

Marrying LLMs with Domain Expert Validation for Causal Graph Generation

Alessandro Castelnovo^{1,2}, Riccardo Crupi¹, Fabio Mercorio^{3,4}, Mario Mezzanzanica^{3,4}, Daniele Poterti³ and Daniele Regoli¹

¹*Data Science and Artificial Intelligence, Intesa Sanpaolo S.p.A., Italy*

²*Dept. of Informatics, Systems and Communication, Univ. of Milan-Bicocca, Italy*

³*Dept. of Statistics and Quantitative Methods, Univ. of Milan-Bicocca, Italy*

⁴*CRISP Research Centre crispresearch.eu, University of Milano Bicocca, Italy*

Abstract

In the era of rapid growth and transformation driven by artificial intelligence across various sectors, which is catalyzing the fourth industrial revolution, this research is directed toward harnessing its potential to enhance the efficiency of decision-making processes within organizations. When constructing machine learning-based decision models, a fundamental step involves the conversion of domain knowledge into causal-effect relationships that are represented in causal graphs. This process is also notably advantageous for constructing explanation models. We present a method for generating causal graphs that integrates the strengths of Large Language Models (LLMs) with traditional causal theory algorithms. Our method seeks to bridge the gap between AI's theoretical potential and practical applications. In contrast to recent related works that seek to exclude the involvement of domain experts, our method places them at the forefront of the process. We present a novel pipeline that streamlines and enhances domain-expert validation by providing robust causal graph proposals. These proposals are enriched with transparent reports that blend foundational causal theory reasoning with explanations from LLMs.

Keywords

Causal Discovery, LLMs, Human-AI-Interaction

1. Introduction

In the modern era, the realms of Causality and Artificial Intelligence (AI) are converging to foster a significant transformation in various sectors, including business and industry. Understanding causality, a practice deeply rooted in causal inference and causal discovery, has become pivotal in enhancing decision-making processes and fostering innovation in a data-driven society. Causal inference, the practice of determining the cause-and-effect relationship between variables [1], and causal discovery, the identification of these relationships from observational data [2], have taken center stage.

In recent years, the push towards explainable AI has garnered significant momentum, with numerous government initiatives underscoring its importance. Notably, the General Data Protection Regulation (GDPR) in the European Union [3], the Defence Advanced Research Projects Agency (DARPA) XAI program in the United States [4], and the European Commission's proposal for legislation on AI systems (The European Commission 2021) have all been striving to

*3rd Italian Workshop on Artificial Intelligence and Applications for Business and Industries - AIABI
co-located with AI*IA 2023*



© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

foster the development of AI systems whose outputs are both comprehensible and trustworthy for end-users. According to [5], humans perceive an explanation as good when it possesses certain characteristics. These include being contrastive, concise, selected, social, and causal. In particular, to ensure the understanding of the underlying causal relationship between data, Direct Acyclic Graphs (DAGs) emerge as a fundamental tool in the evolution of explainable AI systems. These graphs, which are pivotal in modeling causal relationships and forming the structural basis of causal models [1], represent variables as nodes and causal relationships as edges. This structure not only eliminates feedback loops but also delineates causal effects with clarity, facilitating the transparent articulation of causal assumptions. Moreover, DAGs, or Structural Causal Models (SCMs), play a crucial role in identifying confounding variables, thereby aiding in the construction of accurate and interpretable models. Through these capabilities, DAGs have proven to be an indispensable foundation in the creation of AI-driven decision models, aligning perfectly with the aforementioned government initiatives by contributing to the development of explanations that can be easily understood and trusted by users [6]. Thus, the incorporation of SCMs in AI systems aligns well with ongoing governmental efforts to promote explainable AI, serving as a vital tool in the realization of this significant objective.

However, modeling causal knowledge is complex and challenging since it requires an actual understanding of the relations, beyond statistical correlation. Yet, a critical ingredient in this quest for causality is domain knowledge [1]. Nevertheless, in numerous real-world scenarios, the collaboration between non-technical domain experts and computer scientists in defining a causal graph can be a challenging and time-consuming process [7]. This is evident in the work of [6], in which is developed a method for counterfactual generation that respects the causal influence among variables. This application is essential in business and banking contexts, such as credit lending, where a client who has their loan denied has the right to know what actions they can take to get the loan approved. These actions must be feasible and take into account the causal relationships between variables (e.g., the type of job influences income, not the other way around).

LLMs, such as GPT-3.5 and GPT-4, show promise in bridging the gap between computational efficiency and domain expertise [8]. They excel in causal reasoning tasks, translating between natural language and formal methods, generating causal graphs, and identifying background causal contexts [9]. LLMs can act as proxies for human domain knowledge, enhancing causal analysis while reducing human effort [10]. Their prowess, mainly when used alongside existing causal statistical methods [11] holds the potential to revolutionize how we approach and understand causality. However, despite their benefits, LLMs may misinterpret complex causal models without deep domain expertise, leading to erroneous decisions [1]. Thus, deep domain knowledge remains crucial and indispensable, acting as a guiding force that enables AI technologies like LLMs to realize their full potential in causal discovery, thereby contributing to the ongoing transformative wave in various industries [12].

The primary objective of this work is to harness AI's transformative power in reshaping business and industry processes, particularly focusing on the integration of LLMs with causal discovery techniques. Automating the discovery of causal relationships, allows domain experts to concentrate solely on validating the proposed causal graph [13], thereby potentially achieving higher quality and greater efficiency, akin to the revolution AI has brought in other sectors. This endeavor seeks to significantly reduce the time and human resources traditionally required,

paving the way for a more innovative and efficient approach to business and industrial processes.

Contribution. We present a first attempt at formalizing a causal pipeline to design a causal graph requiring minimal domain knowledge, that seamlessly combines data-driven statistical causality techniques with the insights of LLMs, acting as proxies for domain expertise. Moreover, this comprehensive framework has been encapsulated into an open-source Python package — that will be available in open source — designed to ease its integration into real-world scenarios.

As a contribution, this framework aims to be domain-specific, statistically robust, transparent, and explainable to ensure trust and effective validation by the human-in-the-loop of the generated causal graph. More in particular:

- **domain-specific:** among the results generated by the framework is a set of probable DAGs that optimally depict the causal relationships derived from the given data. Notably, these graphs are tailored not just to the data, but also to the specific domain, thanks to the integration of the LLMs with the Causal Discovery theory.
- **statistically robust:** the framework includes a final statistical sensitivity assessment for each DAG. This assessment evaluates the causal effect of each edge in accordance with the literature on causal inference, ensuring the robustness of the results.
- **transparent and explainable:** for each DAG, a comprehensive report resulting from rigorous causal theory testing, along with explanations provided by LLMs, either reinforces or questions the presence of particular edges and directional relationships within the graph.

The process concludes with the domain expert making an informed decision to select the preferred causal graph from among the proposed options. In our implementation, we place significant emphasis on the aspect of prompt engineering within the overall processing of the proposed causal pipeline.

2. Related Works

In this section, we explore recent developments in the literature on LLMs that (i) focus on their capabilities in the realm of causal reasoning, and (ii) investigate strategies for integrating LLMs with causal discovery techniques.

Exploring the Causal Reasoning Abilities of LLMs. Recent advancements indicate that models like ChatGPT are progressing towards artificial general intelligence (AGI), notably enhancing causal reasoning and high-precision tasks [14]. There is the potential for a paradigm shift in machine learning that could harmonize the strengths of AI with human capabilities, leading to transformative solutions and enhancing human decision-making [15]. Kıcıman et al. [8] further explore LLMs capabilities in causal reasoning, illustrating their prowess in tasks like code generation and complex reasoning. These models demonstrate high performance in causal discovery, achieving up to 97% accuracy on the Tübingen benchmark and showcasing versatility across different domains. Despite this, they can occasionally falter in basic logic tasks, raising reliability concerns. Kıcıman et al. [8] emphasize the importance of incorporating human

domain knowledge into causal analysis, suggesting that LLMs could serve as powerful tools to enhance this process through dynamic conversational interfaces. Yet, it is vital to remain cautiously optimistic about their capabilities due to potential erratic performances. Future research is poised to delve deeper into the capacities and boundaries of LLMs in this domain.

Integrating LLMs in Causal Discovery. Long et al. [16] introduced a novel approach to the challenges of causal discovery by formalizing the use of imperfect experts as an optimization problem, aiming to minimize the Markov equivalence class (MEC) size while ensuring the true graph remains included. They proposed a greedy approach reliant on Bayesian inference to achieve this, incrementally integrating expert knowledge. Empirical evaluations on real data revealed the effectiveness of their method, especially when the expert consistently provided correct orientations. However, when using LLMs as the experts, the results were mixed, suggesting both the potential and the challenges of integrating LLMs into causal discovery. Ban et al. [17] explore the role of LLMs in Causal Structure Learning (CSL), focusing on utilizing LLMs to pinpoint direct causal relations in observed data. This two-stage framework first uses LLMs to identify potential causal connections based on textual data and then applies these insights as constraints in data-driven CSL algorithms. The aim is to merge the intuitive causal understanding of LLMs with the detailed causal analysis found in CSL, potentially increasing its efficiency and accuracy. However, the study acknowledges the potential errors in the causal statements generated by GPT-4, indicating room for future improvements in both understanding and quality of causal relationships.

In contrast to existing approaches, our work distinguishes itself by prioritizing human validation. Unlike conventional methods that aim to autonomously generate causal graphs, we actively engage domain experts. Our pipeline simplifies their tasks with robust causal graph proposals, accompanied by a transparent report featuring both causal theory reasoning and LLM-based explanations.

3. The causal pipeline with LLM and *Human-in-the loop*

Our proposed framework endeavors to deliver DAGs tailored to specific domains, leveraging an explanation-centric approach and undergoing rigorous statistical testing. This framework progresses through three essential phases: Causal Discovery, LLM-based Causal Elaboration, and Causal Inference, as illustrated in Figure 1. Notably, human involvement remains a crucial element throughout this process. This framework is designed to guide users in selecting the appropriate Causal Graph, drawing from the Set of Directed Causal Graphs, Causal Relationship Explanations, and Causal Sensitivity Analysis.

Causal Discovery Step. The causal discovery step begins with a dataset containing the variables of interest and their observed interactions. We process this dataset using four distinct causal discovery algorithms: PC (a constraint-based method) [2], GES (a score-based method) [18], FCI (an extension of the PC algorithm) [2], and LiNGAM (a functional causal model) [19]. Each of these algorithms represents a distinct method within causal discovery.

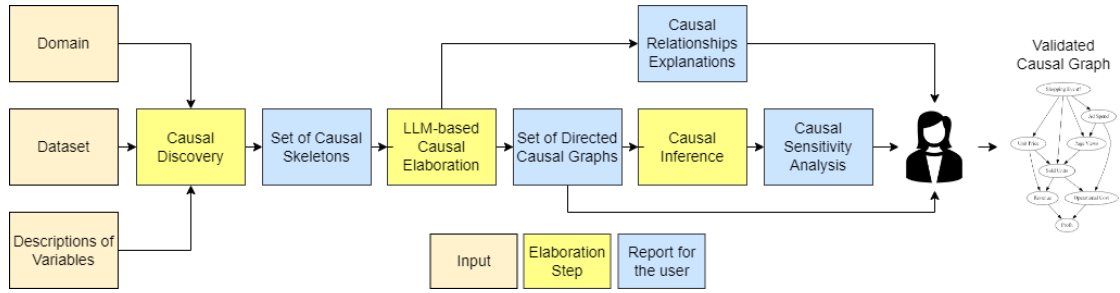


Figure 1: The proposed pipeline.

The output of this process is a set of four Causal Skeletons, each one derived from a causal discovery algorithm. Each skeleton will be further processed by the LLM in the LLM-based Causal Elaboration Step.

LLM-based Causal Elaboration Step. Once we retrieve the sets of Causal Skeletons from the previous step it is time to determine the direction of the relationships performing the task of Pairwise Causal Discovery [20]. In this task, the goal is to determine the causal relationship between two variables. In this stage of the framework, we undertake two primary tasks: initially, we delve into understanding the inherent nature of the relationships depicted in each Causal Skeleton, and subsequently, we explore potential additional relationships through conditional independence tests. One of the fundamental challenges arises from the presumption that the LLM can infer the domain context purely based on variable values¹. In this stage, the success of the task hinges largely on the quality of the prompt. We are in the process of creating a novel prompt that aligns with causal theory and maintains consistent performance across various LLMs, including GPT-3.5, GPT-4, and LLaMA. The outcome of this stage is a set of DAG: a directional graph where the sequence of cause-and-effect relations is so structured that it never loops back on itself. Furthermore, in this stage, we generate explanations using carefully crafted prompts to justify the outputs of the LLMs.

Causal Inference Step. The last phase of our pipeline is designed for the statistical validation of the obtained DAGs from the preceding step. In pursuit of this objective, we adhere to the four key steps of causal inference as outlined in [1]. To facilitate this process, we leverage the capabilities of DoWhy, a widely recognized open-source Python library. What distinguishes DoWhy is its strong foundation in causal assumptions, firmly rooted in the well-established framework of causal graphs [21].

At the conclusion of the pipeline, we generate a transparent report for the user. This report includes all the proposed DAGs and the rationale behind their generation. It encompasses the results of causal discovery tests, LLMs prompt outcomes, and causal inference, empowering users to make an informed choice regarding their preferred causal graph.

¹For the LLM to make meaningful inferences and not merely fabricate a domain, the expert needs to provide at least some foundational information about the domain.

4. Conclusion and Next Step

We presented a novel pipeline for constructing a causal graph that effectively combines well-established causal theory algorithms from the literature with LLMs. While related works focus on generating a final causal graph without the involvement of domain experts, our approach is distinct as it places human validation at the core of the process, thereby aligning with the broader trend where AI aids decision-makers in organizations, enhancing innovation processes and managerial tasks. Our pipeline is designed to streamline and improve the expert's work by providing robust causal graph proposals. It accompanies these proposals with a transparent report that includes the underlying causal theory reasoning and explanations derived from LLMs. This approach not only sets us apart from existing methodologies but also aligns with the ongoing revolution, where AI is a pivotal tool in enhancing business decisions. Our next steps involve refining our prompts, open-source our code, and enhancing clarity in our LLM-based Causal Elaboration step through detailed insights and prompt examples. Additionally, we will further explore the limitations of LLMs, especially in sophisticated domains necessitating profound expertise, it is also acknowledged that our method could derive further validation and credibility through a more detailed evaluation segment. This includes highlighting the pipeline's efficacy through application on real-world datasets and engaging in a comparative analysis with alternate methodologies.

References

- [1] J. Pearl, et al., *Models, reasoning and inference*, Cambridge, UK: CambridgeUniversityPress 19 (2000) 3.
- [2] P. Spirtes, C. N. Glymour, R. Scheines, *Causation, prediction, and search*, MIT press, 2000.
- [3] The European Union, *EU General Data Protection Regulation (GDPR): Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation)*, Official Journal of the European Union, 2016. <http://data.europa.eu/eli/reg/2016/679/2016-05-04>.
- [4] D. Gunning, D. Aha, *Darpa's explainable artificial intelligence (xai) program*, *AI magazine* 40 (2019) 44–58.
- [5] T. Miller, *Explanation in artificial intelligence: Insights from the social sciences*, *Artificial intelligence* 267 (2019) 1–38.
- [6] R. Crupi, A. Castelnovo, D. Regoli, B. San Miguel Gonzalez, *Counterfactual explanations as interventions in latent space*, *Data Mining and Knowledge Discovery* (2022) 1–37.
- [7] X. Xie, F. Du, Y. Wu, *A visual analytics approach for exploratory causal analysis: Exploration, validation, and applications*, *IEEE Transactions on Visualization and Computer Graphics* 27 (2020) 1448–1458.
- [8] E. Kıcıman, R. Ness, A. Sharma, C. Tan, *Causal reasoning and large language models: Opening a new frontier for causality*, *arXiv preprint arXiv:2305.00050* (2023).
- [9] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee,

- Y. Li, S. Lundberg, Sparks of artificial general intelligence: Early experiments with gpt-4, preprint arXiv:2303.12712 (2023).
- [10] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, Language models are unsupervised multitask learners, OpenAI blog (2019).
 - [11] M. Hernán, J. Robins, Causal Inference: What If, Chapman & Hall/CRC, Boca Raton, 2020.
 - [12] D. Danks, Unifying the mind: Cognitive representations as graphical models, Mit Press, 2014.
 - [13] B. Youngmann, M. Cafarella, B. Salimi, A. Zeng, Causal data integration, arXiv preprint arXiv:2305.08741 (2023).
 - [14] Z. C. et al., Understanding causality with large language models: Feasibility and opportunities, arXiv preprint arXiv:2304.05524 (2023).
 - [15] C. T. Wolf, Reprogramming the american dream: from rural america to silicon valley—making ai serve us all by kevin scott and greg shaw, Information & Culture 56 (2021) 113–114.
 - [16] S. Long, A. Piché, V. Zantedeschi, T. Schuster, A. Drouin, Causal discovery with language models as imperfect experts, arXiv preprint arXiv:2307.02390 (2023).
 - [17] T. Ban, L. Chen, X. Wang, H. Chen, From query tools to causal architects: Harnessing large language models for advanced causal discovery from data, arXiv preprint arXiv:2306.16902 (2023).
 - [18] D. M. Chickering, Optimal structure identification with greedy search, Journal of machine learning research 3 (2002) 507–554.
 - [19] S. Shimizu, P. O. Hoyer, A. Hyvärinen, A. Kerminen, M. Jordan, A linear non-gaussian acyclic model for causal discovery., Journal of Machine Learning Research 7 (2006).
 - [20] J. M. Mooij, J. Peters, D. Janzing, J. Zscheischler, B. Schölkopf, Distinguishing cause from effect using observational data: methods and benchmarks, The Journal of Machine Learning Research 17 (2016) 1103–1204.
 - [21] A. Sharma, E. Kiciman, Dowhy: An end-to-end library for causal inference, arXiv preprint arXiv:2011.04216 (2020).