shaobo Cui, Zhijing Jin, Bernhard Schölkopf, and Boi Faltings. 2024. The odyssey of commonsense causality: From foundational benchmarks to cutting-edge reasoning. *CoRR*, abs/2406.19307. 129
- Andrew M. Dai and Quoc V. Le. 2015. Semi-supervised sequence learning. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 3079–3087. 79
- P. Daniušis, D. Janzing, J. Mooij, J. Zscheischler, B. Steudel, K. Zhang, and B. Schölkopf. 2010. Inferring deterministic causal relations. In *26th Conference on Uncertainty in Artificial Intelligence*, pages 143–150, Corvallis, OR. AUAI Press. Best student paper award. 73
- Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. 2023. Analyzing transformers in embedding space. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, *ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 16124–16170. Association for Computational Linguistics. 40, 48

- Ernest Davis and Scott Aaronson. 2023. Testing GPT-4 with Wolfram Alpha and Code Interpreter plug-ins on math and science problems. *arXiv preprint arXiv:2308.05713*. 33
- Shane Dawson, Dragan Gašević, George Siemens, and Srecko Joksimovic. 2014. Current state and future trends: A citation network analysis of the learning analytics field. In *Proceedings of the fourth international conference on learning analytics and knowledge*, pages 231–240. 127
- Gaston De Serres, France Markowski, Eveline Toth, Monique Landry, Danielle Auger, Marlène Mercier, Philippe Bélanger, Bruno Turmel, Horacio Arruda, Nicole Boulianne, et al. 2013. Largest measles epidemic in North America in a decade—Quebec, Canada, 2011: Contribution of susceptibility, serendipity, and superspreading events. *The Journal of infectious diseases*, 207(6):990–998. 20
- David DeFranza, Mike Lindow, Kevin Harrison, Arul Mishra, and Himanshu Mishra. 2020. Religion and reactance to covid-19 mitigation guidelines. *American Psychologist*. 99
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009)*, 20-25 June 2009, Miami, Florida, USA, pages 248–255. IEEE Computer Society. 121
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. 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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics. 1, 7, 14, 20, 32, 122, 125, 176
- Quang Do, Yee Seng Chan, and Dan Roth. 2011. Minimally supervised event causality identification. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 294–303, Edinburgh, Scotland, UK. Association for Computational Linguistics. 81, 150
- Diego Doimo, Aldo Glielmo, Alessio Ansuini, and Alessandro Laio. 2020. Hierarchical nucleation in deep neural networks. *Advances in Neural Information Processing Systems*, 33:7526–7536. 38
- Richard Doll and A Bradford Hill. 1950. Smoking and carcinoma of the lung. *British medical journal*, 2(4682):739. 20
- Richard Doll and A Bradford Hill. 1954. The mortality of doctors in relation to their smoking habits. *British medical journal*, 1(4877):1451. 20
- Naveen Donthu, Satish Kumar, Debmalya Mukherjee, Nitesh Pandey, and Weng Marc Lim. 2021. How to conduct a bibliometric analysis: An overview and guidelines. *Journal of business research*, 133:285–296. 126
- Jesse Dunietz, Lori Levin, and Jaime Carbonell. 2017. The BECauSE corpus 2.0: Annotating causality and overlapping relations. In *Proceedings of the 11th Linguistic Annotation Workshop*, pages 95–104, Valencia, Spain. Association for Computational Linguistics. 81, 150

- Sergey Edunov, Myle Ott, Marc'Aurelio Ranzato, and Michael Auli. 2020. On the evaluation of machine translation systems trained with back-translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2836–2846, Online. Association for Computational Linguistics. 69
- Naoki Egami, Christian J. Fong, Justin Grimmer, Margaret E. Roberts, and Brandon M. Stewart. 2018. How to make causal inferences using texts. *CoRR*, abs/1802.02163. 81
- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard H. Hovy, Hinrich Schütze, and Yoav Goldberg. 2021. Measuring and improving consistency in pretrained language models. *Trans. Assoc. Comput. Linguistics*, 9:1012–1031. 38
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*. <https://transformer-circuits.pub/2021/framework/index.html>. 36, 37, 38, 39
- Jordan Ellenberg. 2021. Coronavirus vaccines work. But this statistical illusion makes people think they don't. *The Washington Post*. 33
- Susan A. Elmore. 2018. The altmetric attention score: What does it mean and why should i care? *Toxicologic Pathology*, 46(3):252–255. PMID: 29448902. 127
- Seymour Epstein. 1994. Integration of the cognitive and the psychodynamic unconscious. *American psychologist*, 49(8):709. 84, 87
- Xing Fang and Justin Zhijun Zhan. 2015. Sentiment analysis using product review data. *Journal of Big Data*, 2:1–14. 176
- Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E. Roberts, Brandon M. Stewart, Victor Veitch, and Diyi Yang. 2021a. Causal inference in natural language processing: Estimation, prediction, interpretation and beyond. *CoRR*, abs/2109.00725. 62, 81
- Amir Feder, Nadav Oved, Uri Shalit, and Roi Reichart. 2021b. CausaLM: Causal model explanation through counterfactual language models. *Computational Linguistics*, 47(2):333–386. 62
- Hao Fei, Bobo Li, Qian Liu, Lidong Bing, Fei Li, and Tat-Seng Chua. 2023. Reasoning implicit sentiment with chain-of-thought prompting. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1171–1182, Toronto, Canada. Association for Computational Linguistics. 83
- Justin S Feinstein. 2013. Lesion studies of human emotion and feeling. *Current opinion in neurobiology*, 23(3):304–309. 83, 93
- Matthew Finlayson, Aaron Mueller, Sebastian Gehrmann, Stuart Shieber, Tal Linzen, and Yonatan Belinkov. 2021. Causal analysis of syntactic agreement mechanisms in neural

- language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1828–1843, Online. Association for Computational Linguistics. 51, 56, 62
- Ronald A. Fisher and E. B. Ford. 1927. The spread of a gene in natural conditions in a colony of the moth panaxia dominula l. *Heredity*, 11:143–174. Early work by Fisher on the application of randomization in agricultural experiments. 7
- Christian Fong and Justin Grimmer. 2016. Discovery of treatments from text corpora. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics. 81
- Markus Freitag, Isaac Caswell, and Scott Roy. 2019. APE at scale and its implications on MT evaluation biases. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 34–44, Florence, Italy. Association for Computational Linguistics. 69
- Jörg Frohberg and Frank Binder. 2022. CRASS: A novel data set and benchmark to test counterfactual reasoning of large language models. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 2126–2140, Marseille, France. European Language Resources Association. 150
- Shana Kushner Gadarian, Sara Wallace Goodman, and Thomas B Pepinsky. 2021. Partisanship, health behavior, and policy attitudes in the early stages of the covid-19 pandemic. *Plos one*, 16(4):e0249596. 99, 100, 106
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*. 57
- Eugene Garfield. 1964. Science citation index – a new dimension in indexing. *Science*, 144(3619):649–654. 126
- Eugene Garfield. 1972. Citation analysis as a tool in journal evaluation: Journals can be ranked by frequency and impact of citations for science policy studies. *Science*, 178(4060):471–479. 126
- Eugene Garfield, Irving H Sher, Richard J Torpie, et al. 1964. The use of citation data in writing the history of science. 126
- Gary Rosenberg Gary Holden and Kathleen Barker. 2005. Bibliometrics. *Social Work in Health Care*, 41(3-4):67–92. 113
- Oguzhan Gencoglu and Mathias Gruber. 2020. Causal modeling of twitter activity during covid-19. *Computation*, 8(4):85. 100
- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting recall of factual associations in auto-regressive language models. *CoRR*, abs/2304.14767. 36, 37, 38, 40, 43, 44, 47

- Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. 2022. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 30–45, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. 37, 40, 48
- Mor Geva, Ankit Gupta, and Jonathan Berant. 2020. Injecting numerical reasoning skills into language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 946–958, Online. Association for Computational Linguistics. 62
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 36, 37, 38
- Thomas A Glass, Steven N Goodman, Miguel A Hernán, and Jonathan M Samet. 2013. Causal inference in public health. *Annual review of public health*, 34:61–75. 33
- Clark Glymour, Kun Zhang, and Peter Spirtes. 2019. Review of causal discovery methods based on graphical models. *Frontiers in Genetics*, 10:524. 8, 9, 10
- Madelyn Glymour, Judea Pearl, and Nicholas P Jewell. 2016. *Causal inference in statistics: A primer*. John Wiley and Sons. 8, 10, 21, 135
- Moisés Goldszmidt and Judea Pearl. 1992. Rank-based systems: A simple approach to belief revision, belief update, and reasoning about evidence and actions. *KR*, 92:661–672. 22, 23
- Olga Golovneva, Moya Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2022. Roscoe: A suite of metrics for scoring step-by-step reasoning. 31, 143
- Charles J Gomez, Andrew C Herman, and Paolo Parigi. 2022. Leading countries in global science increasingly receive more citations than other countries doing similar research. *Nature Human Behaviour*, 6(7):919–929. 113
- Mingming Gong, Kun Zhang, Tongliang Liu, Dacheng Tao, Clark Glymour, and Bernhard Schölkopf. 2016. Domain adaptation with conditional transferable components. In *International conference on machine learning*, pages 2839–2848. PMLR. 72
- Andrew Gordon, Zornitsa Kozareva, and Melissa Roemmele. 2012. SemEval-2012 task 7: Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In **SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 394–398, Montréal, Canada. Association for Computational Linguistics. 8, 150
- Yvette Graham, Barry Haddow, and Philipp Koehn. 2020. Statistical power and translationese in machine translation evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 72–81, Online. Association for Computational Linguistics. 69

- Giovanni Grano, Andrea Di Sorbo, Francesco Mercaldo, Corrado A Visaggio, Gerardo Canfora, and Sebastiano Panichella. 2017. Software applications user reviews. 87, 89
- Guy Grossman, Soojong Kim, Jonah M Rexer, and Harsha Thirumurthy. 2020. Political partisanship influences behavioral responses to governors' recommendations for covid-19 prevention in the united states. *Proceedings of the National Academy of Sciences*, 117(39):24144–24153. 99, 100, 105
- Peter D Grünwald. 2007. *The minimum description length principle*. MIT press. 74
- Pedro Henrique Calais Guerra, Wagner Meira Jr., and Claire Cardie. 2014. Sentiment analysis on evolving social streams: how self-report imbalances can help. In *Seventh ACM International Conference on Web Search and Data Mining, WSDM 2014, New York, NY, USA, February 24-28, 2014*, pages 443–452. ACM. 103
- Çaglar Gülçehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, Loïc Barrault, Huei-Chi Lin, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2015. On using monolingual corpora in neural machine translation. *CoRR*, abs/1503.03535. 81
- Nizar Habash. 2008. Four techniques for online handling of out-of-vocabulary words in arabic-english statistical machine translation. In *ACL 2008, Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics, June 15-20, 2008, Columbus, Ohio, USA, Short Papers*, pages 57–60. The Association for Computer Linguistics. 160
- Danny Halawi, Jean-Stanislas Denain, and Jacob Steinhardt. 2023. Overthinking the truth: Understanding how language models process false demonstrations. *CoRR*, abs/2307.09476. 40, 48
- Joseph Y. Halpern and Judea Pearl. 2005a. Causes and explanations: A structural-model approach. part i: Causes. *The British Journal for the Philosophy of Science*, 56(4):843–887. 135
- Joseph Y Halpern and Judea Pearl. 2005b. Causes and explanations: A structural-model approach. part ii: Explanations. *The British journal for the philosophy of science*. 135
- Michael Hanna, Ollie Liu, and Alexandre Variengien. 2023. How does GPT-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model. *CoRR*, abs/2305.00586. 38, 41
- Yuval Noah Harari. 2014. *Sapiens: A brief history of humankind*. Random House. 20
- Vasileios Hatzivassiloglou and Kathleen R McKeown. 1997. Predicting the semantic orientation of adjectives. In *Proceedings of the 35th annual meeting of the association for computational linguistics and eighth conference of the european chapter of the association for computational linguistics*, pages 174–181. Association for Computational Linguistics. 93
- Vasileios Hatzivassiloglou and Janyce M Wiebe. 2000. Effects of adjective orientation and gradability on sentence subjectivity. In *Proceedings of the 18th conference on Computational linguistics-Volume 1*, pages 299–305. Association for Computational Linguistics. 83, 93
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. DeBERTa: Decoding-enhanced Bert with disentangled attention. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. 14

- James J Heckman, Lance J Lochner, and Petra E Todd. 2006. Earnings functions, rates of return and treatment effects: The mincer equation and beyond. *Handbook of the Economics of Education*, 1:307–458. 20
- Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2010. SemEval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 33–38, Uppsala, Sweden. Association for Computational Linguistics. 150
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*. 25, 28
- Miguel A Hernán and James M Robins. 2010. Causal inference. 114
- John Hewitt and Percy Liang. 2019. Designing and interpreting probes with control tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, Hong Kong, China. Association for Computational Linguistics. 38
- John Hewitt and Christopher D. Manning. 2019. A structural probe for finding syntax in word representations. 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)*, pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics. 38
- Christopher Hidey and Kathy McKeown. 2016. Identifying causal relations using parallel Wikipedia articles. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1424–1433, Berlin, Germany. Association for Computational Linguistics. 32
- Keisuke Hirano and Guido W Imbens. 2004. The propensity score with continuous treatments. *Applied Bayesian modeling and causal inference from incomplete-data perspectives*, 226164:73–84. 111, 173
- Jorge E Hirsch. 2005. An index to quantify an individual’s scientific research output. *Proceedings of the National academy of Sciences*, 102(46):16569–16572. 126
- Matthew Ho, Aditya Sharma, Justin Chang, Michael Saxon, Sharon Levy, Yujie Lu, and William Yang Wang. 2022. Wikiwhy: Answering and explaining cause-and-effect questions. *arXiv preprint arXiv:2210.12152*. 20
- Mickel Hoang, Oskar Alija Bihorac, and Jacobo Rouces. 2019. Aspect-based sentiment analysis using BERT. In *Proceedings of the 22nd Nordic Conference on Computational Linguistics, NoDaLiDa 2019, Turku, Finland, September 30 - October 2, 2019*, pages 187–196. Linköping University Electronic Press. 93
- Paul W Holland. 1986. Statistics and causal inference. *Journal of the American statistical Association*, 81(396):945–960. 108

- Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear. 90
- Mark Hopkins and Judea Pearl. 2007. Causality and counterfactuals in the situation calculus. *Journal of Logic and Computation*, 17(5):939–953. 135
- Pedram Hosseini, David A. Broniatowski, and Mona T. Diab. 2021. Predicting directionality in causal relations in text. *CoRR*, abs/2103.13606. 81
- Jianhua Hou. 2017. Exploration into the evolution and historical roots of citation analysis by referenced publication year spectroscopy. *Scientometrics*, 110:1437–1452. 126
- Patrik O. Hoyer, Dominik Janzing, Joris M. Mooij, Jonas Peters, and Bernhard Schölkopf. 2008. Nonlinear causal discovery with additive noise models. In *Advances in Neural Information Processing Systems 21, Proceedings of the Twenty-Second Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 8-11, 2008*, volume 21, pages 689–696. Curran Associates, Inc. 10, 94
- Yan Cathy Hua, Paul Denny, Katerina Taskova, and Jörg Wicker. 2023. A systematic review of aspect-based sentiment analysis (ABSA): domains, methods, and trends. *CoRR*, abs/2311.10777. 95
- Yimin Huang and Marco Valtorta. 2006. Pearl’s calculus of intervention is complete. In *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence*, pages 217–224. 23, 24
- Paul Hünermund and Elias Bareinboim. 2019. Causal inference and data fusion in econometrics. *CoRR*, abs/1912.09104. 23, 25, 33
- Dieuwke Hupkes, Sara Veldhoen, and Willem H. Zuidema. 2018. Visualisation and ‘diagnostic classifiers’ reveal how recurrent and recursive neural networks process hierarchical structure. *J. Artif. Intell. Res.*, 61:907–926. 38
- Ferenc Huszár. 2023. We May be Surprised Again: Why I take LLMs seriously. 21
- B Ian Hutchins, Xin Yuan, James M Anderson, and George M Santangelo. 2016. Relative citation ratio (rcr): a new metric that uses citation rates to measure influence at the article level. *PLoS biology*, 14(9):e1002541. 127
- Oana Ignat, Zhijing Jin, Artem Abzaliev, Laura Biester, Santiago Castro, Naihao Deng, Xinyi Gao, Aylin Gunal, Jacky He, Ashkan Kazemi, Muhammad Khalifa, Namho Koh, Andrew Lee, Siyang Liu, Do June Min, Shinka Mori, Joan Nwatu, Verónica Pérez-Rosas, Siqi Shen, Zekun Wang, Winston Wu, and Rada Mihalcea. 2024. Has it all been solved? Open NLP research questions not solved by large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*. International Committee on Computational Linguistics. 1, 20, 32
- Hal Daumé III and Jagadeesh Jagarlamudi. 2011. Domain adaptation for machine translation by mining unseen words. In *The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference, 19-24 June, 2011, Portland, Oregon, USA - Short Papers*, pages 407–412. The Association for Computer Linguistics. 160

- Guido W Imbens and Donald B Rubin. 2015. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press. 173
- Sehrish Iqbal, Saeed-Ul Hassan, Naif Radi Aljohani, Salem Alelyani, Raheel Nawaz, and Lutz Bornmann. 2021. A decade of in-text citation analysis based on natural language processing and machine learning techniques: An overview of empirical studies. *Scientometrics*, 126(8):6551–6599. 127
- Dominik Janzing. 2007. On causally asymmetric versions of occam’s razor and their relation to thermodynamics. *arXiv preprint arXiv:0708.3411*. 94
- Dominik Janzing. 2019. The cause-effect problem: Motivation, ideas, and popular misconceptions. In Isabelle Guyon, Alexander R. Statnikov, and Berna Bakir Batu, editors, *Cause Effect Pairs in Machine Learning*, pages 3–26. Springer. 85, 94
- Dominik Janzing, Rafael Chaves, and Bernhard Schölkopf. 2016. Algorithmic independence of initial condition and dynamical law in thermodynamics and causal inference. *New Journal of Physics*, 18(9):093052. 94
- Dominik Janzing, Joris M. Mooij, Kun Zhang, Jan Lemeire, Jakob Zscheischler, Povilas Daniusis, Bastian Steudel, and Bernhard Schölkopf. 2012. Information-geometric approach to inferring causal directions. *Artif. Intell.*, 182-183:1–31. 73, 94
- Dominik Janzing and Bernhard Schölkopf. 2010. Causal inference using the algorithmic Markov condition. *IEEE Transactions on Information Theory*, 56(10):5168–5194. 70, 74, 86, 94
- Sébastien Jean, Orhan Firat, Kyunghyun Cho, Roland Memisevic, and Yoshua Bengio. 2015. Montreal neural machine translation systems for wmt’15. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 134–140. 81
- Chaker Jebari, Enrique Herrera-Viedma, and Manuel Jesus Cobo. 2021. The use of citation context to detect the evolution of research topics: a large-scale analysis. *Scientometrics*, 126(4):2971–2989. 127
- Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438. 21, 38
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. Is BERT really robust? A strong baseline for natural language attack on text classification and entailment. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 8018–8025. AAAI Press. 17
- Zhihua Jin, Xin Jiang, Xingbo Wang, Qun Liu, Yong Wang, Xiaozhe Ren, and Huamin Qu. 2021a. NumGPT: Improving numeracy ability of generative pre-trained models. *arXiv preprint arXiv:2109.03137*. 56
- Zhijing Jin, Yuen Chen, Felix Leeb, Luigi Gresele, Ojasv Kamal, Zhiheng Lyu, Kevin Blin, Fernando Gonzalez, Max Kleiman-Weiner, Mrinmaya Sachan, and Bernhard Schölkopf. 2023a. CLadder: Assessing causal reasoning in language models. In *NeurIPS*. 62, 129

- Zhijing Jin, Amir Feder, and Kun Zhang. 2022a. CausalNLP tutorial: An introduction to causality for natural language processing. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, pages 17–22, Abu Dhabi, UAE. Association for Computational Linguistics. 62
- Zhijing Jin, Abhinav Lalwani, Tejas Vaidhya, Xiaoyu Shen, Yiwen Ding, Zhiheng Lyu, Mrinmaya Sachan, Rada Mihalcea, and Bernhard Schölkopf. 2022b. Logical fallacy detection. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7180–7198, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. 18, 19, 32
- Zhijing Jin, Jiarui Liu, Zhiheng Lyu, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona T. Diab, and Bernhard Schölkopf. 2024. Can large language models infer causation from correlation? In *The Twelfth International Conference on Learning Representations, ICLR 2024*. OpenReview.net. 29, 32, 62, 129, 150
- Zhijing Jin, Zhiheng Lyu, Yiwen Ding, Mrinmaya Sachan, Kun Zhang, Rada Mihalcea, and Bernhard Schölkopf. 2023b. AI Scholars: A dataset for NLP-involved causal inference. 62
- Zhijing Jin, Zeyu Peng, Tejas Vaidhya, Bernhard Schölkopf, and Rada Mihalcea. 2021b. Mining the cause of political decision-making from social media: A case study of COVID-19 policies across the US states. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 288–301, Punta Cana, Dominican Republic. Association for Computational Linguistics. 62, 81
- Zhijing Jin, Julius von Kügelgen, Jingwei Ni, Tejas Vaidhya, Ayush Kaushal, Mrinmaya Sachan, and Bernhard Schölkopf. 2021c. Causal direction of data collection matters: Implications of causal and anticausal learning for NLP. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9499–9513, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 62, 84, 86, 94
- Thorsten Joachims. 1999. Transductive inference for text classification using support vector machines. In *Proceedings of the Sixteenth International Conference on Machine Learning (ICML 1999)*, Bled, Slovenia, June 27 - 30, 1999, pages 200–209. Morgan Kaufmann. 78
- Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics. 76, 160
- Daniel Kahneman. 2011. *Thinking, Fast and Slow*. Macmillan. 84, 87, 94
- Daniel Kahneman and Shane Frederick. 2002. *Representativeness revisited: Attribute substitution in intuitive judgment*, pages 49–81. Cambridge University Press. 87
- Daniel Kahneman, Barbara L Fredrickson, Charles A Schreiber, and Donald A Redelmeier. 1993. When more pain is preferred to less: Adding a better end. *Psychological science*, 4(6):401–405. 84, 88, 94
- Immanuel Kant. 1781. *Critique of Pure Reason*. Cambridge University Press. 7

- Karim S Kassam and Wendy Berry Mendes. 2013. The effects of measuring emotion: Physiological reactions to emotional situations depend on whether someone is asking. *PloS one*, 8(6):e64959. 83, 94
- Daniel Martin Katz, Michael James Bommarito, Shang Gao, and Pablo Arredondo. 2023. Gpt-4 passes the bar exam. *Available at SSRN* 4389233. 29, 32
- Chhinder Kaur and Anand Sharma. 2020. Twitter sentiment analysis on coronavirus using textblob. Technical report, EasyChair. 102
- Harleen Kaur, Shafqat Ul Ahsaan, Bhavya Alankar, and Victor Chang. 2021. A proposed sentiment analysis deep learning algorithm for analyzing covid-19 tweets. *Information Systems Frontiers*, pages 1–13. 102
- Divyansh Kaushik, Eduard H. Hovy, and Zachary Chase Lipton. 2020. Learning the difference that makes A difference with counterfactually-augmented data. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net. 56
- Katherine A. Keith, David Jensen, and Brendan O’Connor. 2020. Text and causal inference: A review of using text to remove confounding from causal estimates. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 5332–5344. Association for Computational Linguistics. 81
- Armin Kekić, Jonas Dehning, Luigi Gresele, Julius von Kügelgen, Viola Priesemann, and Bernhard Schölkopf. 2023. Evaluating vaccine allocation strategies using simulation-assisted causal modeling. *Patterns*. 20
- Phillip Keung, Yichao Lu, György Szarvas, and Noah A. Smith. 2020. The multilingual Amazon reviews corpus. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4563–4568, Online. Association for Computational Linguistics. 83, 85, 87, 89, 93
- Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. 2023. Causal reasoning and large language models: Opening a new frontier for causality. *arXiv preprint arXiv:2305.00050*. 8, 18, 32, 33
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1746–1751. ACL. 93
- Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization. 102
- Rodney Kinney, Chloe Anastasiades, Russell Authur, Iz Beltagy, Jonathan Bragg, Alexandra Buraczynski, Isabel Cachola, Stefan Candra, Yoganand Chandrasekhar, Arman Cohan, Miles Crawford, Doug Downey, Jason Dunkelberger, Oren Etzioni, Rob Evans, Sergey Feldman, Joseph Gorney, David Graham, Fangzhou Hu, Regan Huff, Daniel King, Sebastian Kohlmeier, Bailey Kuehl, Michael Langan, Daniel Lin, Haokun Liu, Kyle Lo, Jaron Lochner, Kelsey MacMillan, Tyler Murray, Chris Newell, Smita Rao, Shaurya Rohatgi, Paul Sayre, Zejiang Shen, Amanpreet Singh, Luca Soldaini, Shivashankar Subramanian, Amber Tanaka, Alex D. Wade, Linda Wagner, Lucy Lu Wang,

- Chris Wilhelm, Caroline Wu, Jiangjiang Yang, Angele Zamarron, Madeleine van Zuylen, and Daniel S. Weld. 2023. The semantic scholar open data platform. *CoRR*, abs/2301.10140. 115, 121, 122, 176
- Teun Kloek and Herman K. van Dijk. 1976. Bayesian estimates of equation system parameters, an application of integration by monte carlo. *Econometrica*, 46:1–19. 121
- Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *MT summit*, volume 5, pages 79–86. Citeseer. 76
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *CoRR*, abs/2205.11916. 62
- Andrei N Kolmogorov. 1965. Three approaches to the quantitative definition of information. *Problems of information transmission*, 1(1):1–7. 74
- Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. 2016. MAWPS: A math word problem repository. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1152–1157, San Diego, California. Association for Computational Linguistics. 57
- Moritz UG Kraemer, Chia-Hung Yang, Bernardo Gutierrez, Chieh-Hsi Wu, Brennan Klein, David M Pigott, Louis Du Plessis, Nuno R Faria, Ruoran Li, William P Hanage, et al. 2020. The effect of human mobility and control measures on the covid-19 epidemic in china. *Science*, 368(6490):493–497. 100
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*. 72
- Yash Kumar Lal, Nathanael Chambers, Raymond Mooney, and Niranjan Balasubramanian. 2021. TellMeWhy: A dataset for answering why-questions in narratives. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 596–610, Online. Association for Computational Linguistics. 150
- Jeffrey R Lax and Justin H Phillips. 2009. Gay rights in the states: Public opinion and policy responsiveness. *American Political Science Review*, 103(3):367–386. 98, 100
- Jeffrey R Lax and Justin H Phillips. 2012. The democratic deficit in the states. *American Journal of Political Science*, 56(1):148–166. 98, 100
- Joseph E LeDoux. 1998. *The emotional brain: The mysterious underpinnings of emotional life*. Simon and Schuster. 87
- Jan Lemeire and Erik Dirkx. 2006. Causal models as minimal descriptions of multivariate systems. 72
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics. 14

- Ming Li, Paul Vitányi, et al. 2008. *An introduction to Kolmogorov complexity and its applications*, volume 3. Springer. 74
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692. 14, 122, 126
- Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. S2ORC: The semantic scholar open research corpus. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4969–4983, Online. Association for Computational Linguistics. 115, 122, 176
- Joshua R Loftus, Chris Russell, Matt J Kusner, and Ricardo Silva. 2018. Causal reasoning for algorithmic fairness. *arXiv preprint arXiv:1805.05859*. 33
- Steven Loria. 2018. TextBlob documentation. *Release 0.15*, 2. 102
- Haoyan Luo and Lucia Specia. 2024. From understanding to utilization: A survey on explainability for large language models. *CoRR*, abs/2401.12874. 1
- Zhiheng Lyu, Zhijing Jin, Fernando Gonzalez, Rada Mihalcea, Bernhard Schölkopf, and Mrinmaya Sachan. 2024. On the causal nature of sentiment analysis. *CoRR*, abs/2404.11055. 62
- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pages 142–150. 85
- Bill MacCartney and Christopher D. Manning. 2008. Modeling semantic containment and exclusion in natural language inference. In *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pages 521–528, Manchester, UK. Coling 2008 Organizing Committee. 18
- Cristian R. Machuca, Cristian Gallardo, and Renato M. Toasa. 2021. Twitter sentiment analysis on coronavirus: Machine learning approach. *Journal of Physics: Conference Series*, 1828(1):012104. 102
- Goutam Majumder, Partha Pakray, Alexander Gelbukh, and David Pinto. 2016. Semantic textual similarity methods, tools, and applications: A survey. *Computación y Sistemas*, 20. 114, 121
- Paolo Manghi, Andrea Mannocci, Francesco Osborne, Dimitris Sacharidis, Angelo Salatino, and Thanasis Vergoulis. 2021. New trends in scientific knowledge graphs and research impact assessment. 127
- Kamran H Manguri, Rebaz N Ramadhan, and Pshko R Mohammed Amin. 2020. Twitter sentiment analysis on worldwide covid-19 outbreaks. *Kurdistan Journal of Applied Research*, pages 54–65. 102
- Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. Introduction to information retrieval. In *J. Assoc. Inf. Sci. Technol.* 114, 121

- Muvazima Mansoor, Kirthika Gurumurthy, Anantharam R. U, and V. R. Badri Prasad. 2020. Global sentiment analysis of COVID-19 tweets over time. *CoRR*, abs/2010.14234. 102
- Micol Marchetti-Bowick and Nathanael Chambers. 2012. Learning for microblogs with distant supervision: Political forecasting with Twitter. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 603–612, Avignon, France. Association for Computational Linguistics. 101
- Alexander Marx and Jilles Vreeken. 2021. Formally justifying mdl-based inference of cause and effect. *CoRR*, abs/2105.01902. 81
- Callum McDougall, Arthur Conmy, Cody Rushing, Thomas McGrath, and Neel Nanda. 2023. Copy suppression: Comprehensively understanding an attention head. *CoRR*, abs/2310.04625. 45, 47
- Matt McGue, Merete Osler, and Kaare Christensen. 2010. Causal inference and observational research: The utility of twins. *Perspectives on psychological science*, 5(5):546–556. 118
- Brendan D. McKay and Adolfo Piperno. 2014. Practical graph isomorphism, II. *J. Symb. Comput.*, 60:94–112. 12
- Ian R. McKenzie, Alexander Lyzhov, Michael Pieler, Alicia Parrish, Aaron Mueller, Ameya Prabhu, Euan McLean, Aaron Kirtland, Alexis Ross, Alisa Liu, Andrew Gritsevskiy, Daniel Wurgaft, Derik Kauffman, Gabriel Recchia, Jiacheng Liu, Joe Cavanagh, Max Weiss, Sicong Huang, The Floating Droid, Tom Tseng, Tomasz Korbak, Xudong Shen, Yuhui Zhang, Zhengping Zhou, Najoung Kim, Samuel R. Bowman, and Ethan Perez. 2023. Inverse scaling: When bigger isn’t better. *CoRR*, abs/2306.09479. 47
- Sergio Hernan Garrido Mejia, Elke Kirschbaum, and Dominik Janzing. 2022. Obtaining causal information by merging datasets with MAXENT. In *International Conference on Artificial Intelligence and Statistics, AISTATS 2022, 28-30 March 2022, Virtual Event*, volume 151 of *Proceedings of Machine Learning Research*, pages 581–603. PMLR. 94
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. *arXiv preprint arXiv:2202.05262*. 36, 37, 38, 40, 41, 51, 62, 92, 161
- Osman Ali Mian, Alexander Marx, and Jilles Vreeken. 2021. Discovering fully oriented causal networks. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 8975–8982. AAAI Press. 81
- Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. 2020. A diverse corpus for evaluating and developing English math word problem solvers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 975–984, Online. Association for Computational Linguistics. 57
- Tomás Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances*

- in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, pages 3111–3119. 46
- Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Noisy channel language model prompting for few-shot text classification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5316–5330, Dublin, Ireland. Association for Computational Linguistics. 94
- Paramita Mirza, Rachele Sprugnoli, Sara Tonelli, and Manuela Speranza. 2014. Annotating causality in the TempEval-3 corpus. In *Proceedings of the EACL 2014 Workshop on Computational Approaches to Causality in Language (CAtoCL)*, pages 10–19, Gothenburg, Sweden. Association for Computational Linguistics. 150
- Paramita Mirza and Sara Tonelli. 2014. An analysis of causality between events and its relation to temporal information. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 2097–2106, Dublin, Ireland. Dublin City University and Association for Computational Linguistics. 81
- Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay Rajpurohit, Oyvind Tafjord, Ashish Sabharwal, Peter Clark, et al. 2022a. Lila: A unified benchmark for mathematical reasoning. *arXiv preprint arXiv:2210.17517*. 52
- Swaroop Mishra, Arindam Mitra, Neeraj Varshney, Bhavdeep Sachdeva, Peter Clark, Chitta Baral, and Ashwin Kalyan. 2022b. NumGLUE: A suite of fundamental yet challenging mathematical reasoning tasks. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3505–3523, Dublin, Ireland. Association for Computational Linguistics. 50
- Henk F Moed. 2006. *Citation analysis in research evaluation*, volume 9. Springer Science & Business Media. 113
- Saif M. Mohammad, Xiaodan Zhu, Svetlana Kiritchenko, and Joel D. Martin. 2015. Sentiment, emotion, purpose, and style in electoral tweets. *Inf. Process. Manag.*, 51(4):480–499. 101
- Martin M Monti, Lawrence M Parsons, and Daniel N Osherson. 2012. Thought beyond language: Neural dissociation of algebra and natural language. *Psychological science*, 23(8):914–922. 53
- Joris M. Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. 2014. Distinguishing cause from effect using observational data: Methods and benchmarks. *CoRR*, abs/1412.3773. 94
- Rodrigo Moraes, João Francisco Valiati, and Wilson P Gavião Neto. 2013. Document-level sentiment classification: An empirical comparison between svm and ann. *Expert Systems with Applications*, 40(2):621–633. 93
- Jeffrey Morris. 2021. Israeli data: How can efficacy vs. severe disease be strong when 60% of hospitalized are vaccinated. *Covid Data Science*. Accessed: 27th of October 2023. 33

- John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. TextAttack: A framework for adversarial attacks, data augmentation, and adversarial training in NLP. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 119–126, Online. Association for Computational Linguistics. 17
- Nasrin Mostafazadeh, Alyson Grealish, Nathanael Chambers, James Allen, and Lucy Vanderwende. 2016. CaTeRS: Causal and temporal relation scheme for semantic annotation of event structures. In *Proceedings of the Fourth Workshop on Events*, pages 51–61, San Diego, California. Association for Computational Linguistics. 150
- John E Mueller. 1970. Presidential popularity from truman to johnson. *The American Political Science Review*, 64(1):18–34. 100
- John E Mueller. 1973. *War, presidents, and public opinion*. New York: Wiley. 100
- Martin Müller, Marcel Salathé, and Per Egil Kummervold. 2020. Covid-twitter-bert: A natural language processing model to analyse COVID-19 content on twitter. *CoRR*, abs/2005.07503. 102
- Daniel Naber et al. 2003. A rule-based style and grammar checker. 27
- Razieh Nabi and Ilya Shpitser. 2018. Fair inference on outcomes. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 1931–1940. AAAI Press. 33
- Neel Nanda and Joseph Bloom. 2022. Transformerlens. <https://github.com/neelnanda-io/TransformerLens>. 41
- Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. 2023. Progress measures for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. 38, 39
- Tetsuya Nasukawa and Jeonghee Yi. 2003. Sentiment analysis: Capturing favorability using natural language processing. In *Proceedings of the 2nd international conference on Knowledge capture*, pages 70–77. ACM. 83, 93
- Brady Neal. 2020. *Introduction to causal inference*. 135
- Jingwei Ni, Zhijing Jin, Markus Freitag, Mrinmaya Sachan, and Bernhard Schölkopf. 2022. Original or translated? A causal analysis of the impact of translationese on machine translation performance. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5303–5320, Seattle, United States. Association for Computational Linguistics. 62, 80, 86
- Kamal Nigam and Rayid Ghani. 2000. Analyzing the effectiveness and applicability of co-training. In *Proceedings of the 2000 ACM CIKM International Conference on Information and Knowledge Management, McLean, VA, USA, November 6-11, 2000*, pages 86–93. ACM. 80

- Kamal Nigam, Andrew McCallum, and Tom M. Mitchell. 2006. Semi-supervised text classification using EM. In Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien, editors, *Semi-Supervised Learning*, pages 32–55. The MIT Press. 79
- Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of GPT-4 on medical challenge problems. *CoRR*, abs/2303.13375. 29, 32
- Nostalgebraist. 2020. interpreting gpt: the logit lens. Accessed: Nov 2023. 37, 40
- Maxwell I. Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Augustus Odena. 2021. Show your work: Scratchpads for intermediate computation with language models. *CoRR*, abs/2112.00114. 28
- Jong-Hoon Oh, Kentaro Torisawa, Chikara Hashimoto, Ryu Iida, Masahiro Tanaka, and Julien Kloetzer. 2016. A semi-supervised learning approach to why-question answering. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*, pages 3022–3029. AAAI Press. 81
- Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. 2020. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001. 38, 39
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2022. In-context learning and induction heads. *CoRR*, abs/2209.11895. 36, 37, 38, 39
- OpenAI. 2023. GPT-4 technical report. *CoRR*, abs/2303.08774. 1, 7, 14, 20, 29, 30, 32, 36, 91
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. *CoRR*, abs/2203.02155. 7, 14, 29, 30, 50, 51, 58, 61, 91
- Alexander Pak and Patrick Paroubek. 2010. Twitter as a corpus for sentiment analysis and opinion mining. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC’10)*, Valletta, Malta. European Language Resources Association (ELRA). 101
- Kuntal Kumar Pal and Chitta Baral. 2021. Investigating numeracy learning ability of a text-to-text transfer model. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3095–3101, Punta Cana, Dominican Republic. Association for Computational Linguistics. 62
- Georgios Paltoglou and Mike Thelwall. 2012. Twitter, myspace, digg: Unsupervised sentiment analysis in social media. *ACM Trans. Intell. Syst. Technol.*, 3(4):66:1–66:19. 101

- Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment classification using machine learning techniques. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing, EMNLP 2002, Philadelphia, PA, USA, July 6-7, 2002*, pages 79–86. 93
- Kiho Park, Yo Joong Choe, Yibo Jiang, and Victor Veitch. 2024. The geometry of categorical and hierarchical concepts in large language models. 38
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online. Association for Computational Linguistics. 50, 52, 57, 156
- Judea Pearl. 1988. *Probabilistic reasoning in intelligent systems: Networks of plausible inference*. Morgan Kaufmann. 9, 23
- Judea Pearl. 1995. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688. 21, 22, 23, 26, 51, 53, 54, 100, 107, 108, 141
- Judea Pearl. 2001. Direct and indirect effects. In *UAI '01: Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence, University of Washington, Seattle, Washington, USA, August 2-5, 2001*, pages 411–420. Morgan Kaufmann. 51, 54, 55, 92
- Judea Pearl. 2009a. Causal inference in statistics: An overview. *Statistics Surveys*, 3(none):96 – 146. 135
- Judea Pearl. 2009b. *Causality: Models, reasoning and inference (2nd ed.)*. Cambridge University Press. 7, 21, 22, 23, 25, 26, 61, 69, 85, 117, 135, 140, 141
- Judea Pearl. 2011. The algorithmization of counterfactuals. *Annals of Mathematics and Artificial Intelligence*, 61:29–39. 22
- Judea Pearl. 2022. Comment: Understanding simpson’s paradox. In Hector Geffner, Rina Dechter, and Joseph Y. Halpern, editors, *Probabilistic and Causal Inference: The Works of Judea Pearl*, volume 36 of *ACM Books*, pages 399–412. ACM. 21
- Judea Pearl and Elias Bareinboim. 2022. External validity: From *Do*-calculus to transportability across populations. In Hector Geffner, Rina Dechter, and Joseph Y. Halpern, editors, *Probabilistic and Causal Inference: The Works of Judea Pearl*, volume 36 of *ACM Books*, pages 451–482. ACM. 23
- Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. 2016. *Causal inference in statistics: A primer*. John Wiley & Sons. 22
- Judea Pearl and Dana Mackenzie. 2018. *The book of why: The new science of cause and effect*. Basic books. 2, 20, 21, 22, 23, 25, 26, 28, 29, 32, 135
- Derek C Penn and Daniel J Povinelli. 2007. Causal cognition in human and nonhuman animals: A comparative, critical review. *Annu. Rev. Psychol.*, 58:97–118. 20
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. 2010. Identifying cause and effect on discrete data using additive noise models. In *Proceedings of the Thirteenth International*

- Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*, volume 9 of *JMLR Proceedings*, pages 597–604. JMLR.org. 94
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. 2017. *Elements of causal inference: Foundations and learning algorithms*. The MIT Press. 7, 21, 25, 61, 69, 70, 72, 85, 135
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China. Association for Computational Linguistics. 21, 38
- Piotr Piękos, Mateusz Malinowski, and Henryk Michalewski. 2021. Measuring and improving BERT’s mathematical abilities by predicting the order of reasoning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 383–394, Online. Association for Computational Linguistics. 62
- Fredrik Niclas Piro and Gunnar Sivertsen. 2016. How can differences in international university rankings be explained? *Scientometrics*, 109(3):2263–2278. 113
- Ferran Pla and Lluís-F. Hurtado. 2014. Political tendency identification in Twitter using sentiment analysis techniques. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 183–192, Dublin, Ireland. Dublin City University and Association for Computational Linguistics. 101
- Barbara Plank, Anders Johannsen, and Anders Søgaard. 2014. Importance weighting and unsupervised domain adaptation of POS taggers: a negative result. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 968–973. ACL. 82
- Drago Plecko and Elias Bareinboim. 2022. Causal fairness analysis. *arXiv preprint arXiv:2207.11385*. 33
- Stanley A Plotkin. 2005. Vaccines: Past, present and future. *Nature medicine*, 11(Suppl 4):S5–S11. 20
- Maria Pontiki, Dimitris Galanis, John Pavlopoulos, Harris Papageorgiou, Ion Androutsopoulos, and Suresh Manandhar. 2014. Semeval-2014 task 4: Aspect based sentiment analysis. In *Proceedings of the 8th International Workshop on Semantic Evaluation, SemEval@COLING 2014, Dublin, Ireland, August 23-24, 2014*, pages 27–35. The Association for Computer Linguistics. 95
- Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, and Rada Mihalcea. 2023. Beneath the tip of the iceberg: Current challenges and new directions in sentiment analysis research. *IEEE Trans. Affect. Comput.*, 14(1):108–132. 83, 90
- Lutz Prechelt, Daniel Graziotin, and Daniel Méndez Fernández. 2018. A community’s perspective on the status and future of peer review in software engineering. *Information and Software Technology*, 95:75–85. 113

- George Psacharopoulos and Harry Anthony Patrinos. 2004. Returns to investment in education: A further update. *Education economics*, 12(2):111–134. 20
- Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is chatgpt a general-purpose natural language processing task solver? *CoRR*, abs/2302.06476. 32
- Lianhui Qin, Antoine Bosselut, Ari Holtzman, Chandra Bhagavatula, Elizabeth Clark, and Yejin Choi. 2019. Counterfactual story reasoning and generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5043–5053, Hong Kong, China. Association for Computational Linguistics. 8, 18, 32, 150
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. 176
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8). 1, 7, 14, 20, 29, 32, 41, 57, 76, 91, 160, 176
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67. 93
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics. 72
- Alan Ramponi and Barbara Plank. 2020. Neural unsupervised domain adaptation in NLP—A survey. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6838–6855, Barcelona, Spain (Online). International Committee on Computational Linguistics. 81
- Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A. Smith, and Yejin Choi. 2018. Event2Mind: Commonsense inference on events, intents, and reactions. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 463–473, Melbourne, Australia. Association for Computational Linguistics. 150
- Tilman R  uker, Anson Ho, Stephen Casper, and Dylan Hadfield-Menell. 2023. Toward transparent AI: A survey on interpreting the inner structures of deep neural networks. In *2023 IEEE Conference on Secure and Trustworthy Machine Learning, SaTML 2023, Raleigh, NC, USA, February 8-10, 2023*, pages 464–483. IEEE. 1
- Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. 2022. Impact of pretraining term frequencies on few-shot numerical reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 840–854, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. 29, 50, 62

- Andrew J. Reagan, Lewis Mitchell, Dilan Kiley, Christopher M. Danforth, and Peter Sheridan Dodds. 2016. The emotional arcs of stories are dominated by six basic shapes. *EPJ Data Sci.*, 5(1):31. 88, 164
- Radim Řehůřek and Petr Sojka. 2010. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta. ELRA. <http://is.muni.cz/publication/884893/en>. 46
- Roi Reichart and Ari Rappoport. 2007. Self-training for enhancement and domain adaptation of statistical parsers trained on small datasets. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 616–623, Prague, Czech Republic. Association for Computational Linguistics. 82
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Conference on Empirical Methods in Natural Language Processing*. 179
- David B Resnik, Christina Gutierrez-Ford, and Shyamal Peddada. 2008. Perceptions of ethical problems with scientific journal peer review: an exploratory study. *Science and engineering ethics*, 14(3):305–310. 113
- Parker Riley, Isaac Caswell, Markus Freitag, and David Grangier. 2020. Translationese as a language in “multilingual” NMT. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7737–7746, Online. Association for Computational Linguistics. 69, 80
- Jorma Rissanen. 1984. Universal coding, information, prediction, and estimation. *IEEE Trans. Inf. Theory*, 30(4):629–636. 75
- Margaret E Roberts, Brandon M Stewart, and Richard A Nielsen. 2020. Adjusting for confounding with text matching. *American Journal of Political Science*, 64(4):887–903. 62, 81
- Stephen E. Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. *Found. Trends Inf. Retr.*, 3:333–389. 121
- Anna Rogers, Marzena Karpinska, Jordan Boyd-Graber, and Naoaki Okazaki. 2023. Program chairs’ report on peer review at acl 2023. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages xl–lxxv, Toronto, Canada. Association for Computational Linguistics. 113
- Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-resolution image synthesis with latent diffusion models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10674–10685. 176
- Paul R Rosenbaum and Donald B Rubin. 1983. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55. 114, 118
- Sara Rosenthal, Noura Farra, and Preslav Nakov. 2017. SemEval-2017 task 4: Sentiment analysis in Twitter. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 502–518, Vancouver, Canada. Association for Computational Linguistics. 102

- Sara Rosenthal, Preslav Nakov, Svetlana Kiritchenko, Saif Mohammad, Alan Ritter, and Veselin Stoyanov. 2015. Semeval-2015 task 10: Sentiment analysis in twitter. In *Proceedings of the 9th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2015, Denver, Colorado, USA, June 4-5, 2015*, pages 451–463. The Association for Computer Linguistics. 101
- Kenneth J Rothman and Sander Greenland. 2005. Causation and causal inference in epidemiology. *American journal of public health*, 95(S1):S144–S150. 33
- Donald B Rubin. 1974. Estimating causal effects of treatments in randomized and non-randomized studies. *Journal of educational Psychology*, 66(5):688. 108, 172
- Donald B. Rubin. 1980. Randomization analysis of experimental data: The fisher randomization test comment. *Journal of the American Statistical Association*, 75(371):591–593. 7
- Donald B Rubin. 2005. Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 100(469):322–331. 173
- Jürgen Rudolph, Samson Tan, and Shannon Tan. 2023. Chatgpt: Bullshit spewer or the end of traditional assessments in higher education? *Journal of Applied Learning and Teaching*, 6(1). 32
- Mukund Rungta, Janvijay Singh, Saif M. Mohammad, and Diyi Yang. 2022. Geographic citation gaps in NLP research. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1371–1383, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. 113
- Bertrand Russell. 2004. *History of western philosophy*. Routledge. 7
- Mrinmaya Sachan, Kumar Dubey, and Eric Xing. 2017. From textbooks to knowledge: A case study in harvesting axiomatic knowledge from textbooks to solve geometry problems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 773–784. 50
- Mrinmaya Sachan, Kumar Avinava Dubey, Tom M Mitchell, Dan Roth, and Eric P Xing. 2018. Learning pipelines with limited data and domain knowledge: A study in parsing physics problems. *Advances in Neural Information Processing Systems*, 31. 50
- Mrinmaya Sachan and Eric Xing. 2017. Learning to solve geometry problems from natural language demonstrations in textbooks. In *Proceedings of the 6th Joint Conference on Lexical and Computational Semantics (* SEM 2017)*, pages 251–261. 50
- Peter Salovey and John D Mayer. 2004. *Emotional intelligence*. Dude publishing. 83, 93
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: Smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108. 14, 57
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. 2019a. ATOMIC: an atlas of machine commonsense for if-then reasoning. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in*

- Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 3027–3035. AAAI Press. 8, 18, 32, 81, 150
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019b. Social iqa: Commonsense reasoning about social interactions. In *EMNLP 2019*. 8, 18, 32, 150
- Ryoma Sato, Makoto Yamada, and Hisashi Kashima. 2022. Twin papers: A simple framework of causal inference for citations via coupling. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*, pages 4444–4448. ACM. 118
- Ajay B Satpute, Jocelyn Shu, Jochen Weber, Mathieu Roy, and Kevin N Ochsner. 2013. The functional neural architecture of self-reports of affective experience. *Biological psychiatry*, 73(7):631–638. 83, 94
- Timo Schick and Hinrich Schütze. 2022. True few-shot learning with Prompts—A real-world perspective. *Transactions of the Association for Computational Linguistics*, 10:716–731. 94
- Bruce Schneier. 1996. Applied cryptography. *John Willey and Sons Inc.*, 77
- Bernhard Schölkopf. 2022. Causality for machine learning. In Hector Geffner, Rina Dechter, and Joseph Y. Halpern, editors, *Probabilistic and Causal Inference: The Works of Judea Pearl*, volume 36 of *ACM Books*, pages 765–804. ACM. 84, 86
- Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun Zhang, and Joris M. Mooij. 2012. On causal and anticausal learning. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress. 70, 71, 72, 73, 77, 86
- Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. 2021. Towards causal representation learning. *CoRR*, abs/2102.11107. 84
- Lukas Schott, Jonas Rauber, Matthias Bethge, and Wieland Brendel. 2018. Towards the first adversarially robust neural network model on mnist. *arXiv preprint arXiv:1805.09190*. 86
- Rico Sennrich and Biao Zhang. 2019. Revisiting low-resource neural machine translation: A case study. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 211–221, Florence, Italy. Association for Computational Linguistics. 81
- Minjoon Seo, Hannaneh Hajishirzi, Ali Farhadi, Oren Etzioni, and Clint Malcolm. 2015. Solving geometry problems: Combining text and diagram interpretation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1466–1476, Lisbon, Portugal. Association for Computational Linguistics. 50
- Eleni Sgouritsa, Dominik Janzing, Philipp Hennig, and Bernhard Schölkopf. 2015. Inference of cause and effect with unsupervised inverse regression. In *Artificial intelligence and statistics*, volume 38 of *JMLR Workshop and Conference Proceedings*, pages 847–855. PMLR, JMLR.org. 72

- Nihar B Shah. 2022. An overview of challenges, experiments, and computational solutions in peer review. *Communications of the ACM*, 65(6):76–87. 113
- Naji Shajarisales, Dominik Janzing, Bernhard Schölkopf, and Michel Besserve. 2015. Telling cause from effect in deterministic linear dynamical systems. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 285–294. JMLR.org. 94
- Jiajun Shen, Peng-Jen Chen, Matthew Le, Junxian He, Jiatao Gu, Myle Ott, Michael Auli, and Marc’Aurelio Ranzato. 2021a. The source-target domain mismatch problem in machine translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1519–1533, Online. Association for Computational Linguistics. 69, 80
- Jianhao Shen, Yichun Yin, Lin Li, Lifeng Shang, Xin Jiang, Ming Zhang, and Qun Liu. 2021b. Generate & rank: A multi-task framework for math word problems. *arXiv preprint arXiv:2109.03034*. 62
- Shohei Shimizu, Patrik O. Hoyer, Aapo Hyvärinen, and Antti J. Kerminen. 2006. A linear non-gaussian acyclic model for causal discovery. *J. Mach. Learn. Res.*, 7:2003–2030. 10
- Hidetoshi Shimodaira. 2000. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of statistical planning and inference*, 90(2):227–244. 73
- Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4222–4235, Online. Association for Computational Linguistics. 21
- Sam Shleifer and Alexander M. Rush. 2020. Pre-trained summarization distillation. *CoRR*, abs/2010.13002. 14
- Ilya Shpitser and Judea Pearl. 2006a. Identification of conditional interventional distributions. In *22nd Conference on Uncertainty in Artificial Intelligence, UAI 2006*, pages 437–444. 23, 24
- Ilya Shpitser and Judea Pearl. 2006b. Identification of joint interventional distributions in recursive semi-markovian causal models. In *Proceedings, The Twenty-First National Conference on Artificial Intelligence and the Eighteenth Innovative Applications of Artificial Intelligence Conference, July 16-20, 2006, Boston, Massachusetts, USA*, pages 1219–1226. AAAI Press. 141
- Mrityunjay Singh, Amit Kumar Jakhar, and Shivam Pandey. 2021a. Sentiment analysis on the impact of coronavirus in social life using the bert model. *Social Network Analysis and Mining*, 11(1):1–11. 102
- Shikhar Singh, Nuan Wen, Yu Hou, Pegah Alipoormolabashi, Te-lin Wu, Xuezhe Ma, and Nanyun Peng. 2021b. COM2SENSE: A commonsense reasoning benchmark with complementary sentences. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 883–898, Online. Association for Computational Linguistics. 150

- Surya Nath Singh. 2014. Sampling techniques & determination of sample size in applied statistics research : an overview. 121
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Kumar Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Schärli, Aakanksha Chowdhery, Philip Andrew Mansfield, Blaise Agüera y Arcas, Dale R. Webster, Gregory S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle K. Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. 2022. Large language models encode clinical knowledge. *CoRR*, abs/2212.13138. 32
- Abir Smiti. 2020. A critical overview of outlier detection methods. *Computer Science Review*, 38:100306. 126, 177
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642. 93
- Anders Søgaard. 2013. Semi-supervised learning and domain adaptation in natural language processing. *Synthesis Lectures on Human Language Technologies*, 6(2):1–103. 80, 81, 160
- Anders Søgaard and Christian Rishøj. 2010. Semi-supervised dependency parsing using generalized tri-training. In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pages 1065–1073, Beijing, China. Coling 2010 Organizing Committee. 82
- Ray J Solomonoff. 1964. A formal theory of inductive inference. part ii. *Information and control*, 7(2):224–254. 74
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. MpNet: Masked and permuted pre-training for language understanding. *arXiv preprint arXiv:2004.09297*. 121, 122, 175
- Peter Spirtes, Clark Glymour, and Richard Scheines. 1993. Causation, prediction, and search. 7
- Peter Spirtes, Clark Glymour, and Richard Scheines. 2000. *Causation, Prediction, and Search, Second Edition*. Adaptive computation and machine learning. MIT Press. 7
- Peter Spirtes, Clark Glymour, and Richard Scheines. 2001. *Causation, prediction, and search*. MIT press. 7, 8, 10, 19, 21, 25
- Peter Spirtes and Kun Zhang. 2016. Causal discovery and inference: Concepts and recent methodological advances. In *Applied informatics*, volume 3, pages 1–28. SpringerOpen. 8, 10, 87
- Daniel Spokoyny, Ivan Lee, Zhao Jin, and Taylor Berg-Kirkpatrick. 2021. Masked measurement prediction: Learning to jointly predict quantities and units from textual context. *arXiv preprint arXiv:2112.08616*. 62

- Kathrin Spreyer and Jonas Kuhn. 2009. Data-driven dependency parsing of new languages using incomplete and noisy training data. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning (CoNLL-2009)*, pages 12–20, Boulder, Colorado. Association for Computational Linguistics. 82
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew K. Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakas, and et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *CoRR*, abs/2206.04615. 32
- Manfred Stede. 2008. Connective-based local coherence analysis: A lexicon for recognizing causal relationships. In *Semantics in Text Processing. STEP 2008 Conference Proceedings*, pages 221–237. College Publications. 32
- Mark Steedman, Miles Osborne, Anoop Sarkar, Stephen Clark, Rebecca Hwa, Julia Hockenmaier, Paul Ruhlen, Steven Baker, and Jeremiah Crim. 2003. Bootstrapping statistical parsers from small datasets. In *10th Conference of the European Chapter of the Association for Computational Linguistics*, Budapest, Hungary. Association for Computational Linguistics. 81
- James A Stimson, Michael B MacKuen, and Robert S Erikson. 1995. Dynamic representation. *American political science review*, pages 543–565. 98, 100
- Alessandro Stolfo, Zhijing Jin, Kumar Shridhar, Bernhard Schölkopf, and Mrinmaya Sachan. 2023. A causal framework to quantify the robustness of mathematical reasoning with language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Toronto, Canada. Association for Computational Linguistics. 25, 29
- Masashi Sugiyama and Motoaki Kawanabe. 2012. *Machine learning in non-stationary environments: Introduction to covariate shift adaptation*. MIT press. 73
- Songbo Tan, Xueqi Cheng, Yuefen Wang, and Hongbo Xu. 2009. Adapting naive bayes to domain adaptation for sentiment analysis. In *Advances in Information Retrieval, 31th European Conference on IR Research, ECIR 2009, Toulouse, France, April 6-9, 2009. Proceedings*, volume 5478 of *Lecture Notes in Computer Science*, pages 337–349. Springer. 93
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023a. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca. 14, 30, 51, 61
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023b. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca. 91

- Robert F Tate. 1954. Correlation between a discrete and a continuous variable. point-biserial correlation. *The Annals of mathematical statistics*, 25(3):603–607. 124
- Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najaoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das, and Ellie Pavlick. 2019. What do you learn from context? probing for sentence structure in contextualized word representations. *ArXiv*, abs/1905.06316. 38
- Avijit Thawani, Jay Pujara, Filip Ilievski, and Pedro Szekely. 2021. Representing numbers in NLP: a survey and a vision. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–656, Online. Association for Computational Linguistics. 62
- Mike Thelwall, Kevan Buckley, and Georgios Paltoglou. 2011. Sentiment in twitter events. *J. Assoc. Inf. Sci. Technol.*, 62(2):406–418. 101
- Jörg Tiedemann and Lars Nygaard. 2004. The OPUS corpus - parallel and free: <http://logos.uio.no/opus>. In *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)*, Lisbon, Portugal. European Language Resources Association (ELRA). 76, 160
- Erik Tjong Kim Sang and Johan Bos. 2012. Predicting the 2011 Dutch senate election results with Twitter. In *Proceedings of the Workshop on Semantic Analysis in Social Media*, pages 53–60, Avignon, France. Association for Computational Linguistics. 101
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971. 14, 30, 36, 51, 60, 91
- Ciprian-Octavian Truică, Elena-Simona Apostol, Maria-Luiza Șerban, and Adrian Paschke. 2021. Topic-based document-level sentiment analysis using contextual cues. *Mathematics*, 9(21). 83
- Ruibo Tu, Chao Ma, and Cheng Zhang. 2023. Causal-discovery performance of chatgpt in the context of neuropathic pain diagnosis. *arXiv preprint arXiv:2301.13819*. 8
- Peter D. Turney. 2002. Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*, pages 417–424. ACL. 83, 93
- Amos Tversky and Daniel Kahneman. 1974. Judgment under uncertainty: Heuristics and biases: Biases in judgments reveal some heuristics of thinking under uncertainty. *Science*, 185(4157):1124–1131. 88
- Muhammad Umer, Saima Sadiq, Malik Muhammad Saad Missen, Zahid Hameed, Zahid Aslam, Muhammad Abubakar Siddique, and Michele Nappi. 2021. Scientific papers citation analysis using textual features and smote resampling techniques. *Pattern Recognition Letters*, 150:250–257. 127

- Marco Valenzuela, Vu Ha, and Oren Etzioni. 2015. Identifying meaningful citations. In *Scholarly Big Data: AI Perspectives, Challenges, and Ideas, Papers from the 2015 AAAI Workshop, Austin, Texas, USA, January, 2015*, volume WS-15-13 of *AAAI Technical Report*. AAAI Press. 113, 122, 123, 127, 176
- Lucrezia Valeriani, Diego Doimo, Francesca Cuturello, Alessandro Laio, Alessio Ansuini, and Alberto Cazzaniga. 2024. The geometry of hidden representations of large transformer models. *Advances in Neural Information Processing Systems*, 36. 38
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008. 76, 125, 160, 176
- Victor Veitch, Dhanya Sridhar, and David M. Blei. 2020. Adapting text embeddings for causal inference. In *Proceedings of the Thirty-Sixth Conference on Uncertainty in Artificial Intelligence, UAI 2020, virtual online, August 3-6, 2020*, volume 124 of *Proceedings of Machine Learning Research*, pages 919–928. AUAI Press. 62, 81
- Pranav Venkit, Mukund Srinath, Sanjana Gautam, Saranya Venkatraman, Vipul Gupta, Rebecca Passonneau, and Shomir Wilson. 2023. The sentiment problem: A critical survey towards deconstructing sentiment analysis. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13743–13763, Singapore. Association for Computational Linguistics. 83
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Simas Sakenis, Jason Huang, Yaron Singer, and Stuart Shieber. 2020a. Causal mediation analysis for interpreting neural NLP: The case of gender bias. *arXiv preprint arXiv:2004.12265*. 92
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart M. Shieber. 2020b. Investigating gender bias in language models using causal mediation analysis. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. 51, 62
- Tanmay Vijay, Ayan Chawla, Balan Dhanka, and Purnendu Karmakar. 2020. Sentiment analysis on covid-19 twitter data. In *2020 5th IEEE International Conference on Recent Advances and Innovations in Engineering (ICRAIE)*, pages 1–7. IEEE. 102
- Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. 2008. Extracting and composing robust features with denoising autoencoders. In *Machine Learning, Proceedings of the Twenty-Fifth International Conference (ICML 2008), Helsinki, Finland, June 5-9, 2008*, volume 307 of *ACM International Conference Proceeding Series*, pages 1096–1103. ACM. 78
- Elena Voita and Ivan Titov. 2020. Information-theoretic probing with minimum description length. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 183–196. Association for Computational Linguistics. 75, 81

- Julius von Kügelgen, Alexander Mey, and Marco Loog. 2019. Semi-generative modelling: Covariate-shift adaptation with cause and effect features. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1361–1369. PMLR. 72, 81
- Julius von Kügelgen, Alexander Mey, Marco Loog, and Bernhard Schölkopf. 2020. Semi-supervised learning, causality, and the conditional cluster assumption. In *Conference on Uncertainty in Artificial Intelligence*, pages 1–10. PMLR. 72, 81
- Eric Wallace, Yizhong Wang, Sujian Li, Sameer Singh, and Matt Gardner. 2019. Do NLP models know numbers? probing numeracy in embeddings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5307–5315, Hong Kong, China. Association for Computational Linguistics. 62
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 3261–3275. 32
- Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 billion parameter autoregressive language model. 58
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. 37, 38, 40, 41
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022. Super-naturalinstructions: Generalization via declarative instructions on 1600+ NLP tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 5085–5109. Association for Computational Linguistics. 32
- Jason Wei, Najoung Kim, Yi Tay, and Quoc V. Le. 2023. Inverse scaling can become u-shaped. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 15580–15591. Association for Computational Linguistics. 48
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022a. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*. 50, 60
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022b. Chain of thought prompting elicits reasoning

- in large language models. In *Advances in Neural Information Processing Systems*. 22, 28, 62
- Bernard L Welch. 1947. The generalization of ‘student’s’ problem when several different population variances are involved. *Biometrika*, 34(1/2):28–35. 79
- Sean Welleck, Peter West, Jize Cao, and Yejin Choi. 2022. Symbolic brittleness in sequence models: on systematic generalization in symbolic mathematics. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 8629–8637. 50
- Cynthia G Whitney, Fangjun Zhou, James Singleton, and Anne Schuchat. 2014. Benefits from immunization during the vaccines for children program era—united states, 1994–2013. *Morbidity and Mortality Weekly Report*, 63(16):352. 20
- Janyce M Wiebe. 1994. Tracking point of view in narrative. *Computational Linguistics*, 20(2):233–287. 83, 93
- Moritz Willig, Matej Zečević, Devendra Singh Dhami, and Kristian Kersting. 2023. Probing for correlations of causal facts: Large language models and causality. 18
- James Wilsdon. 2016. The metric tide: Independent review of the role of metrics in research assessment and management. 113
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics. 14, 41, 58, 90, 102, 122, 132, 160
- Marta Natalia Wróblewska. 2021. Research impact evaluation and academic discourse. *Humanities and Social Sciences Communications*, 8(1):1–12. 127
- Yuxi Xie, Guanzhen Li, and Min-Yen Kan. 2023. Echo: Event causality inference via human-centric reasoning. *arXiv preprint arXiv:2305.14740*. 8
- Xiaoyu Xing, Zhijing Jin, Di Jin, Bingning Wang, Qi Zhang, and Xuanjing Huang. 2020. Tasty burgers, soggy fries: Probing aspect robustness in aspect-based sentiment analysis. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3594–3605, Online. Association for Computational Linguistics. 93, 95
- Jinghang Xu, Wanli Zuo, Shining Liang, and Xianglin Zuo. 2020. A review of dataset and labeling methods for causality extraction. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 1519–1531, Barcelona, Spain (Online). International Committee on Computational Linguistics. 32
- Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. 2023. Harnessing the power of llms in practice: A survey on chatgpt and beyond. *ArXiv*, abs/2304.13712. 122

- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *CoRR*, abs/1906.08237. 90, 93
- Bei Yu, Yingya Li, and Jun Wang. 2019. Detecting causal language use in science findings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4664–4674, Hong Kong, China. Association for Computational Linguistics. 32
- Qinan Yu, Jack Merullo, and Ellie Pavlick. 2023. Characterizing mechanisms for factual recall in language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9924–9959, Singapore. Association for Computational Linguistics. 38, 40, 47
- Nurulhuda Zainuddin and Ali Selamat. 2014. Sentiment analysis using support vector machine. *2014 International Conference on Computer, Communications, and Control Technology (I4CT)*, pages 333–337. 176
- Robert B Zajonc. 1980. Feeling and thinking: Preferences need no inferences. *American psychologist*, 35(2):151. 87
- Matej Zečević, Moritz Willig, Devendra Singh Dhami, and Kristian Kersting. 2023. Causal parrots: Large language models may talk causality but are not causal. *Transactions on Machine Learning Research*. 7, 20, 21, 32
- Cheng Zhang, Stefan Bauer, Paul Bennett, Jiangfeng Gao, Wenbo Gong, Agrin Hilmkil, Joel Jennings, Chao Ma, Tom Minka, Nick Pawlowski, et al. 2023. Understanding causality with large language models: Feasibility and opportunities. *arXiv preprint arXiv:2304.05524*. 20, 32
- Jiajun Zhang and Chengqing Zong. 2016. Exploiting source-side monolingual data in neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1535–1545, Austin, Texas. Association for Computational Linguistics. 79
- Kun Zhang, Mingming Gong, and Bernhard Schölkopf. 2015a. Multi-source domain adaptation: A causal view. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA*, pages 3150–3157. AAAI Press. 72
- Kun Zhang and Aapo Hyvärinen. 2009. Causality discovery with additive disturbances: An information-theoretical perspective. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2009, Bled, Slovenia, September 7-11, 2009, Proceedings, Part II* 20, pages 570–585. Springer. 7, 10, 87, 94
- Kun Zhang, Bernhard Schölkopf, Krikamol Muandet, and Zhikun Wang. 2013. Domain adaptation under target and conditional shift. In *International Conference on Machine Learning*, pages 819–827. PMLR. 72
- Li Zhang, Qing Lyu, and Chris Callison-Burch. 2020a. Reasoning about goals, steps, and temporal ordering with WikiHow. In *Proceedings of the 2020 Conference on Empirical*

- Methods in Natural Language Processing (EMNLP)*, pages 4630–4639, Online. Association for Computational Linguistics. 18, 32, 150
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. OPT: open pre-trained transformer language models. *CoRR*, abs/2205.01068. 1, 7, 20, 32
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015b. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28. 83, 85, 87, 89, 91, 93
- Xikun Zhang, Deepak Ramachandran, Ian Tenney, Yanai Elazar, and Dan Roth. 2020b. Do language embeddings capture scales? In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4889–4896, Online. Association for Computational Linguistics. 62
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2023. Large language models are human-level prompt engineers. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. 161
- Zhi-Hua Zhou and Ming Li. 2005. Tri-training: Exploiting unlabeled data using three classifiers. *IEEE Trans. Knowl. Data Eng.*, 17(11):1529–1541. 78
- Xiao-Dan Zhu, Peter D. Turney, Daniel Lemire, and André Vellino. 2015. Measuring academic influence: Not all citations are equal. *Journal of the Association for Information Science and Technology*, 66. 113, 115, 123
- Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul F. Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. *CoRR*, abs/1909.08593. 29
- Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2024. Can large language models transform computational social science? *Comput. Linguistics*, 50(1):237–291. 29, 32
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. 2023. Representation engineering: A top-down approach to AI transparency. *CoRR*, abs/2310.01405. 1