**SUPSI**

# Causal Graph Identification by Large Language Models

Studente/i

**Piqué Gregorio**

Relatore

**Antonucci Alessandro**

Correlatore

**Zaffalon Marco**

Committente

Corso di laurea

**Ingegneria informatica**

Codice progetto

**C10681**

Anno

**2023**

Data

**August 31, 2023**

STUDENTSUPSI

# Contents

## Abstract

Advances in causal inference are vital across multiple fields and contexts. A correct and complete understanding of the causal relationships behind the system of interest is a fundamental requirement for making accurate decisions. Several methods and techniques can be used to identify causal relationships with the task called causal discovery, but many of these approaches present different flaws and weaknesses. *Large Language Models* (LLMs) can be used as a new assistant to aid human efforts and contributing to the task of causal analysis. This project aims to evaluate the ability of LLMs in identifying causal relationships and causal graphs from natural language texts. Specifically, the LLM used in the project is the *Generative Pre-trained Transformer* (GPT) model, which is one of the most important and advanced LLMs recently developed. We will also present a set of techniques that can be used to improve the accuracy of LLM results. The project required the implementation of a software infrastructure to collect textual data and interact with the GPT API to process it to conduct causal discovery. The results showed that LLMs provide non-trivial contribution in helping with the identification of causal graphs. The application of LLMs to tasks of this nature is still in its early stages and has some limitations, but it has achieved some promising results and revealed new opportunities.

## Riassunto

I progressi nell'inferenza causale sono fondamentali in diversi campi e contesti. Una comprensione corretta e completa delle relazioni causali alla base del sistema di interesse è un requisito fondamentale per prendere decisioni accurate. Per identificare le relazioni causali si può usare l'operazione chiamata "scoperta causale", applicando diversi metodi e tecniche, ma molti di questi approcci presentano vari difetti e debolezze. I modelli linguistici di grandi dimensioni o *Large Language Models* (LLM) possono essere utilizzati come un nuovo assistente per aiutare gli sforzi umani e contribuire al compito dell'analisi causale. Questo progetto mira a valutare la capacità degli LLM di identificare relazioni e grafi causali da testi in linguaggio naturale. Nello specifico, il LLM usato in questo progetto è quello *Generative Pre-trained Transformer* (GPT), uno dei modelli più importanti e avanzati sviluppati di recente. Verrà inoltre presentata una serie di tecniche che possono essere utilizzate per migliorare l'accuratezza dei risultati degli LLM. Il progetto ha richiesto l'implementazione di un'infrastruttura software per raccogliere dati testuali e interagire con l'API GPT per elaborarli e condurre l'operazione di scoperta causale. I risultati hanno dimostrato che gli LLM forniscono un contributo non banale nell'identificazione di grafi causali. L'applicazione degli LLM a compiti di questa natura è ancora agli inizi e presenta varie limitazioni, ma ha ottenuto risultati promettenti e ha rivelato nuove opportunità.

## Progetto assegnato

### Descrizione

Il progetto consiste primariamente in uno studio empirico delle possibilità di apprendere grafi causali mediante query su LLMs (large language models). Il lavoro richiede lo sviluppo dell'infrastruttura software (codice Python che interagisce con API) per l'apprendimento del grafo da testi in linguaggio naturale. Una serie di benchmark esistenti verranno usate per la validazione e per il confronto contro altre tecniche allo stato dell'arte.

### Compiti

Sviluppo codice. Benchmarking. Valutazione risultati e sintesi.

### Obbiettivi

Confronto con lo stato dell'arte.

### Tecnologie

Python + API GPT.

# Chapter 1

# Introduction

Advances in causal inference are crucial for many fields and contexts. In the realm of artificial intelligence systems, a key challenge lies in their predominant reliance on statistical approaches and lack the ability to reason. The need for trustworthy machine learning tools has led to a growing interest in causality as a potential solution to this problem. Causality is the study of cause and effect, and causal understanding is the foundation of sound decision-making [1, 2]. Without it, decisions are likely to be ineffective or even harmful. Causality is a fundamental way of understanding the world, and it can be used to build more accurate and reliable AI systems. Machines equipped with a model of reality, similar to that used in causality, could have the potential to achieve strong AI and artificial general intelligence [3, 4].

These principles are equally relevant in the medical field, where understanding causality plays a crucial role in advancing medical knowledge. In medicine, most of the asked research questions are not associational, i.e., modelling statistical correlations between various quantities, but causal in nature; with these questions, researchers try to uncover the cause-and-effect relationships between variables, such as treatments, interventions, and outcomes. These questions can be hardly answered from observed data alone and could require specific and expert domain knowledge [3].

Although expert opinion remains one of if not the best tool for causal analysis (e.g., causal discovery for building causal graphs), it can be very time and resource consuming, as the amount of research data becomes larger and larger reaching dimensions that limit the possibility of parsing through the enormity of available evidence. The human factor also increases the likelihood of introducing potential errors or overlooking critical graphs details.

These difficulties could be partially solved by using large language models [5], which have been trained on vast volumes of textual data [3]. One remarkable example of such a language model is the *Generative Pre-trained Transformer* (GPT) language model. GPT LLMs are designed to understand and generate human-like text based on the given prompts and are accessible, for example, with the GPT API (see Chapter 3.3 for more details) or with the

LLM-based chatbot ChatGPT (Figure 1.1).



Figure 1.1: ChatGPT prompt example.

## 1.1 Project Objective

The primary goal of this project is to conduct an empirical study to assess the possibility of performing causal analysis using LLMs. The project focuses on the operation of causal discovery, which is the task of learning the structure of causal relationships between variables and entities; its output is a directed graph that represents the underlying data-generation process and provides insight into the true causal relationships between variables. The generated graph forms the basis for many, if not all, fundamental tasks in causal analysis, such as effect inference, prediction, and causal attribution [6].

Figure 1.2 shows an example of a simple causal graph that represents the relationship between its entities, which are encoded as graph nodes.

Figure 1.2: Causal graph example

## 1.2 Document Overview

This section provides a succinct overview of the content and objectives of each chapter in this document.

To start, Chapter 2 discusses the state of the art concerning the current causal discovery techniques. Chapter 3 then presents the methods, approaches and tools used in the project. Chapter 4 delves deeper into the implementation details of the developed software for the project. The subsequent Chapter 5 presents the benchmark tests and the evaluation metrics used to assess the quality of the results, analyses and discusses them, and addresses the limitations found. To conclude, Chapter 6 presents the final considerations on the project and topic covered, briefly indicating what was addressed, highlighting the limitations of the project, and then leaving indications of possible future work.

Causal Graph Identification by Large Language Models

# Chapter 2

# State of the Art

This chapter presents the current state-of-the-art methods and techniques used for causal analysis and causal graph identification from data, which have limitations and challenges. A different approach is the knowledge-based method, which focuses on the variables' metadata, rather than their data values. LLMs, although based on models learned from data, emerge as a third method that behaves like a virtual expert by answering questions and guiding the construction of graphs.

## 2.1 Causal Discovery from Data

The main objective of causal discovery is understanding the underlying cause-and-effect relationships in the system of interest. The current state-of-the-art techniques for causal discovery from data can be divided into two main classes: constraint-based methods and score-based methods [7, 8].

Constraint-based methods use the patterns of conditional independence among variables to deduce or determine the underlying causal relationships within a system, examining how the presence or absence of one variable affects the likelihood of another variable.

Score-based methods, on the other hand, use scoring functions to assess the quality of causal structures and then search for the optimal structure that maximises the assigned score.

Both these methods include different algorithms and techniques. While these have been shown to be effective across multiple applications and scenarios, they are not without challenges. Noise and unobserved confounders can complicate causal relationships identification. Noise is random variation in the data that can obscure the true relationships between variables. Unobserved confounders are variables that are not measured in the study, but that could be affecting the relationship between the variables of interest. Additionally, the limited availability of data further challenges the task of causal discovery.

Moreover, when dealing only with real-world observational data, it is generally not possible to

perform causal discovery and find the exact causal graph. The reason relies on the Markov equivalence class property, where multiple graphs structures are equally likely to be found, given the same data distribution [6, 9].

Different approaches and efforts have tried to overcome these limitations, yet the identification of causal graphs continues to pose challenges, especially when working with real-world observational data, revealing a concerning assessment of their effectiveness [6].

An alternative approach, known as knowledge-based method, can overcome these limitations by adopting a different focus.

## 2.2  Knowledge-Based Method

The knowledge-based method relies on the metadata associated with variables, rather than their data values. This metadata-based reasoning is typically done by human domain experts when constructing causal graphs, who use their general or specialized domain knowledge and common sense.

However, relying solely on an expert-opinion-based method can prove to be both time and resource-intensive. This approach is also susceptible to errors, as even experts can inadvertently overlook important graph details or make mistakes [3].

Large language models emerge as promising tools for tackling tasks of this nature.

## 2.3  LLM Approach

Language models provide a fresh perspective to causal discovery by adopting the same metadata-based method: having been trained on vast volumes of textual data, they "reason" through the metadata of variables and the contextual information expressed in natural language. Unlike score-based causal discovery methods, LLMs use their training knowledge combined with additional input data to identify the causal relationships between variables [6].

Multiple recent research papers and publications had a similar interest as this project, wanting to assess the causal capabilities and reasoning abilities of LLMs [1, 3, 6, 10, 11].

While this project concentrated on analyzing the causal capabilities of LLMs specifically for uncovering causal graphs, the other papers had distinct focuses or used alternative approaches. One presents an algorithm that rejects graphs from the same *Markov Equivalence Class*, while controlling the probability of excluding the true graph [1]. One tested various prompt-engineering techniques (see Chapter 3.3.2 for more details on this subject) by using different causation verbs in the prompt messages, incorporating multiple levels of detail in describing the variables and concepts, and referencing various medical authorities [3]. Another analyzed different causal reasoning capabilities in more detail, distinguishing between necessity, sufficiency, normality, and responsibility [6]. Another argued that LLMs are similar to parrots, implying that they only mimic what was seen during the training phase, raising

questions about the extent to which one can truly understand the real world's functioning by merely observing its "shadows" [10]. Lastly, another paper identifies different kinds of causal questions, noting that LLMs seem to perform well primarily with a specific group of these - namely, identifying causal relationships using domain knowledge [11].

# Chapter 3

# Background

This chapter presents a background knowledge of the base methods, approaches, and tools used in the project, which focuses on the operation of causal discovery. This task will be performed starting from natural language, that is textual data, such as scientific papers and research publications. The data will then be processed to extract the main textual entities. The key elements of the text can be identified and isolated using appropriate *Natual Language Processing* (NLP) operations, such as *Named Entity Recognition* (NER). A LLM-based discovery procedure will then be used to find the causal relationship between these entities. The final operation creates the causal graph using the causal relationships found in the previous step and plots the resulting graph. It is though necessary to present a set of definitions [12] for the key concepts on causality and graph theory before moving on with the process details. These definitions are provided in the following Section 3.1.

## 3.1 Graph Theory and Causality

Causality refers to the relationship between cause and effect, where one event (the cause) brings about another event (the effect). It is a fundamental concept in understanding the mechanisms that drive relationships within a system. Causal analysis is an operation consisting of identifying the causal graphs from a given dataset and context, by uncovering the cause-and-effect relationships and dependencies between the variables and entities of the system of interest: this is done by answering questions such as "*Which variables directly affect each other?*" or "*What is the causal directionality between variables?*". Causality relationships in systems can be represented and studied using graphs.

**Definition 1** (Graph). A graph $G = (V, E)$ is a mathematical object used to model relationships between objects, represented by a tuple of two sets: a finite set of nodes V in one-to-one correspondence with the considered objects, and a finite set of edges $E \subseteq V \times V$. In particular we consider *directed* graphs (DG), this meaning that we distinguish between the arc (X, Y) and the arc (Y, X).

**Definition 2** (Path). A path $\pi = (X - \cdots - Y)$ is a sequence of vertices such that each vertex is connected to the next vertex in the sequence by an edge. The path can be either directed or undirected. A directed path $\pi = (X \rightarrow \cdots \rightarrow Y)$ is a path where all the edges are directed.

**Definition 3** (Cycle). A cycle is a path that starts and ends at the same vertex. In causality, one of the most used type of graphs is the *Directed Acyclic Graph* (DAG), which is a directed graph with no cycles.

A causal graph is a graphical description of a system in terms of cause-and-effect relationships, i.e., the causal mechanism. Causal graphs, usually in the form of directed acyclic graphs, encode contextual knowledge of variables (both observable and unobservable) and their causal dependency. In a directed causal graph, the nodes represent the context entities and variables (e.g., in a medical context they would be symptoms, illnesses, diseases, treatments, medications, outcomes, etc. . . . ) while the edges represent the causal relationship between said entities (e.g., a medical treatment can cause a particular outcome or side effect), indicating the direction of causality.

**Definition 4** (Causal edge direction). In a causal DAG, the relationships between pairs of entities are represented using graph edges, which can take the form of directed, bidirected, or non-existing edges.

- The **directed edge** (A $\rightarrow$ B) denotes a direct causal dependency between the two features A and B, where A is a direct cause of B, without excluding the possible presence of a common cause of both A and B.

- The **bi-directed edge** (A $\leftrightarrow$ B or both A $\rightarrow$ B and A $\leftarrow$ B) represents a causal relationship where A and B are causally correlated, and the two variables have an unobserved or latent common cause.

- A **non-existent edge** denotes that no causal relationship exists between the two variables.

**Definition 5** (Direct and Indirect Cause). For each directed edge (X, Y) $\in$ E, X is a direct cause of Y and Y is a direct effect of X. Recursively, every cause of X that is not a direct cause of Y, is an indirect cause of Y.

Figure 3.1 shows a simple example of a causal graph. In this example, the variable "smoking" is the cause for both "lung cancer" and "tumors". Neither of these variables is the cause of the other; instead, they share a common cause—namely, "smoking".
The acyclicity property in causal graphs is crucial for ensuring their interpretability and for preserving the causal relationships they represent. This property is fundamental for many reasons: it ensures logical consistency, temporal ordering, identifiability of causal effects, facilitates counterfactual reasoning, and aids in prediction and intervention tasks.

Figure 3.1: Causal graph example.

Although causal graphs are mostly represented as acyclic, there are cases where cycles are accepted and even necessary to encode complex dynamics that exist in the system of interest, like for economic scenarios. Figure 3.2 shows a simple example of a case where a cycle is necessary to represent the complex relationship between investments, economic growth, and inflation. In this project, cyclical graphs will be considered incorrect.



Figure 3.2: Graph example with cycle.

Central items to the causal graphs are the nodes, which represent essential elements within the domain of interest. In the context of this project, the domain pertains to the medical and health-related field, so the nodes would be symptoms, diseases, treatments, outcomes, etc. . . . The data used for constructing these causal graphs is collected in a textual form, from

abstracts of medical publications and research papers. For this reason, it is necessary to process the textual data in order to extract the entities. This can be achieved with appropriate *Natual Language Processing* (NLP) operations, such as *Named Entity Recognition* (NER).

## 3.2   Named Entity Recognition

*Natual Language Processing* (NLP) is a field of artificial intelligence that makes human language intelligible to machines. It involves the use of algorithms and statistical models to enable computers to understand and process human language. NLP has many applications, including machine translation, sentiment analysis, text summarization, and speech recognition. It is used in a wide range of industries, including healthcare, finance, and marketing [13]. Large language models can be considered among the most advanced and significant expressions of NLP applications. These models are at the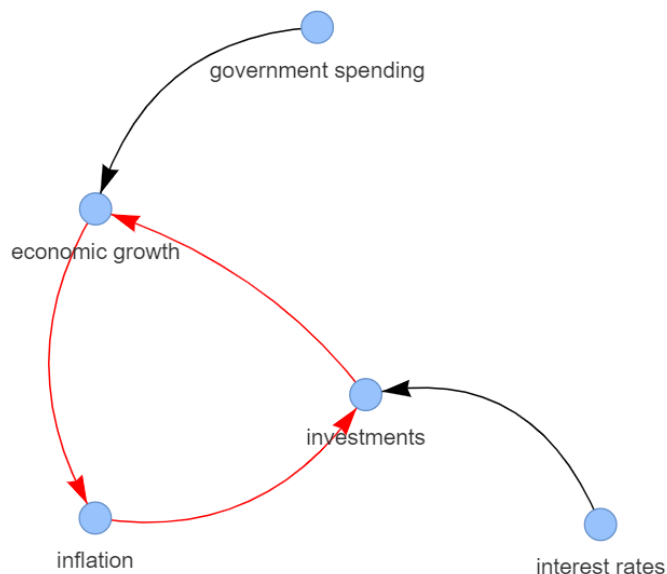 forefront of NLP research and technology due to their ability to generate coherent and contextually relevant human-like text.

Named entity recognition (NER) is a crucial NLP task that aims to identify and classify named entities within text, as shown in Figure 3.3. In the context of medical texts, NER plays a vital role in extracting specific medical entities such as diseases, symptoms, treatments, drugs, anatomical terms, and medical procedures. Medical texts pose challenges for NER due to their specialized terminology, which often includes abbreviations and multiple names and entities that are used interchangeably as synonyms. Additionally, the complex language structures found in medical texts, along with the diverse sources from which they originate, further complicate the NER process.



Figure 3.3: NER example.

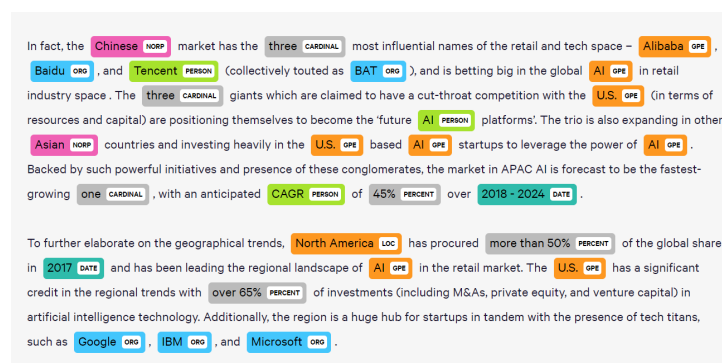The NER operation can be conducted in several ways, with many tools and packages, such as *spacy*[1], *flair*[2], and the LLM *Bert*[3].

---

[1]spaCy · Industrial-strength NLP in Python. https:// spacy.io/. (accessed: 23.08.2023).
[2]flair. https://flairnlp.github.io/. (accessed: 23.08.2023).
[3]J. Devlin et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding".

## 3.3  GPT API

Part of the project's activities are performed using the GPT LLM through the GPT API, including tasks like causal relationship identification. While numerous methods exist for NER, this particular step is also performed using the GPT LLM to follow the project's general idea of employing LLMs, specifically the GPT model.

The GPT LLMs are developed by OpenAI, an American artificial intelligence research company. The original paper on generative pre-training of a transformer-based language model (GPT) was published in 2018. GPT-2 was released in 2019, and GPT-3 was released in 2020. On March 14, 2023, OpenAI announced the release of GPT-4, capable of accepting text or image inputs [14].

The GPT API is a tool that provides access to OpenAI's GPT models, allowing the integration of natural language processing capabilities into applications. It is a RESTful API, meaning that it is a type of web-based application programming interface that follows the principles of Representational State Transfer (REST) to enable communication and interaction between different software systems over the internet [15]. Developers can send requests to the API using any programming language, but there are official OpenAI packages and libraries that simplify the development process. For this project, the Python package was used [16].

When using the API, the user specifies the model to use and provides the necessary parameters such as the prompt and optional messages to contextualize the use and behavior of the model, i.e., how the model should answer to requests. The API then processes the request and returns as a response the generated context-aware text based on the provided input.

The GPT API can be employed in a variety of applications and use cases. It can be used to generate conversational agents, draft emails or other pieces of writing, provide language translation, answer questions, or assist with content creation [17].

### 3.3.1  Using the GPT API

The GPT API is used by specifying the model to use (e.g., GPT-4) and by providing additional messages, including the system message and the user message. These serve as instructions to the model, with the system message being a system level instruction to guide the model's behavior throughout the conversation (e.g., asking the model to answer or to act in a specific way), and the user message functioning as the actual request the model is required to answer.

The following example [18] shows how the model's behavior can be guided throughout the conversation, by stating in the system message that the model is "an assistant that speaks like Shakespeare", and asking it to tell a joke.

```
1  'messages': [
```

```
2    {'role':'system', 'content':'You are an assistant that speaks like
     Shakespeare.'},
3    {'role':'user', 'content':'Tell me a joke.'}
4  ]
```

The model answers the request given in the user message by following the instructions of the system message.

```
1  {...
2  'message': {'role':'assistant',
3             'content':'Why did the chicken cross the road? To get to
     the other side, but verily, the other side was full of peril and
     danger, so it quickly returned from whence it came, forsooth!'}
4  ...}
```

In particular, the system message is used to contextualize the model and its behavior, to make it more useful and accurate for the required operation: for the project's causal analysis tasks, for example, the system message was "*You are a helpful assistant for causal reasoning and cause-and-effect relationship discovery*", to try steering the output space to more causally consistent answers. This was shown being an effective prompt-engineering technique that results in more accurate answers [6]. Prompt-engineering will be discussed in more detail in the next Section 3.3.2.

The system message helps set the behavior of the assistant, by modifying the *personality* of the assistant or providing specific instructions about how it should behave throughout the conversation. However, this message is optional and the model's behavior without a system message is likely to resemble that of using a generic one, such as "*You are a helpful assistant.*" [17].

Another parameter that can be optionally set is the *temperature*: this represents the degree of exploration or randomness of the model's output. It ranges from 0 to 2, with a default setting of 1. A higher value (e.g., 0.8) increases creativity and diversity but might be less focused. A lower value (e.g., 0.2) produces more deterministic output following patterns [17, 19]. The default value for the temperature parameter used in this project was 0.2.

Follows a complete example of a GPT API chat completion request that uses the GPT-4 model, and specifies system, assistant, and user messages.

```
1  openai.ChatCompletion.create(
2    model="gpt-4",
3    messages=[
4        {"role": "system", "content": "You are a helpful assistant."},
5        {"role": "user", "content": "Who won the world series in 2020?"},
6        {"role": "assistant", "content": "The Los Angeles Dodgers won the
     World Series in 2020."},
7        {"role": "user", "content": "Where was it played?"}
8    ],
9    temperature=0.3,
10 )
```

STUDENTSUPSI

An example of response looks as follows:

```
{
  "choices": [
    {
      "finish_reason": "stop",
      "index": 0,
      "message": {
        "content": "The 2020 World Series was played in Texas at Globe
    Life Field in Arlington.",
        "role": "assistant"
      }
    }
  ],
  "created": 1677664795,
  "id": "chatcmpl-7QyqpwdfhqwajicIEznoc6Q47XAyW",
  "model": " gpt-4",
  "object": "chat.completion",
  "usage": {
    "completion_tokens": 17,
    "prompt_tokens": 57,
    "total_tokens": 74
  }
}
```

In Python, when utilizing the OpenAI API to interact with the assistant, the response received is structured as a dictionary, in JSON format. The content of the response can be accessed by first navigating through the dictionary keys: the "choices" key corresponds to the list of responses provided by the assistant. Using *[ 'choices' ][0]* it is possible to retrieve the *response* from the list. Inside the first response item, the key labeled "*message*" contains the details of the assistant's message. By using *[ 'message' ][ 'content' ]* it is possible to access the message content of the response.

The following code snippet can be used to extract the assistant's reply.

```
response['choices'][0]['message']['content']
```

Varying temperature values and system messages deeply influences the response of the LLM, even when the user message remains the same. Table 3.1 provides an example using the GPT-4 model. In this example, the LLM was queried with the same user message ("*Complete the sentence: I'll fly a starship...*"), but with two different system messages. Each of these system messages was used with two different temperature values: $0$ and $1.9$.

### 3.3.2 GPT Prompt Engineering

Depending on the task, the GPT LLM can typically produce satisfactory results when asked to answer a question. However, there are some expedients that have shown to be beneficial and to increase the results accuracy when querying the LLM. These techniques are part

| System message | Temperature | Response |
|---|---|---|
| You are a creative marketing manager. | 0.0 | ... across the universe of innovative ideas, navigating through the galaxies of consumer needs and desires, to land on the perfect marketing strategy that will skyrocket our brand to unparalleled heights. |
| | 1.9 | ... across the universe of ideas, innovating and turning impossibilities into massive revenue streams for our organization. |
| You are an aspiring country musician with a warm and inviting personality. | 0.0 | ... across the universe divide, just to find that sweet country rhythm that makes my soul come alive. |
| | 1.9 | ... across the Universe divide, strumming my guitar beneath those endless Skies, captivating hearts while relayin' wholesome country melody-vibe. |

Table 3.1: Example of same user message with different system message and temperature values.

of the discipline known as *prompt-engineering*, which comprises a set of rules and instructions that serve as guidelines to enhance the capabilities of language models (particularly large language models, such as the used GPT-3.5 and GPT-4 models) on a wide range of common and complex tasks [20, 21]. Prompt engineering involves crafting system and user messages that guide the model's responses and shape its output to meet specific requirements [22].

The goal of prompt engineering is to provide contextual cues and instructions to the language model, enabling it to generate more accurate, relevant, and desired responses. By designing prompts, it is possible to tailor the behavior of LLMs, making them more suitable for various tasks, domains, and user needs.

Among the many prompt engineering techniques [19, 20, 22, 23], one of the most important strategies is improving the clarity and precision of the prompt text by providing clear and specific instructions. Delimiters like brackets, tags or quotes can segregate sections within the prompt, aiding in a more organized interpretation of the input. Furthermore, prompting for structured output by specifying the desired response format guides the model in generating well-organized results.

Checking task conditions also ensures the necessary assumptions are met. For instance, the prompt can verify whether essential information is available to complete the task and provide alternative instructions if this information is missing.

```
1 You will be provided with text delimited by triple backticks.
2 Read it carefully to understand the context and content. If it contains a
```

```
      sequence of instructions , re -write those instructions in the
      following format :
3
4 Step 1 - ...
5 Step 2 - ...
6 ...
7 Step N - ...
8
9 If the text does not contain a sequence of instructions , then simply
      write "No steps provided."
10
11 '''{text_1}'''
```

The technique called *few-shot prompting* instead involves showing successful task examples to the model before requesting similar ones. This helps the model understand the context better, preparing it to deliver pertinent and accurate responses.

```
1 This is awesome !          // Negative
2 This is bad !              // Positive
3 Wow that movie was rad !   // Positive
4 What a horrible show !     // ?
```

Another principle of prompt engineering consists in providing the model with enough time to "think". This involves requesting the model to answer with a step-by-step explanation of its thought process before providing the final answer. By doing so, the model is directed to work out its own solution rather than rushing to conclusions. This principle applies, for example, when the model needs to verify the accuracy of a given solution: in this case, it is prompted to formulate its solution first and then compare it to the provided one [19, 23].

```
1 Your task is to determine if the student's solution is correct or not.
2 To solve the problem do the following :
3 - First , work out your own solution to the problem.
4 - Then compare your solution to the student's solution and evaluate if
      the student's solution is correct or not.
5 Don't decide if the student's solution is correct until you have done the
       problem yourself.
```

Appendix B shows the user messages used in the project.

# Chapter 4

# Methodology and Implementation

This chapter presents the implementation details of the project. The implementation work performed within this thesis project can be divided in two main steps: data collection and data analysis. The former one consisted in collecting the necessary data for the latter, which can itself be divided into multiple other sub-operations. This section first presents the data collection process, highlighting the usage of the *National Center for Biotechnology Information* (NCBI) API for requesting the necessary textual data we use for evaluate our methods. It then delves into the operations of data processing and causal analysis.

## 4.1 Scraping

The first step of the project consists in collecting the necessary textual data for testing the causal discovery capabilities of the GPT LLM.

This involved implementing a series of Python scripts to automate a data acquisition process. Considering the project's focus on medical papers and health-related publications, this general scraping process was adapted to suit the relevant context. As a result, the data was sourced from the PubMed database. PubMed serves as a freely accessible search engine mainly utilizing the *MEDLINE* database, which contains references and abstracts related to life sciences and biomedical subjects [24].

The scraping process only extracts the article abstract and extra details, such as title, publication date, and keywords.

### 4.1.1 Scraping Pipeline

A pipeline handling the essential operations was created for extracting the necessary textual data from the PubMed database. To automate this extraction process, a Python script was written using the public API provided by the NCBI as interface into its query and database system. The main flow is shown in Figure 4.1. The pipeline allows the user to extract textual data from PubMed by searching for specific terms.
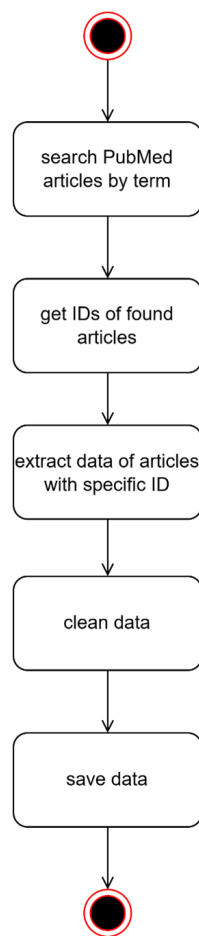
Figure 4.1: PubMed data extraction flow.

The pipeline's main operations are handled by the *search_by_terms*, *get_articles_data*, and *clean_data* procedures. Appendix A.1 presents the source code of these functions.

### 4.1.2  Collecting the Data

The *search_by_terms* procedure is the first operation of the pipeline. As the name suggests, it allows the user to search for articles in the PubMed database containing the specified search terms. These terms are joined as query parameters in the request URL. An API key is also sent in the request URL, to ensure smooth and supported access to the desired resources.

The response is in a xml format, and it is processed to extract all ID numbers of the articles found in the specified NCBI database, which in this case is PubMed.

The implemented function returns the extracted IDs. However, the script allows users to utilize the NCBI *Entrez History* feature, which proves to be significantly more efficient when dealing with tasks that involve searching for or downloading a substantial number of records. This approach helps to streamline the process and optimize the retrieval of records in a more efficient manner, making it possible to upload many IDs or download several hundred records at once [25].

The pipeline then continues with the second step of the data acquisition process, with the *get_articles_data* function. This procedure queries the NCBI for the actual content of the articles with the specified ID.

The NCBI API allows users to query article data with the article ID or by using the Entrez History feature, which can provide a more efficient data retrieval. A URL parameter of the request defines the main data content requested, which, in this case, are the abstracts (*rettype=abstract*).

The returned data is in a xml format, and it is processed and parsed to extract the necessary information. The recovered data from the article includes the abstract, the article ID number, the title, the keywords, and the publication date. Follows an example of data retrieved from PubMed, which is then processed by the scraping pipeline to extract the necessary information.

```
1  <PubmedArticleSet>
2    <PubmedArticle>
3      <MedlineCitation>
4        <PMID>37535727</PMID>
5        <Article>
6          <ArticleTitle>Side effects loom over Alzheimer's drugs.</
     ArticleTitle>
7          <Abstract>
8            <AbstractText>
9              Despite landmark antibody approval, research into brain
     swelling and bleeding lags.
10           </AbstractText>
```

```
11          </Abstract >
12          <AuthorList >
13            <Author >
14              <LastName >Couzin - Frankel </LastName >
15              <ForeName >Jennifer </ForeName >
16              <Initials >J</Initials >
17            </Author >
18          </AuthorList >
19        </Article >
20      </MedlineCitation >
21      <PubmedData >
22        <History >
23          <PubMedPubDate >
24            <Year >2023</Year ><Month >8</Month ><Day >3</Day >
25            <Hour >19</Hour ><Minute >14</Minute >
26          </PubMedPubDate >
27        </History >
28      </PubmedData >
29    </PubmedArticle >
30 </PubmedArticleSet >
```

### 4.1.3  Data Preparation

The *clean_data* procedure is the third and last step of the data acquisition pipeline. It performs cleaning operations on the acquired data, such as removing null abstract values, duplicates, and potentially removing data of articles published in a particular date range.

The cleaned data is eventually saved to a file and returned by the scraping pipeline. This data can also serve as input for the causal analysis pipeline, which will process it to extract causal graphs from the obtained abstracts.

## 4.2  Causal Analysis

After completing the preliminary phase of data collection, the main focus of the project shifted towards the actual analysis operations. The primary objective was to investigate the potential of LLMs in causal discovery.

Causal analysis involves the process of revealing the cause-and-effect relationships and dependencies among variables from a provided dataset and context, by answering questions such as "*Which variables directly affect each other?*".

This step of the project involved processing the collected data from the PubMed database to extract information from the abstracts. This consisted of extracting the main named entities within the textual data, such as the names of diseases, drugs, and genes. These named entities were then used to perform the actual causal analysis. The output of the causal

Figure 4.2: Causal discovery pipeline flow.

relationship discovery process is then used to plot the resulting causal graph represented as a directed graph and possibly depicted by dedicated tools (see Section 4.3).

### 4.2.1   Causal Analysis Pipeline

The components of the causal analysis process, consisting of various sub-steps, have been integrated into a single operational pipeline, called *causal_discovery_pipeline*, whose general flow is shown in Figure 4.2. These operations include extracting entities from the textual data with the *gpt_ner* function, performing the actual causal analysis on the found entities using the *gpt_causal_discovery* procedure, and ultimately generate the resulting causal graph with the *build_graph* function. Appendix A.2 presents the source code of these functions.

### 4.2.2   Extracting Medical Entities from Text

As previously mentioned, the second part of the project consisted in working with the collected data. The first step of the operation involved performing Named Entity Recognition on the abstracts, a fundamental procedure to extract and classify named entities. This step was essential for further processing and analysis.

The NER operation was performed using the GPT LLM. To enhance the performance of the Language Model for the task, both the system and user messages were designed accordingly.

The employed system message was "*You are a helpful assistant for Named Entity Recognition of medical texts*" to provide guidance to the model and improve its understanding of the task at hand.

To further aid the model's comprehension, the user message was crafted using the abstract of the medical text, complemented with additional information about the types of entities to be extracted. In this case, since the texts were focused on medical literature and research publications, the model was explicitly instructed to identify entities with a particular emphasis on diseases, medications, treatments, symptoms, etc.... .

The intention of customizing the user message by providing relevant context and specific entity requirements, was to guide the LLM towards producing more accurate and relevant results for the ongoing NER operation. Appendix B.1 shows the full user message.

The result of the *gpt_ner* function is an array containing all the found entities and it is then used for the subsequent causal analysis.

The *causal_discovery_pipeline* also provides the option to include an entity optimization step through the *optimize_entities* procedure: by using the GPT API, the pipeline operation focuses on identifying synonyms, redundant entities, or entities and names that can be used interchangeably. In the generated output, entities with synonymous or similar meanings are matched together. Appendix B.2 shows the full user message used for the entity optimization function.

### 4.2.3   Causal Discovery

With the completion of the NER operation and the extraction of entities, we now proceed to the central step of the pipeline: the causal discovery operation.

This step consists of the *gpt_causal_discovery* function, which takes the found entities and performs causal discovery. The approach for this operation uses a naïve discovery procedure. This method does not use actual data as input to perform inferences: in this sense, it resembles how humans recall facts by relying on memorized knowledge rather than actively referencing actual data [10]. This approach aims to infer the causal relationship among the various variables by querying the LLM regarding the direction of the pairwise causal relationships for each possible pair combination.

The type of causal relationship between a pair of entities corresponds to the edge orientation in the causal graph: directed edges indicate direct causes, bi-directed edges represent entities that are causally correlated and the two have an unobserved or latent common cause, and non-existent edges indicate the absence of a causal relationship between the variables.

To infer the direction of the causal edge, the pipeline function queries the LLM to determine which cause-and-effect relationship is more likely between the two entities. The system message used for this operation is "*You are a helpful assistant for causal reasoning and cause-and-effect relationship discovery*", to try guiding the output towards more causally consistent answers. On the other hand, the user message introduces the current pair of entities of interest, asking a single question about the direction of the causal dependency. It also requests a step-by-step explanation in response. The possible answers the LLM is expected to choose from are also listed within the user message in the following form:

A. "X" causes "Y";

B. "Y" causes "X";

C. "X" and "Y" are not causally related;

D. there is a common factor that is the cause for both "X" and "Y"

To enhance the accuracy and exploration of cause-and-effect relationships, the prompt uses random verbs of causation when querying the GPT LLM. These are randomly chosen from a set of causation verbs, containing verbs such as "*provokes*", "*triggers*", "*causes*", "*leads to*", "*induces*", "*results in*", "*generates*", and more (see Appendix B.3 for all other causation verbs). For example, if the process is querying about the type of causal relationship between the variables "*cigarette smoking*" and "*lung cancer*", and the randomly chosen verb of causation is "*provokes*", the resulting options in the user message are:

A. "cigarette smoking" provokes "lung cancer";

B. "lung cancer" provokes "cigarette smoking";

C. "cigarette smoking" and "lung cancer" are not causally related;

D. there is a common factor that is the cause for both "cigarette smoking" and "lung cancer"

This approach can be beneficial in terms of coverage of language patterns and potential causal relationships, can reduce the risk of bias that may come from consistently relying on a specific verb, and can encourage the model to explore different relationships between variables, allowing for a more comprehensive analysis of the data [1, 3]. Appendix B.3 shows the full user message used for the causal discovery process.

The LLM is queried for each pair of variables. In a default execution mode, the pipeline examines combinations (without repetition) of all pairs of the identified entities. The total number of queries (one for each pair) is

$$C_k(n) = \binom{k}{n} = \frac{n!}{k! \, (n \, - \, k)!} \; , \tag{4.1}$$

where *n* is the total number of elements in a set (number of entities), and *k* is the number of elements to be chosen each time from the set (two, in this case, as the query is done for each pair of variables). In the case of ten entities, the total queries are $C_2(10) = \binom{2}{10} = \frac{10!}{2! \, (10 - 2)!} = \frac{10 \cdot 9}{2} = 45$.

The pipeline, however, allows to perform a double test for each pair of variables, checking all potential variations without repetition (i.e., relationship "X" - "Y" and "Y" - "X"). As a result, the LLM is queried twice for each pair of entities, with a total number of variations of

$$V_k(n) = \frac{n!}{(n \, - \, k)!} \; . \tag{4.2}$$

Again, *n* is the total number of entities and *k* is two. With ten entities, the total queries are $V_2(10) = \frac{10!}{(10 - 2)!} = \frac{10!}{8!} = 10 \cdot 9 = 90$.

The complexity for this operation has an upper bound of $O(n^2)$.

In this case, the pipeline also performs a compatibility check of the answers: the response to the query regarding the causal relationship between "X" and "Y" is cross-validated with the response to the query about the relationship between "Y" and "X", and the two answers must be consistent with each other. This validation process is managed by the *check_invalid_answers* function, which assesses response consistency and distinguishes between valid and invalid edge directions.

The edge direction and causal relationship associated with "invalid" answers are then re-queried using the *correct_invalid_edges* function. In this process, the LLM is queried again by also adding, in the user message, the inconsistent answers obtained earlier, seeking the most likely relationship between the given variables. The newly acquired answers are then appended to the previously identified "valid" edges.

The output of this pipeline step involving the *gpt_causal_discovery* function is an array containing the type of causal relationship between each pair of entities. Appendix B.4 shows the user message used for the invalid edge correction process.

## 4.3   Analysing the Graph

The upcoming and final operations of the project involve preparing for the construction and analysis of the resulting causal graph, using libraries and packages such as *NetworkX*, *Pyvis*, and *graph-tool*. Appendix A.3 presents the source code of the used functions.

### 4.3.1 Graph Preprocessing

Before plotting the graph, the pipeline performs an intermediate operation of edge and node preprocessing. This step assists the following ones by decoding the LLM answers and converting them into sets of nodes and normalized directed edges, represented in the form of "X → Y".

The main procedure responsible for this operation is the *preprocess_edges* function. This function generates various components: a set consisting of all graph nodes (entities previously extracted), an array containing the normalized directed edges and one containing all bidirectional edges.

The *preprocess_edges* function relies on another procedure for the normalization of edges, which is the *normalize_edge_direction* function. This procedure takes as input the nodes involved and the LLM's response regarding their causal relationship. The function processes the output of the LLM and returns the corresponding edge in the form of "X → Y", representing the causal dependency between the nodes.

An additional step is taken before plotting the graph, where a check is performed to determine whether the resulting graph is acyclic.

### 4.3.2 Cycle Check

To guarantee logical consistency the resulting graph should be a DAG, which is a directed graph without cycles. In the context of causal analysis and causality in general, the acyclic property of graphs is crucial to maintain the logical meaning and coherence encoded within the graph. The absence of cycles ensures that there are no circular dependencies or contradictory relationships, allowing for a clear and meaningful representation of causality in the system.

The *find_cycles* function is the dedicated procedure designed to determine whether the constructed graph contains any cycles. It accepts an array of nodes and an array of edges as input parameters, which together represent the graph. The function makes use of the *graph-tool*[1] Python package, a highly efficient module for graph manipulation and analysis. The underlying components of this package are primarily implemented in C++ to optimize performance.

The cycles are identified using the *all_circuits* function from the graph-tool package. This function is based on the Johnson's algorithm for finding all elementary circuits of a directed graph; it is similar to other algorithms (e.g., by Tarjan[2]), but it is faster because it considers each edge at most twice between any one circuit and the next in the output sequence. The algorithm is time bounded by $O((n+e)(c+1))$ and space bounded by $O(n+e)$, for a graph with *n* nodes, *e* edges and *c* cycles [26].

---

[1]T. P. Peixoto. "The graph-tool python library". In: figshare (2014).
[2]R. Tarjan. "Depth-First Search and Linear Graph Algorithms". In: SIAM Journal on Computing 1.2 (1972), pp. 146–160.
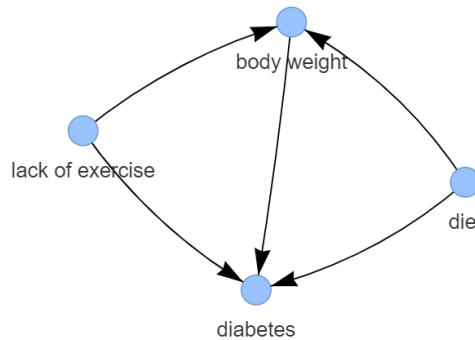
Figure 4.3: Interactive graph generated with Pyvis for the conducted causal analysis.

In a directed graph, the number of potential cycles has a tendency to grow exponentially as the number of nodes increases. As a result, even relatively small graphs can contain a substantial number of them. Since listing and returning all cycles of a graph is not the main focus of this project, the *find_cycles* function only returns a symbolic number of the first 100 found, assuming there are any.

In case cycles exist, they are represented as lists of graph nodes and returned as an output parameter of the *find_cycles* function.

### 4.3.3   Plotting the Graph

The operation for plotting the causal graph is the last step of the casual discovery pipeline, and it is processed by the *build_graph* procedure. This function is designed to construct and visualise a directed graph using libraries for graph and network creation, manipulation, and analysis, such as the NetworkX and Pyvis libraries. NetworkX is dedicated to general-purpose graph operations [27], while Pyvis serves as a visualisation library suitable for generating interactive network graphs [28].

Within a graph, bidirected edges are coloured in grey, and the function allows for highlighting cycles, as well as creating interactive plots. Cycles within the graph can be highlighted by coloring relevant edges in red. It also supports both static and interactive modes of graph presentation, simplifying the visualisation and analysis of entity relationships. The final interactive graph is then exported as an *.html* file.

Figure 4.3 presents an example of the outcome of the conducted causal analysis. It displays the resulting interactive causal graph, which has been plotted using the Pyvis package.

## 4.4 Running the Process

A separate Python script called *causal_analysis.py* was implemented to allow running the entire process from a terminal. The script can be used to execute both the scraping and causal discovery operations, either sequentially (using the data extracted from the scraping step as input to the causal discovery step) or independently.

The script also allows running the causal discovery on data specified by the user or on the benchmarks, where the user can also choose to use the random baseline algorithm or the GPT LLM.

The following examples show two of the possible commands to run the causal analysis process on a specified data file (defined with the *data-path* option) or on the benchmark tests using the GPT LLM (instead of the baseline algorithm).

```
1 $ python causal_llms.py c --data-path=../data/abstracts.csv
2 $ python causal_llms.py b --algorithm=gpt
```

Appendix D shows all command line actions and options.

# Chapter 5

# Analysis and Evaluation

This chapter discusses the evaluation of the used causal discovery process, by presenting and analysing the achieved results. Various benchmarks and metrics are used to evaluate the performance of this approach. This analysis aims to provide a comprehensive understanding of the strengths and limitations of this method and approach.

## 5.1 Benchmarks and Metrics

This section introduces the benchmarks used to evaluate the capabilities of our LLM approach to causal discovery. These include tests for two types of tasks: pairwise causal relationships and full graph discovery. Additionally, real-world text evaluations are conducted, focusing specifically on abstracts from medical papers, which are relevant to this project. The chapter then shows the set of evaluation metrics used to assess the quality of the achieved results, like commonly used *precision*, *recall*, *F1 score*, and *Structural Hamming Distance* (SHD).

Because of its relatively recent popularity, this subject has witnessed many contributions over the past few months. Many publications and research papers have tested the general capabilities of LLMs using standardized exams and tests written to assess human aptitude and knowledge across various domains [29, 30, 31, 32, 33, 34]. For causal reasoning capabilities, researchers have used widely known and established benchmarks with datasets from multiple domains, including medicine and climate science [6].

### 5.1.1 Pairwise Causal Relationship Discovery

In the context of causal discovery abilities, these benchmark datasets mainly consist of lists of variable pairs, where each pair represents a causal relationship that can be encoded as a directed edge in a causal DAG.

The assessment of LLMs' causal discovery capabilities involves tasks that focus on identifying pairwise causal relationships, determining whether variable *A* causes variable *B* or

| Variable A | Variable B | Domain |
|---|---|---|
| Age of Abalone | Shell weight | Zoology |
| Cement | Compressive strength of concrete | Engineering |
| Alcohol | Mean corpuscular volume | Biology |
| Organic carbon in soil | Clay content in soil | Pedology |
| Photosynthetic Photon Flux Density | Net Ecosystem productivity | Physics |
| Drinking water access | Infant mortality | Epidemiology |

Table 5.1: Example cause-effect pairs from the Tubingen benchmark.

vice-versa. These tasks involve both well-known scenarios that an average non-expert can correctly address using common sense and basic field knowledge (e.g., Tubingen cause-effect pairs dataset [35]), as well as more specialized domains that require expertise in a specific field to ensure accurate understanding and interpretation (e.g., Neuropathic pain dataset [36]). Table 5.1 [6] shows some examples of the cause-effect pairs and respective domains of interest from the the Tubingen benchmark, where the task is to determine which variable is the cause and which one the effect.

As previously mentioned, it has been noted that prompt engineering significantly increases the accuracy of results when querying the LLM for causal dependencies and edge directions [21]. Furthermore, using advanced LLMs, such as GPT-4, along with these prompt engineering techniques results even in higher accuracy.

### 5.1.2 Full Causal Graph Identification

Since the primary objective of this project is to evaluate the abilities of LLMs in identifying complete causal graphs, the LLM was tested using slightly different benchmarks compared to the ones mentioned earlier. Extending the task from simple identification of pairwise causal relationships to full graph discovery introduces additional challenges that are not present in the former task. These include, for example, the need to avoid introducing edges between unrelated variables and distinguishing between direct and indirect causes [6]. The adopted strategy, as discussed in the previous chapters, involves enumerating all possible pairs of variables and performing the pairwise test for each pair combination.

For this project, the LLM was tested against existing causal graphs, which served as benchmarks representing the ground truth. The graphs used as ground truth [3, 37] predominantly revolve around medical and health-related subjects, as the project focused on identifying causal relationships and uncovering causal graph structures within the medical context. Figure 5.1 shows one of ground truth causal graphs used in the tests. Appendix C.1 shows all other ones.
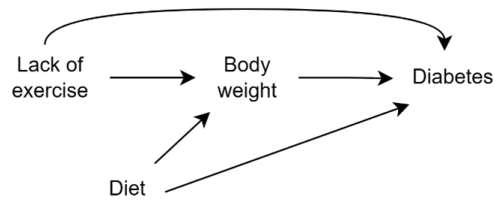
Figure 5.1: 'Diabetes' benchmark.

### 5.1.3   Evaluation Metrics

Various evaluation metrics are used to assess the quality of the obtained causal discovery results. These metrics aim to identify shared patterns between the ground truth model and the one generated from the process. Given that the ground truth when dealing with causality and causal discovery is commonly represented in a graph form (e.g., DAGs), these metrics are also related with network ones.

These include commonly used metrics like precision, recall which are calculated using classification ones, such as *True Positives*, *True Negatives*, *False Positives*, and *False Negatives*. In the context of causal discovery and causal relationship identification, True Positives (TP) represent the cases where the graph discovery process correctly identifies the direction of a causal edge as present in the ground-truth graph. True Negatives (TN) instead, refer to instances where the process correctly identifies the absence of a causal edge in the ground-truth graph. False Positives (FP) occur when the process incorrectly identifies the presence of a causal edge that is not present in the ground-truth graph. Ultimately, False Negatives (FN) represent cases where the process fails to identify a causal edge that is actually present in the ground-truth graph.

Other evaluation metrics used are F1 score, accuracy, Structural Hamming Distance (SHD) (calculated with the *causal-discovery-toolbox*[1] package) and more. Table 5.2 [38] lists the evaluation metrics used for the benchmark tests.

## 5.2   Results and Discussion

This section presents the results of the causal discovery tests conducted on benchmarks and real medical texts, and addresses the limitations found. The results are discussed, showing how well the tests performed compared to the ground truth graphs. Among the metrics used to evaluate the performance of the tests are precision, recall, and F1 score.

---

[1]D. Kalainathan and O. Goudet. "Causal Discovery Toolbox: Uncover causal relationships in Python". In: arXiv e-prints, arXiv:1903.02278 (Mar. 2019).

| Metric | Description |
|---|---|
| Missing edges | Number of edges that are present in the ground truth graph but not in the generated one |
| Extra edges | Number of edges that are present in the generated graph but not in the ground truth one |
| Correctly directed edges | Number of edges present in the generated graph that were correctly directed |
| Incorrectly directed edges | Number of edges present in the generated graph that were incorrectly directed |
| Structural hamming distance | Sum of missing edges, extra edges, and incorrectly directed edges |
| Precision ($\mathrm{Pr}$) | Measures how many of the identified causal relationships are correct out of the total relationships identified. $\mathrm{Pr} = \frac{\mathrm{TP}}{\mathrm{TP+FP}}$ |
| Recall ($\mathrm{Re}$) | Measures the ability to identify all actual causal relationships. $\mathrm{Re} = \frac{\mathrm{TP}}{\mathrm{TP+FN}}$ |
| F1 score ($\mathrm{F1}$) | Harmonic mean of precision and recall. $\mathrm{F1} = 2 \cdot \frac{\mathrm{Pr \cdot Re}}{\mathrm{Pr+Re}}$ |
| Precision-Recall Curve | Depicts the trade-off between the precision and recall of the identified causal relationships. |
| Area Under PR Curve | Quantifies the overall performance by summarizing the precision-recall trade-off across different thresholds. |

Table 5.2: Evaluation metrics.

STUDENTSUPSI

| Model | SHD | Accuracy | Precision | Recall | F1 score |
|---|---|---|---|---|---|
| Random | 8.3 | 0.20 | 0.33 | 0.38 | 0.36 |
| GPT-3.5-turbo | 4.17 | 0.61 | 0.71 | 0.62 | 0.66 |
| GPT-4 | 1.67 | 0.86 | 0.89 | 0.98 | 0.93 |

Table 5.3: Benchmark results.

### 5.2.1 Causal Discovery from Benchmark Tests

The benchmark evaluation allows to quantify the capabilities of the LLM by comparing its predictions against a ground-truth reference, with commonly used evaluation metrics, like accuracy, precision, and recall.

The benchmark dataset consists of a collection of relatively small causal graphs [3] (less than $10$ nodes), each represented by a set of nodes and a set of arcs. The causal discovery process begins by considering the graph nodes from the ground truth dataset and proceeds to identify the causal relationships among these nodes. The results are then evaluated through a comparison with the ground truth graphs.

The benchmarks are run with different LLM-based methods and algorithms. An algorithm returning either edge at random for each pair of entities was used as a baseline. The randomness is uniform, meaning that each possible outcome (in this case, each edge direction) has an equal probability of being selected. Table 5.3 shows the results (with SHD, precision, recall, and F1 score) of the different methods and algorithms applied to the benchmark tests. Appendix C.2 shows additional results using all the evaluation metrics defined in Section 5.1.3.

As expected, the baseline algorithm (indicated by the "Random" row in the table) presents the poorest performance, with an average structural hamming distance exceeding $8$ errors, a precision of $0.33$, a recall of $0.38$ and a resulting F1 score of approximately $0.35$.

The *GPT-3.5-turbo* model instead shows better results, with an average SHD that is half of the one achieved by the random baseline. Both precision and recall values show notable improvements compared to the baseline method, with a F1 score of about $0.66$. The highest performance is achieved by the most advanced LLM, the *GPT-4* model. The results show an important reduction in SHD to an average of $1.6$ errors, while precision rises to $0.89$ and recall achieves an impressive $0.98$, leading to a F1 score of $0.93$. This score nearly triples the F1 score achieved by the random baseline, showing the non-trivial contribution of LLM outputs in facilitating the identification of causal graphs. The GPT-4 results achieved an average accuracy of about $0.86$, which is $0.20$ higher than that achieved in the paper where the majority of the ground truth graphs were taken from [3].

### 5.2.2   Causal Discovery from Real Medical Abstracts

The previously presented tests used existing causal graphs as ground-truth benchmarks. It is thought important to test this approach within real-world contexts. Causal graph identification is a challenging task, especially when the text is complex and contains a lot of technical jargon. As previously discussed, medical texts are a particularly difficult case, as they are often full of abbreviations, acronyms, and other specialized terms. This makes it difficult for language models to identify the main entities in the text and to build accurate causal graphs. For this reason, the process was evaluated by extending its application to real-world scenarios, using actual medical texts.

As expected, the anticipated accuracy of the causal graphs tends to decrease as the complexity of the text increases. The results suggest that this approach is a promising one, but that there is still room for improvement.

LLMs are a powerful tool, but they are not perfect. They can be biased, make errors, and even hallucinate, by answering questions about obscure topics and making things up that sound plausible but are not actually true. Therefore, it is important to use them in conjunction with other methods, such as human judgment and domain knowledge.

The results of LLMs used for causal discovery can serve as a valuable head start, as they can be used to identify potential causal relationships by distinguishing words and phrases that are often used together in sentences that describe such relationships [10].

The following section presents a full example of the causal analysis process applied to the abstract of a real medical article extracted from the PubMed repository.

### 5.2.3   Example with Real Medical Text

The complete causal analysis process was run on multiple abstracts from PubMed. This section presents one example; others can be found in Appendix E.

The title of this article is "*Research progress on the protective mechanism of a novel soluble epoxide hydrolase inhibitor TPPU on ischemic stroke*" [39] and the abstract is:

*Arachidonic Acid (AA) is the precursor of cerebrovascular active substances in the human body, and its metabolites are closely associated with the pathogenesis of cerebrovascular diseases. In recent years, the cytochrome P450 (CYP) metabolic pathway of AA has become a research hotspot. Furthermore, the CYP metabolic pathway of AA is regulated by soluble epoxide hydrolase (sEH). 1-trifluoromethoxyphenyl-3(1-propionylpiperidin-4-yl) urea (TPPU) is a novel sEH inhibitor that exerts cerebrovascular protective activity. This article reviews the mechanism of TPPU's protective effect on ischemic stroke disease.*
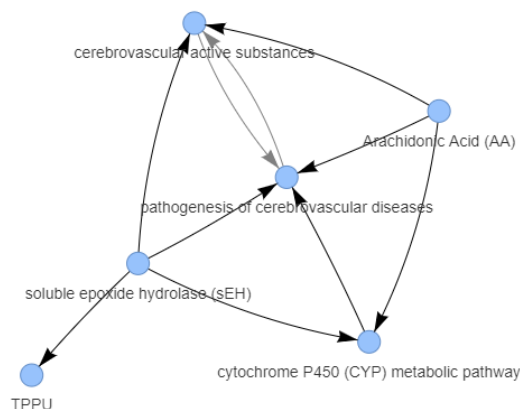
Figure 5.2: Causal graph from real text example.

The entities are extracted from the text using the GPT-based NER process described in Section 4.2.2, and these are:

[*'Arachidonic Acid (AA)', 'cerebrovascular active substances', 'pathogenesis of cerebrovascular diseases', 'cytochrome P450 (CYP) metabolic pathway', 'soluble epoxide hydrolase (sEH)', TPPU'*].

In the original abstract, these are:

*Arachidonic Acid (AA) is the precursor of cerebrovascular active substances in the human body, and its metabolites are closely associated with the pathogenesis of cerebrovascular diseases. In recent years, the cytochrome P450 (CYP) metabolic pathway of AA has become a research hotspot. Furthermore, the CYP metabolic pathway of AA is regulated by soluble epoxide hydrolase (sEH). 1-trifluoromethoxyphenyl-3(1-propionylpiperidin-4-yl) urea (TPPU) is a novel sEH inhibitor that exerts cerebrovascular protective activity. This article reviews the mechanism of TPPU's protective effect on ischemic stroke disease.*

The process then queries the GPT LLM to find the causal relationships between the variables, as previously discussed. The resulting relationships are encoded in the corresponding causal graph, as shown in Figure 5.2.

## 5.3 Limitations

LLMs are very powerful tools. Despite their impressive potential, it's crucial to recognize the limitations they introduce in the context of causal analysis and the identification of causal

graphs.

LLMs heavily depend on the textual data they are trained on, much of which is sourced from internet uploads. As a result, their performance can be influenced by biases and inaccuracies in the training data. Additionally, the language commonly used across the broader internet often features casual and colloquial language patterns that include causal terminology used in a way that differs from the precise and formal language found in medical academic literature. This can introduce further bias and lead the LLM to produce inaccurate or misleading answers. Furthermore, the complex nature of causal relationships, particularly in domains with intricate interdependencies like the medical field, may present challenges for LLMs to consistently capture the nuanced dynamics of causation [3].

LLMs can also be computationally expensive, resulting in potentially slow performance. This can be a limitation in some applications that demand real-time identification of causal graphs. When testing the process with actual medical abstracts, particularly those involving very large texts, the pipeline could take several hours to complete. The operational bottleneck was observed to be the causal query step, during which the LLM was queried about the causal relationship between all pairs of variables. Depending on the text's size and the resulting number of extracted entities, the volume of queries could extend to thousands of GPT requests. Given an average performance of several seconds per GPT API request, the cumulative execution time for the entire pipeline process could extend to hours.

By addressing these limitations and aligning our expectations with the LLMs capabilities, researchers can make informed decisions about using them as valuable tools in the pursuit of identifying intricate cause-and-effect relationships within complex textual data.

# Chapter 6

# Conclusion

This conclusion chapter summarizes the main findings of the study and discusses their implications. It also highlights the limitations of the study and suggests directions for future work.

## 6.1   Project Summary

This project investigated the capabilities of large language models, with a specific focus on their abilities in causal discovery and the identification of causal relationships.

The project's results, which show an average F1 score of 0.93 on the benchmarks, suggest that LLMs have the potential to significantly assist human efforts and contribute to tasks such as the identification of causal graphs.

The results also reveal that LLMs are not exempt from flaws or weaknesses. These imperfections highlight that the output and work of LLMs for building casual graphs should be verified by experts at this time. LLMs can be useful in extracting common knowledge from medical text, and when combined with expert insights, they offer the potential to efficiently generate more comprehensive causal graphs.

Large language models represent a remarkable advancement in artificial intelligence research. They are getting closer to human-level language capabilities than ever before [11]. LLMs represent a new and intriguing opportunity to extract common knowledge from medical literature and can help to complement and speed up causal analysis and causal graph identification. However, more research is essential to address and overcome the limitations of LLMs.

## 6.2   Future Work

Since this project focused primarily on identifying causal graphs, future efforts could prioritize the development of additional functionalities within the main causal analysis pipeline.

For example, the process could ensure the acyclicity of the generated graphs, another important property of DAGs. This might involve developing methods to detect and rectify incorrect or misplaced causal relationships (graph edges).

Another potential feature could involve labeling causal relationships (in the form of graph edges) with citations to the sources that the LLM referenced to establish the direction of causality between specific pairs of entities. This would aide human experts in assessing the validity and accuracy of LLM-generated answers by verifying the information against reliable and legitimate sources.

The findings and results of this project suggest several potential avenues for future development and research. As the field of large language models continues to evolve, exploring how these models can be further fine-tuned and tailored to domain-specific contexts could achieve even more accurate and robust causal analysis results. Investigating ways to mitigate the identified limitations of LLMs, such as improving the accuracy of extracting nuanced causal relationships or refining the verification process by experts, would help to maximize the potential of LLMs in causal graph identification.

STUDENTSUPSI

# Appendix A

# Implementation Details

This appendix section presents the source code used for this project.

## A.1  Scraping

**Search PubMed articles by terms**

```
1  def search_by_terms(terms, db, retmax, use_history):
2
3      terms_string = '+AND+'.join([s.strip().replace(' ', '+') for s in
       terms])
4      url_params = {
5              'db': db,
6              'term': terms_string,
7              'retmax': retmax,
8              'api_key': api_key,
9          }
10
11     if use_history:
12         url_params['usehistory'] = 'y'
13
14     url = f'{base_url}esearch.fcgi'
15     response = requests.get(url, params=url_params)
16     ids = re.findall(r"<Id>(\d+)</Id>", response.text)
17
18     if use_history:
19         web_match = re.search(r"<WebEnv>(\S+)</WebEnv>", response.text)
20         web = web_match.group(1) if web_match else None
21
22         key_match = re.search(r"<QueryKey>(\d+)</QueryKey>", response.
       text)
23         key = key_match.group(1) if key_match else None
24
25         return ids, web, key
```

```
26
27     return ids
```

**Get data of previously found articles**

```python
1  def get_articles_data(ids, web_env, query_key, db, retmax):
2
3      use_web_env = not ids
4
5      url_params = {
6          'db': db,
7          'rettype': 'abstract',
8          'retmode': 'xml',
9          'api_key': api_key,
10         'retmax': retmax,
11     }
12
13     if use_web_env:
14         url_params['query_key'] = query_key
15         url_params['WebEnv'] = web_env
16     else:
17         ids_string = ','.join(map(str, ids))
18         url_params['id'] = ids_string
19
20     url = f'{base_url}efetch.fcgi'
21     response = requests.get(url, params=url_params)
22
23     soup = BeautifulSoup(response.text, features="xml")
24     articles = soup.find_all('PubmedArticle')
25     if not articles:
26         print('ERROR: No articles found')
27         return None
28
29     data = pd.DataFrame(columns=['id', 'title', 'abstract', 'keywords', '
    pub_date'])
30     for article in articles:
31         article_data = {
32             'id': article.find('PMID').get_text(),
33             'title': article.find('ArticleTitle').get_text(),
34             'abstract': ' '.join([a.get_text() for a in article.find_all(
    'AbstractText')]),
35             'keywords': [[k.get_text() for k in article.find_all('Keyword
    ')]],
36         }
37         pub_date = article.find('PubMedPubDate', {'PubStatus': 'received'
    })
38         if pub_date:
```

```
39          article_data['pub_date'] = datetime.strptime(f"{pub_date.find
      ('Day').get_text()} {pub_date.find('Month').get_text()} {pub_date.find
      ('Year').get_text()}", "%d %m %Y")

40

41      data = pd.concat([data, pd.DataFrame(article_data)]).reset_index(
      drop=True)

42

43    return data
```

**Clean invalid data**

```
1  def clean_data(data, drop_id_duplicates, drop_empty_abstracts,
      drop_nan_abstracts, drop_date_nan, drop_date_before, drop_date_after,
      search_terms):

2

3      if drop_id_duplicates:
4          data = data.drop_duplicates(subset=['id'], inplace=False)
5      if drop_empty_abstracts:
6          data = data[data['abstract'] != '']
7      if drop_nan_abstracts:
8          data = data.dropna(subset=['abstract'])
9      if drop_abstracts_with_matches and drop_abstracts_matches:
10         data = data[~data['abstract'].str.startswith(tuple(
      drop_abstracts_matches))]

11

12     if drop_date_nan:
13         data = data.dropna(subset=['pub_date'])

14

15     if drop_date_before:
16         data = data[data['pub_date'] > drop_date_before]
17     if drop_date_after:
18         data = data[data['pub_date'] < drop_date_after]

19

20     if search_terms:
21         data['search_terms'] = [search_terms]*len(data)

22

23     return data.reset_index(drop=True)
```

## A.2   Causal Discovery

The complete GPT user messages can be found at Appendix B.

**Complete *causal_discovery_pipeline***

```python
1  def causal_discovery_pipeline(text_title, text, entities,
       reverse_edge_for_variable_check, optimize_found_entities,
       search_cycles):
2
3      if entities == []:
4          entities = gpt_ner(text)
5
6      if optimize_found_entities:
7          opt_entities = optimize_entities(entities, text)
8          entities = list(opt_entities.keys())
9
10     graph_edges = gpt_causal_discovery(entities, text)
11
12     edges = extract_edge_answers(graph_edges)
13
14     if reverse_edge_for_variable_check:
15         valid_edges, invalid_edges = check_invalid_answers(edges)
16
17         edge_correction_response = correct_invalid_edges(invalid_edges,
       text)
18         corrected_edges = extract_edge_answers(edge_correction_response)
19
20         valid_edges.extend(corrected_edges)
21         edges = valid_edges
22
23     nodes, processed_edges, bidirected_edges = preprocess_edges(edges)
24
25     cycles = []
26     if search_cycles:
27         cycles = find_cycles(nodes=nodes, edges=processed_edges)
28     build_graph(nodes, processed_edges, bidirected_edges, cycles)
29
30
31     if verbose:
32         print_edges(graph_edges)
33
34     return nodes, processed_edges + bidirected_edges, cycles
```

### Base GPT API request helper function

```python
1  def gpt_request(system_msg, user_msg, model, temperature):
2      if not system_msg or not user_msg:
3          return None
4      try:
5          response = openai.ChatCompletion.create(
6                          model=model,
7                          messages=[
```

```
8                                    {"role": "system", "content": system_msg},
9                                    {"role": "user", "content": user_msg}],
10                              temperature=temperature)
11
12          return response.choices[0].message.content
13      except:
14          return None
```

## Entity extraction function

```
1  def gpt_ner(text):
2
3      system_msg = 'You are a helpful assistant for Named Entity
       Recognition of medical texts.'
4
5      # see Appendix B.1 for user message
6      # user_msg = ''
7
8      response = gpt_request(system_msg, user_msg)
9      if not response:
10          return []
11
12      answer_text = response
13
14      soup = BeautifulSoup(answer_text, 'html.parser')
15      entities = [entity.text for entity in soup.find_all('entity')]
16
17      return entities
```

## Entity optimization function

```
1  def optimize_entities(entities, text=None):
2      system_msg = 'You are a helpful assistant for medical entity
       optimization, by accurately identifying synonyms, redundant entities,
       or entities that can be used interchangeably'
3
4      entities_text = '\n'.join([f'<Entity>{entity}</Entity>' for entity in
        entities])
5
6      # see Appendix B.2 for user message
7      # user_msg = ''
8
9      response = gpt_request(system_msg, user_msg)
10      if response:
11          soup = BeautifulSoup(response, 'html.parser')
12          answer = soup.find('answer').text
```

```
13        try:
14            opt_entities = json.loads(answer)
15            if opt_entities:
16                return opt_entities
17        except (json.JSONDecodeError, TypeError):
18            pass
19
20    return entities
```

### Causal discovery function

```
1  def gpt_causal_discovery(entities, text, use_pretrained_knowledge,
   reverse_variable_check):
2
3    graph_edges = []
4
5    system_msg = 'You are a helpful assistant for causal reasoning and
   cause-and-effect relationship discovery.'
6
7    for i1, e1 in enumerate(entities):
8        for i2, e2 in enumerate(entities):
9            if i1 == i2 or (not reverse_variable_check and i1 >= i2):
10                continue
11
12
13            # see Appendix B.3 for user message
14            # user_msg = ''
15
16            response = gpt_request(system_msg, user_msg)
17            if response:
18                graph_edges.append(((e1, e2), response))
19
20    return graph_edges
```

### Compatibility check for GPT answers

```
1  def check_edge_compatibility(answer1, answer2):
2    return (arrows[answer1], arrows[answer2]) in [(forward_arrow,
   backward_arrow), (backward_arrow, forward_arrow), (no_arrow, no_arrow)
   , (bidirectional_arrow, bidirectional_arrow)]
```

### Find invalid GPT answers, when using double variable check

```
1  def check_invalid_answers(directed_edges):
2    invalid_edges = []
```

```
3    valid_edges = []
4    temp_edges = []
5    answers = {}
6    for (n1, n2), answer in directed_edges:
7
8        if (n1, n2) not in temp_edges and (n2, n1) not in temp_edges:
9            temp_edges.append((n1, n2))
10           answers[(n1, n2)] = answer
11       elif (n1, n2) in temp_edges:
12           if answers[(n1, n2)] != answer:
13               invalid_edges.append((((n1, n2), answer), ((n2, n1),
    answers[(n2, n1)])))
14           else:
15               valid_edges.append(((n1, n2), answer))
16
17           temp_edges.remove((n1, n2))
18       elif (n2, n1) in temp_edges:
19           if check_edge_compatibility(answers[(n2, n1)], answer):
20               valid_edges.append(((n1, n2), answer))
21           else:
22               invalid_edges.append((((n1, n2), answer), ((n2, n1),
    answers[(n2, n1)])))
23
24           temp_edges.remove((n2, n1))
25
26   for n1, n2 in temp_edges:
27       if (n1, n2) not in invalid_edges:
28           invalid_edges.append((((n1, n2), answer), ((n2, n1), answers
    [(n2, n1)])))
29
30   return valid_edges, invalid_edges
```

## Return GPT answers in textual form, when re-querying for incoherent answers

```
1  def get_textual_answers(e1, e2, ans):
2      if ans == forward_arrow_answer:
3          return f'"{e1}" causes "{e2}"'
4      elif ans == backward_arrow_answer:
5          return f'"{e2}" causes "{e1}"'
6      elif ans == no_arrow_answer:
7          return f'"{e1}" and "{e2}" are not causally related'
8      elif ans == bidirectional_arrow_answer:
9          return f'there is a common factor that is the cause for both "{e1
    }" and "{e2}"'
10     else:
11         return None
```

**Correct incoherent GPT answers**

```python
def correct_invalid_edges(invalid_edges, text, use_pretrained_knowledge):
    graph_edges = []

    if not invalid_edges:
        return []

    system_msg = 'You are a helpful assistant for causal reasoning and
    cause-and-effect relationship discovery.'

    for ((e1, e2), answer1), ((e3, e4), answer2) in invalid_edges:

        # see Appendix B.4 for user message
        # user_msg = ''

        response = gpt_request(system_msg, user_msg)
        if response:
            graph_edges.append(((e1, e2), response))

    return graph_edges
```

## A.3   Plotting the Causal Graph

**Decode GPT answers in directed edges**

```python
def preprocess_edges(edges):
    nodes = set()
    directed_edges = []
    bidirected_edges = []

    for (n1, n2), answer in edges:

        nodes.add(n1)
        nodes.add(n2)

        direction = normalize_edge_direction(n1, n2, answer)
        if direction:
            if len(direction) == 2:
                bidirected_edges.extend(direction)
            else:
                directed_edges.extend(direction)

    return list(nodes), directed_edges, bidirected_edges
```

### Return edge in normalized form (A → B)

```
1  def normalize_edge_direction(e1, e2, answer):
2      if answer in arrows:
3          if arrows[answer] == forward_arrow:
4              return [(e1, e2)]
5          elif arrows[answer] == backward_arrow:
6              return [(e2, e1)]
7          elif arrows[answer] == bidirectional_arrow:
8              return [(e2, e1), (e1, e2)]
9      return None
```

### Find cycles in causal graph

```
1  def find_cycles(nodes, edges):
2      if not nodes or not edges:
3          return []
4
5      g = gt.Graph(directed=True)
6
7      nodes_ids = {}
8      v_prop = g.new_vertex_property("string")
9      for n in nodes:
10         v = g.add_vertex()
11         v_prop[v] = n
12         nodes_ids[n] = v
13
14     for (n1, n2) in edges:
15         e = g.add_edge(nodes_ids[n1], nodes_ids[n2])
16
17     cycles = []
18     for i, c in enumerate(gt.all_circuits(g)):
19         if i >= 100:
20             break
21         cycles.append([v_prop[v] for v in c])
22
23     return cycles
```

### Build resulting causal graph

```
1  def build_graph(nodes, edges, bidirected_edges, cycles, plot_static_graph
   , directory_name, graph_name):
2
3      if plot_static_graph:
4          plt.figure()
5      G = nx.DiGraph()
```

```
6
7       G.add_nodes_from(nodes)
8
9       for e1, e2 in edges:
10          G.add_edge(e1, e2, color='black', style='solid')
11
12      for cycle in cycles:
13          for i in range(len(cycle) - 1):
14              G[cycle[i]][cycle[i + 1]]['color'] = 'red'
15          G[cycle[-1]][cycle[0]]['color'] = 'red'
16
17      for e1, e2 in bidirected_edges:
18          G.add_edge(e1, e2, color='grey', style='dashed')
19
20      if plot_static_graph:
21          pos = nx.spring_layout(G)
22          nx.draw_networkx_nodes(G, pos)
23          nx.draw_networkx_labels(G, pos)
24
25          edge_colors = [G.edges[edge]['color'] for edge in G.edges()]
26          edge_styles = [G.edges[edge]['style'] for edge in G.edges()]
27
28          nx.draw(G, pos, node_color='skyblue', node_size=1500,
29                  font_size=10, font_weight='bold', arrowsize=20,
        edge_color=edge_colors, style=edge_styles,
30                  width=2)
31          plt.title(graph_name)
32          plt.show()
33
34      net = Network(directed=True, notebook=True)
35      net.from_nx(G)
36      net.force_atlas_2based()
37      net.show_buttons(filter_=['physics'])
38      os.makedirs(directory_name, exist_ok=True)
39      net.save_graph(f'{directory_name}/{graph_name}.html')
```

Causal Graph Identification by Large Language Models

# Appendix B

# GPT User Messages

The following section contains examples of all user messages used within the project. These were constructed with an iterative prompt development approach [19].

## B.1   Named Entity Recognition (NER)

```
1    You will be provided with an abstract of a medical research paper
     delimited by the <Text></Text> xml tags.
2    Please read the provided abstract carefully to comprehend the context
      and content. Analyze the provided text and identify all the
     meaningful entities that could contribute to understanding cause-and-
     effect relationships between factors such as diseases, medications,
     treatments, interventions, symptoms, outcomes, effects, or risk
     factors.
3
4    Avoid including entities that are synonyms or can be used
     interchangeably to already identified ones. For example, if text
     contains both "hospital" and "medical center" (which are synonyms and
     can be used interchangeably) and you already extracted "hospital" as
     final entity, do not include "medical center" as well.
5
6    Your response should highlight entities that are crucial for
     establishing causal relationships in the medical context.
7
8    Answer listing only the found entities within the tags <Answer><
     Entity>[entity]</Entity><Entity>[entity]</Entity></Answer>
9    (e.g., <Answer><Entity>diabetes</Entity><Entity>hypertension</Entity
     ></Answer>)
10
11   Text:
12       <Text>{text}</Text>
```

## B.2  Entity Optimization

```
1    You will be provided with an abstract of a medical research paper
     delimited by the <Text></Text> xml tags, and a list of named entities
     representing medical entities, each one of them delimited by the <
     Entity></Entity> xml tags.
2
3    Text:
4            <Text>{text}</Text>
5
6    Entities:
7            {entities}
8
9    Your task is to optimize this entity list by identifying synonyms
     within the entities and grouping them accordingly.
10   Your goal is to create a JSON object where the keys represent the
     root word entities, and each key is associated with an array of its
     synonyms, i.e., words or entities that can be used interchangeably to
     the root word.
11   If a root word entity has no synonyms, its value in the JSON should
     be an empty array.
12
13   Ensure that each entity appears only once in the dictionary, either
     as key (i.e. root word) or as element in the value arrays (i.e. the
     synonyms): an entity must not appear as key if it is the synonym (i.e.
      in the value array) of another entity, and the other way around (i.e.
      must not be in the value array of an entity if it is already a key of
      the JSON object).
14   An entity must not be a synonym of itself.
15
16   You should efficiently process the given list of entities and produce
      the desired dictionary structure.
17   The output JSON object should accurately represent the optimized
     entities and their synonyms based on the provided list.
18
19   Then provide your final JSON object answer within the tags <Answer>[
     json_obj]</Answer>, (e.g. <Answer>
20       {{
21           "smoking": ["tobacco", "nicotine", "cigarettes", "cigar"],
22           "cancer": ["lung cancer"],
23           "tumors": []
24       }}
25       </Answer>
26   ).
27
28   Follow the example below to understand the expected output.
29
```

```
30      Example:
31
32      Given the initial list of entities:
33      <Entity >smoking </Entity >
34      <Entity >lung cancer </Entity >
35      <Entity >tumors </Entity >
36      <Entity >cancer </Entity >
37      <Entity >tobacco </Entity >
38      <Entity >nicotine </Entity >
39      <Entity >cigarettes </Entity >
40      <Entity >cigar </Entity >
41
42      You should pair the synonyms , generate the following JSON object and
        provide it as your answer:
43      <Answer >
44      {{
45          "smoking": ["tobacco", "nicotine", "cigarettes", "cigar"],
46          "cancer": ["lung cancer"],
47          "tumors": []
48      }}
49      </Answer >
50
51      Note that every entity appears only once in the output JSON object ,
        either as key or as element in the value arrays.
52
53      After you have finished building the JSON object , check and make sure
         that every entity appears only once in the output JSON object , either
         as key or as element in the value arrays.
```

## B.3   Causal Discovery

The *random_causal_verb* variable is randomly selected from the available causation verbs, which include: 'provokes', 'triggers', 'causes', 'leads to', 'induces', 'results in', 'generates', 'produces', 'stimulates', 'instigates', 'engenders', 'promotes', 'gives rise to', 'sparks' [1].

```
1       You will be provided with an abstract of a medical research paper
        delimited by the <Text ></Text > xml tags , and a pair of entities
        extracted from the given abstract delimited by the <Entity ></Entity >
        xml tags representing medical entities , such as medications ,
        treatments , symptoms , diseases , outcomes , side effects , or other
        medical factors.
2
3           Text:
4                   <Text >{text}</Text >
5
6           Entities:
7                   <Entity >{entity1}</Entity >
```

```
8                       <Entity >{entity2}</Entity >
9

10   Please read the provided abstract carefully to comprehend the context
     and content. Examine the roles, interactions, and details surrounding
     the entities within the abstract.

11   Based on the information in the text and your pretrained knowledge,
     determine the most likely cause-and-effect relationship between the
     entities from the following listed options (A, B, C, D):

12
13       Options:
14               A: "{entity1}" {random_causal_verb} "{entity2}";
15               B: "{entity2}" {random_causal_verb} "{entity1}";
16               C: "{entity1}" and "{entity2}" are not directly causally
     related;
17               D: there is a common factor that is the cause for both "{
     entity1}" and "{entity2}";

18
19   Your response should analyze the situation in a step-by-step manner,
     ensuring the correctness of the ultimate conclusion, which should
     accurately reflect the likely causal connection between the two
     entities based on the information presented in the text and any
     additional knowledge you are aware of.

20   If no clear causal relationship is apparent, select the appropriate
     option accordingly.

21
22   Then provide your final answer within the tags <Answer >[answer]</
     Answer >, (e.g. <Answer >C</Answer >).
```

## B.4   Correct Incoherent GPT Answers

The *get_textual_answers* function returns a textual version of the causal relationship be-
tween the two entities, in the form of "*'A' causes 'B'*". For *random_causal_verb* see Appendix
B.3.

```
1    You will be provided with an abstract of a medical research paper
     delimited by the <Text ></Text > xml tags, and a pair of entities
     extracted from the given abstract delimited by the <Entity ></Entity >
     xml tags representing medical entities (such as medications,
     treatments, symptoms, diseases, outcomes, side effects, or other
     medical factors), and two answers you previously gave to this same
     request that are incoherent with each other, delimited by the <Answer
     ></Answer > xml tags.
2        Text:
3                <Text >{text}</Text >
4
5        Entities:
6                <Entity >{entity1}</Entity >
```

```
 7                    <Entity >{entity2}</Entity >
 8
 9          Previous  incoherent  answers:
10                    <Answer >{get_textual_answers (entity1 , entity2 ,
     answer1)}</Answer >
11                    <Answer >{get_textual_answers (entity2 , entity1 ,
     answer2)}</Answer >
12
13   Please read  the provided  abstract  carefully  to comprehend  the context
      and content.
14   Consider  the previous  answers  you gave  to this  same  request  that are
      incoherent  with each other , and the entities  they refer  to in order  to
      give  a correct  answer.  Examine  the roles , interactions , and details
      surrounding  the entities  within  the abstract.
15   Based  on the information  in the text  and your  pretrained  knowledge ,
     determine  the most  likely  cause -and-effect  relationship  between  the
     entities  from  the following  listed  options  (A, B, C, D):
16
17          Options:
18              A: "{entity1}"  {random_causal_verb}  "{entity2}";
19              B: "{entity2}"  {random_causal_verb}  "{entity1}";
20              C: "{entity1}"  and "{entity2}"  are not  directly  causally
     related;
21              D: there  is a common  factor  that is the cause  for both  "{
     entity1}"  and "{entity2}";
22
23   Your  response  should  accurately  reflect  the likely  causal  connection
     between  the two entities  based  on the information  presented  in the
     text  and any additional  knowledge  you are aware  of.
24    If no clear  causal  relationship  is apparent , select  the appropriate
     option  accordingly.
25   Then  provide  your  final  answer  within  the tags  <Answer >[answer]</
     Answer >
26   (e.g.  <Answer >C</Answer >).
```

Causal Graph Identification by Large Language Models

# Appendix C

# Benchmarks

This appendix section presents the benchmarks in more detail. It shows the ground truth graphs and, using the evaluation metrics discussed at Section 5.1.3, presents the additional results from the different algorithms: the random baseline, the *GPT-3.5-turbo* and *GPT-4* models.

## C.1   Ground Truth Graphs

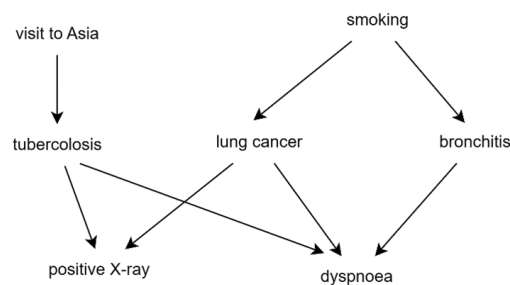The ground truth graphs for the benchmarks are shown at Figures 5.1, C.1, C.2, C.3, C.4, E.9.



Figure C.1: 'Asia' benchmark.

Figure C.2: 'Smoking' benchmark.
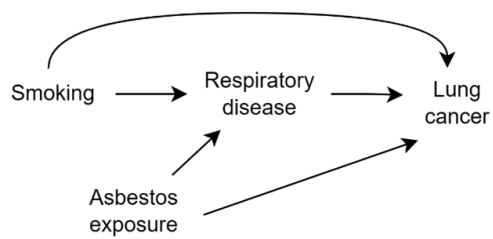


Figure C.3: 'Alcohol' benchmark.
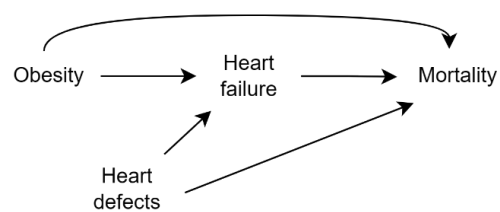


Figure C.4: 'Cancer' benchmark.



Figure C.5: 'Obesity' benchmark.

## C.2   Additional Benchmark Results

The results of the various algorithms applied to the benchmarks are presented in Tables C.1 (baseline), C.2 (*GPT-3.5-turbo*), and C.3 (*GPT-4*).

| Benchmark Name | Missing edges | Extra edges | SHD | Correct edge direction | Incorrect edge direction | Precision | Recall | F1 score | PRC Area |
|---|---|---|---|---|---|---|---|---|---|
| **Asia** | 6 | 15 | 21 | 2 | 15 | 0.12 | 0.25 | 0.16 | 0.25 |
| **Smoking** | 4 | 1 | 5 | 2 | 2 | 0.75 | 0.43 | 0.55 | 0.71 |
| **Alcohol** | 3 | 1 | 4 | 0 | 1 | 0 | 0 | NaN | 0.17 |
| **Cancer** | 2 | 4 | 6 | 2 | 5 | 0.43 | 0.6 | 0.5 | 0.58 |
| **Diabetes** | 2 | 5 | 7 | 1 | 7 | 0.38 | 0.6 | 0.46 | 0.55 |
| **Obesity** | 3 | 4 | 7 | 1 | 5 | 0.33 | 0.4 | 0.36 | 0.46 |
| AVG | 3.33 | 5 | 8.33 | 1.33 | 5.83 | 0.33 | 0.38 | 0.36 | 0.45 |

Table C.1: Benchmark results with random baseline.

| Benchmark Name | Missing edges | Extra edges | SHD | Correct edge direction | Incorrect edge direction | Precision | Recall | F1 score | PRC Area |
|---|---|---|---|---|---|---|---|---|---|
| **Asia** | 3 | 5 | 8 | 5 | 5 | 0.5 | 0.63 | 0.56 | 0.60 |
| **Smoking** | 3 | 0 | 3 | 4 | 0 | 1 | 0.57 | 0.73 | 0.88 |
| **Alcohol** | 2 | 0 | 2 | 1 | 0 | 1 | 0.33 | 0.5 | 0.78 |
| **Cancer** | 0 | 1 | 1 | 4 | 2 | 0.83 | 1 | 0.91 | 0.92 |
| **Diabetes** | 2 | 4 | 6 | 1 | 6 | 0.43 | 0.6 | 0.5 | 0.58 |
| **Obesity** | 2 | 3 | 5 | 2 | 4 | 0.5 | 0.6 | 0.55 | 0.61 |
| AVG | 2 | 2.17 | 4.17 | 2.83 | 2.83 | 0.71 | 0.62 | 0.66 | 0.73 |

Table C.2: Benchmark results with *GPT-3.5-turbo* model.

| Benchmark Name | Missing edges | Extra edges | SHD | Correct edge direction | Incorrect edge direction | Precision | Recall | F1 score | PRC Area |
|---|---|---|---|---|---|---|---|---|---|
| **Asia** | 0 | 8 | 8 | 8 | 8 | 0.5 | 1 | 0.67 | 0.75 |
| **Smoking** | 1 | 0 | 1 | 5 | 1 | 1 | 0.86 | 0.92 | 0.96 |
| **Alcohol** | 0 | 0 | 0 | 3 | 0 | 1 | 1 | 1 | 1 |
| **Cancer** | 0 | 1 | 1 | 4 | 1 | 0.83 | 1 | 0.91 | 0.92 |
| **Diabetes** | 0 | 0 | 0 | 5 | 0 | 1 | 1 | 1 | 1 |
| **Obesity** | 0 | 0 | 0 | 5 | 0 | 1 | 1 | 1 | 1 |
| AVG | 0.17 | 1.5 | 1.67 | 5 | 1.67 | 0.89 | 0.98 | 0.93 | 0.94 |

Table C.3: Benchmark results with *GPT-4* model.

# Appendix D

# Command Line Options

This appendix section presents all command line options for running the whole process using the previously presented *causal_llms.py* script.

```
 1  Usage: causal_llms.py <action> [options]
 2
 3  Description:
 4    This script performs various tasks related to causal discovery.
 5
 6  Actions:
 7    ex        Run the example test.
 8    s         Run the scraping process.
 9    c         Perform causal analysis.
10    sc        Run scraping and causal analysis.
11    b         Run the benchmark tests.
12
13  Options:
14    --help        Show this help message and exit.
15    --data-path   Path of data file for causal analysis.
16    --algorithm   Algorithm to use with benchmarks.
17
18
19  Examples:
20    python causal_llms.py ex
21    python causal_llms.py s
22    python causal_llms.py c --data-path </path/to/data>
23    python causal_llms.py sc
24    python causal_llms.py b --algorithm {b|gpt}
25
26
27  The 'algorithm' parameter sets the algorithm for the benchmark tests.
28  The possible values are:
29  * 'b': Baseline algorithm
30  * 'gpt': GPT LLM
```

# Appendix E

# Examples with Real Medical Texts

This section presents examples of the causal discovery process being applied to real-world abstracts, extracted from PubMed [40, 41, 42, 43, 44, 45, 46, 47, 48].

**Risk stratification of sudden cardiac death: a review. [40]**

*Sudden cardiac death (SCD) is responsible for several millions of deaths every year and remains a major health problem. To reduce this burden, diagnosing and identification of high-risk individuals and disease-specific risk stratification are essential. Treatment strategies include treatment of the underlying disease with lifestyle advice and drugs and decisions to implant a primary prevention implantable cardioverter-defibrillator (ICD) and perform ablation of the ventricles and novel treatment modalities such as left cardiac sympathetic denervation in rare specific primary electric diseases such as long QT syndrome and catecholaminergic polymorphic ventricular tachycardia. This review summarizes the current knowledge on SCD risk according to underlying heart disease and discusses the future of SCD prevention.*
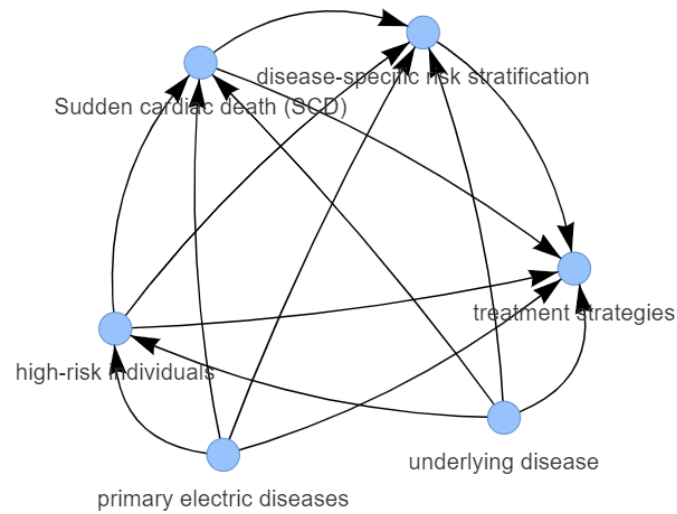
Figure E.1: "Risk stratification of sudden cardiac death: a review" graph.

# Dexamethasone for delayed edema after intracerebral hemorrhage : To be or not to be? [41]

*The pathogenesis of delayed cerebral edema after intracerebral hemorrhage is still unclear. In this case report, we speculate that the formation of subdural effusion or hemorrhage is associated with delayed cerebral edema. By referring to the treatment plan of chronic subdural hematoma, adding dexamethasone to routine medication, certain therapeutic effect has been achieved. Dexamethasone may maintain the stability of blood-brain barrier by directly increasing the expression of ZO-1, and reduce the neuroinflammatory response caused by NF-B pathway by upregulating KLF2 expression, ultimately reducing nerve injury through multiple pathways.*

Figure E.2: "Dexamethasone for delayed edema after intracerebral hemorrhage : To be or not to be?" graph.

## Papillary Glioneuronal Tumor Masquerading as Malignant Brain Tumors: A Case Report. [42]

*Papillary glioneuronal tumor (PGNT) is a low-grade biphasic tumor that is composed of glial fibrillary acidic protein (GFAP)-positive glial cells and synaptophysin-positive neurons. We report a case of PGNT occurring in the right occipital lobe of a 48-year-old woman who presented with acute headache and left homonymous hemianopsia, the latter of which was difficult to distinguish from malignant brain tumors because of peritumoral brain edema, intratumoral hemorrhage, and intraoperative fluorescence staining. PGNT should be included as one of the differential diagnoses in cases where the tumor shows hemorrhagic change despite decreased perfusion in arterial spin labeling MRI.*
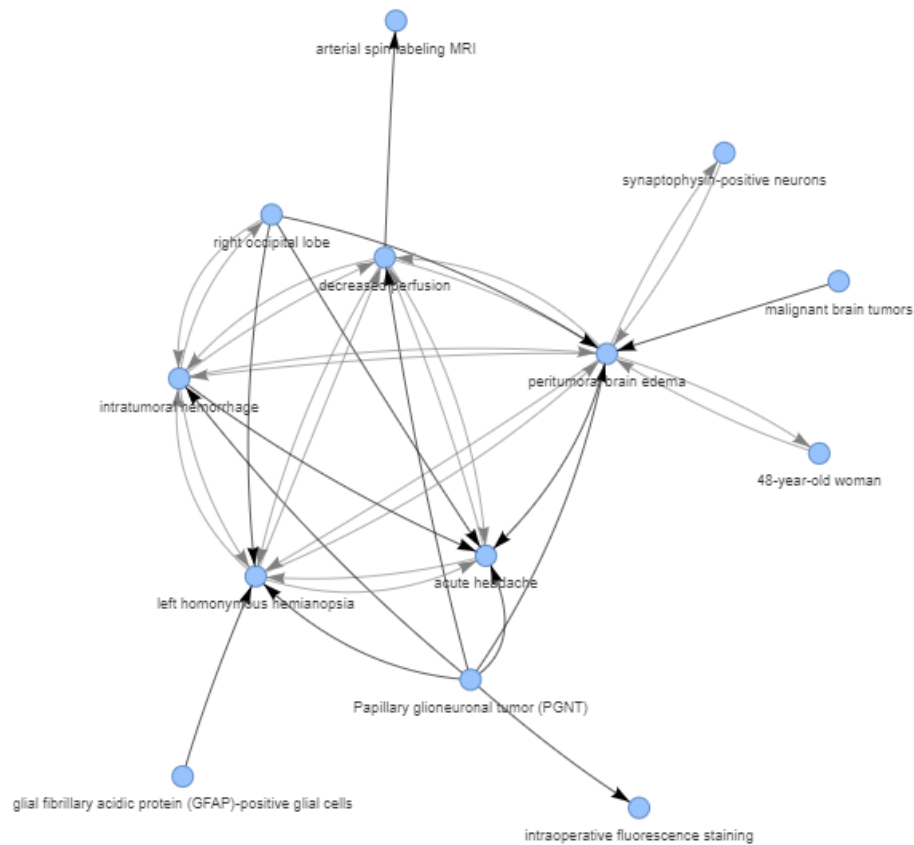
Figure E.3: "Papillary Glioneuronal Tumor Masquerading as Malignant Brain Tumors: A Case Report." graph.

## Gliosarcoma Invading the Temporal Bone, Temporalis Muscle, and Skull Base. [43]

*Gliosarcoma (GS) is a primary central nervous system tumor. It is an unusual type of glioblastoma multiforme (GBM) and rarely invades the skull base. It has a biomorphic tissue pattern with rapid alternation zones of glial and mesenchymal differentiation. We report the case of a 62-year-old male who presented with a one-month history of unsteady gait associated with dizziness. Brain MRI showed a right temporal mass that invaded the skull base with perilesional edema and a significant mass effect on the right lateral ventricle. The patient underwent a right-sided frontotemporal craniotomy with gross total resection. The pathology confirmed the diagnosis of GS. Postoperatively, the patient had an uneventful recovery with no complications and was discharged two days post-surgery.*
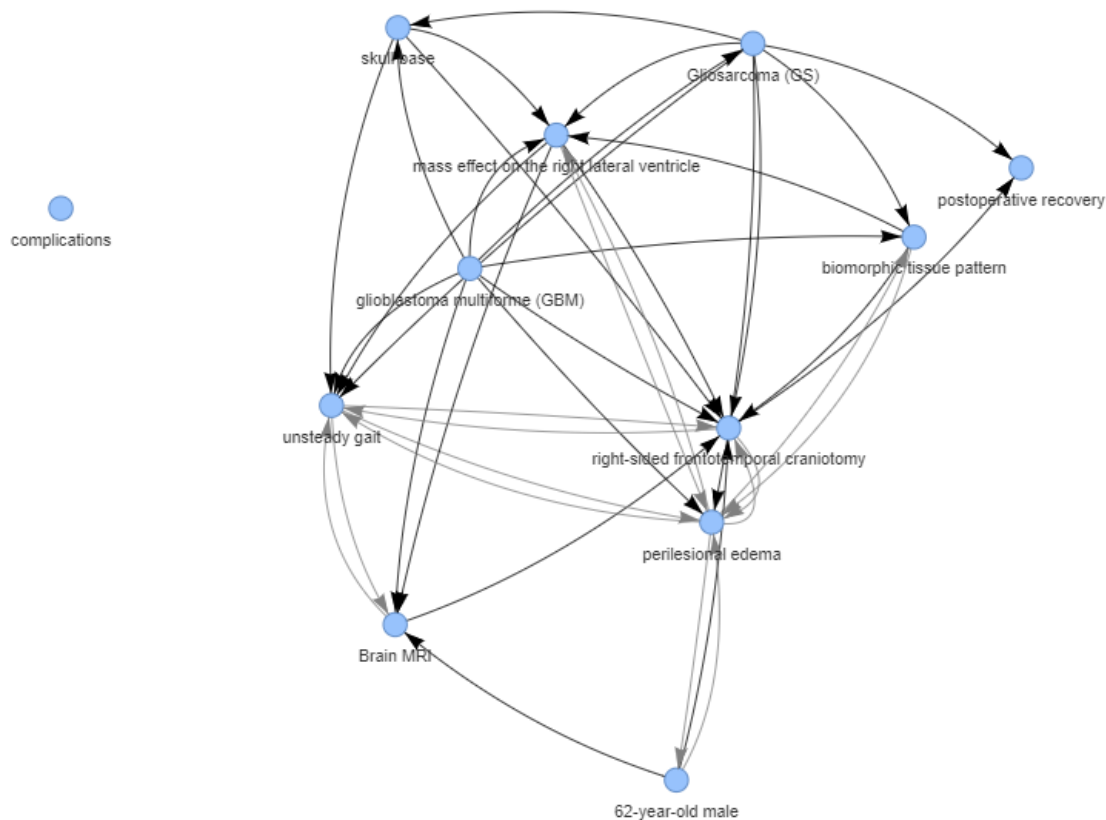
Figure E.4: "Gliosarcoma Invading the Temporal Bone, Temporalis Muscle, and Skull Base." graph.

## A Rare Case of Nontuberculous Mycobacterial Abscess Mimicking Brain Tumor in an Immunocompetent Patient. [44]

*Nontuberculous mycobacteria (NTM) is a type of bacteria that typically infects the pulmonary system, and NTM-central nervous system (CNS) infection, which occurs in the brain, is a very rare disease. A 64-year-old female patient presented with seizures as the main symptom and was found to have a mass of less than 1 cm in the right temporal lobe with accompanying edema. Although diseases such as tumor metastasis and parasitic cyst were suspected, the patient underwent a surgical resection, and NTM-CNS infection with abscess was diagnosed through biopsy. Antibiotic treatment was initiated after surgery, and the patient has been followed up without any significant symptoms. In this report, we review a rare case of NTM-CNS infection and discuss the understanding and treatment of this disease.*
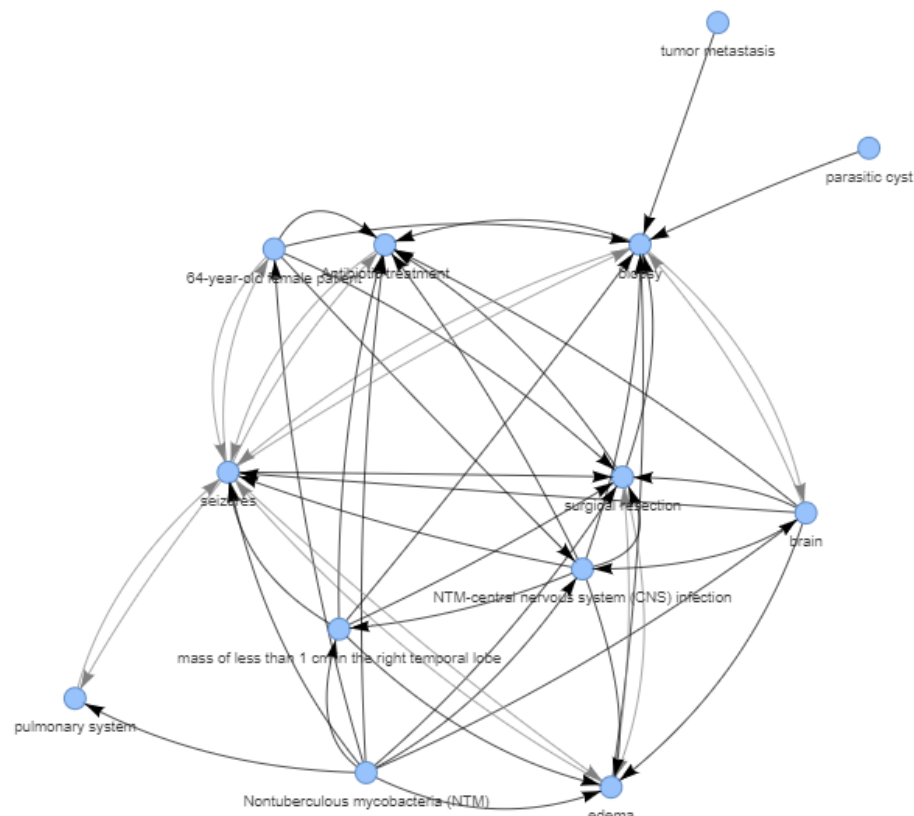
Figure E.5: "A Rare Case of Nontuberculous Mycobacterial Abscess Mimicking Brain Tumor in an Immunocompetent Patient." graph.

## Dramatic Response to Anti-IL-6 Receptor Therapy in Children With Life-Threatening Myelin Oligodendrocyte Glycoprotein-Associated Disease. [45]

*Myelin oligodendrocyte glycoprotein antibody-associated disease (MOGAD) is an immune-mediated neuroinflammatory disorder leading to demyelination of the CNS. Interleukin (IL)-6 receptor blockade is under study in relapsing MOGAD as a preventative strategy, but little is known about the role of such treatment for acute MOGAD attacks. We discuss the cases of a 7-year-old boy and a 15-year-old adolescent boy with severe acute CNS demyelination and malignant cerebral edema with early brain herniation associated with clearly positive serum titers of MOG-IgG, whose symptoms were incompletely responsive to standard acute therapies (high-dose steroids, IV immunoglobulins (IVIGs), and therapeutic plasma exchange). Both boys*

*improved quickly with IL-6 receptor inhibition, administered as tocilizumab.*
*Both patients have experienced remarkable neurologic recovery. We propose*
*that IL-6 receptor therapies might also be considered in acute severe*
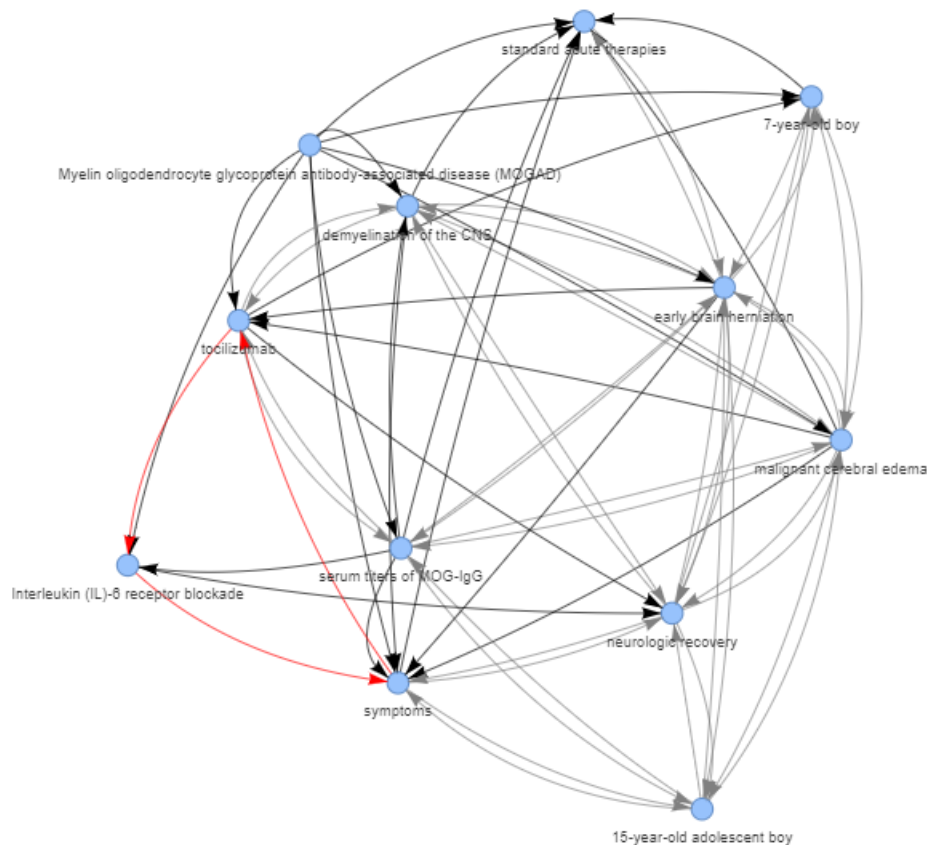*life-threatening presentations of MOGAD.*



Figure E.6: "Dramatic Response to Anti-IL-6 Receptor Therapy in Children With Life-Threatening Myelin Oligodendrocyte Glycoprotein-Associated Disease." graph.

## Role of echocardiography in patients treated with immune checkpoints inhibitors. [46]

*Immune-related adverse events occurring in the heart (cardiac immune-related adverse events; irAEs) by immune checkpoint inhibitors (ICIs) include myocarditis, arrhythmia, conduction disturbance, pericardial diseases, and takotsubo cardiomyopathy. Cardiac irAEs are rare but life-threatening. In cardio-oncology, the study of cardiac disorders caused by cancer treatment has recently attracted attention, and such studies may elucidate the*

*pathophysiology of cardiac irAEs and contribute to management strategies. This review discusses the pathogenic mechanisms underlying cardiac irAEs and the role of echocardiography in patients treated with ICIs.*
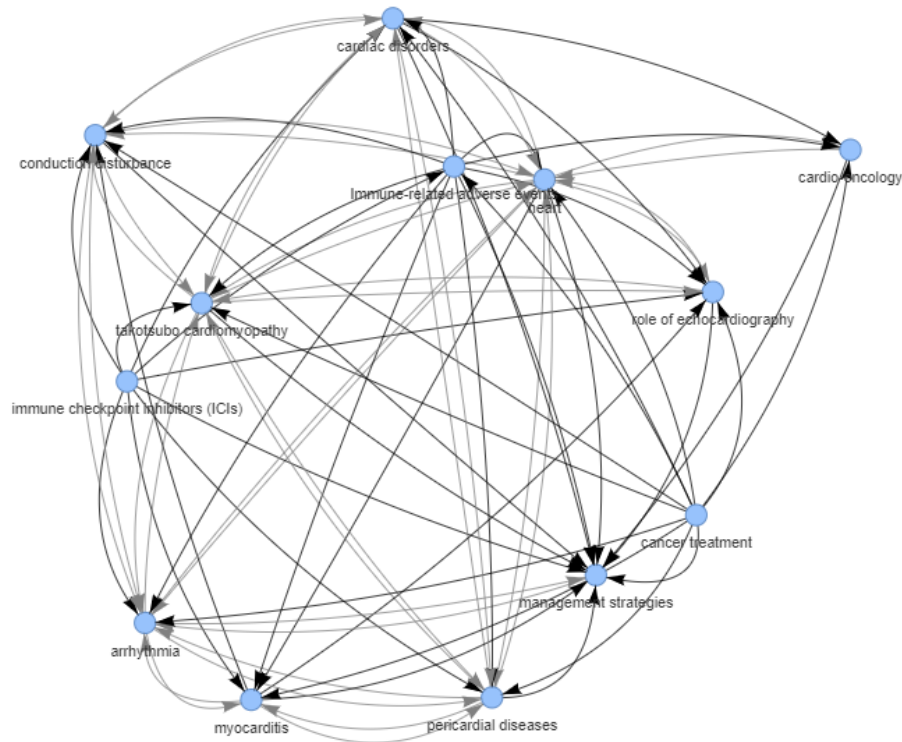


Figure E.7: "Role of echocardiography in patients treated with immune checkpoints inhibitors." graph.

## Advancing the science of management of arrhythmic disease in children and adult congenital heart disease patients within the last 25 years. [47]

*This review article reflects how publications in EP Europace have contributed to advancing the science of management of arrhythmic disease in children and adult patients with congenital heart disease within the last 25 years. A special focus is directed to congenital atrioventricular (AV) block, the use of pacemakers, cardiac resynchronization therapy devices, and implantable cardioverter defibrillators in the young with and without congenital heart disease, Wolff-Parkinson-White syndrome, mapping and ablation technology, and understanding of cardiac genomics to untangle arrhythmic sudden death in the young.*
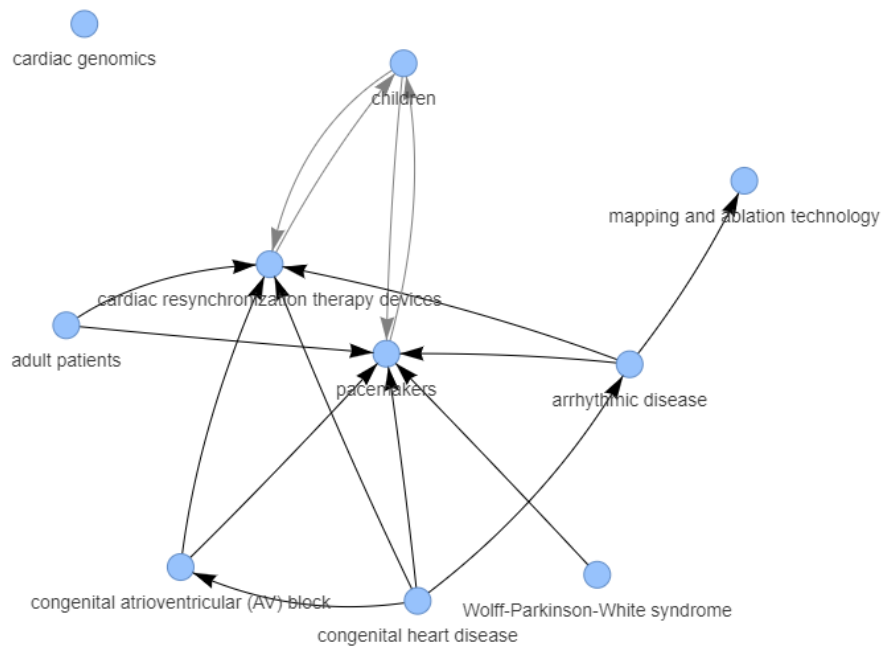
Figure E.8: "Advancing the science of management of arrhythmic disease in children and adult congenital heart disease patients within the last 25 years." graph.

## Left Atrial Thrombus in the Setting of Mitral Stenosis. [48]

*A 56-year-old man with no significant past medical history presented with exertional shortness of breath. Echocardiogram, cardiac magnetic resonance, and computed tomography showed mitral stenosis and a left atrial thrombus. Left atrial thrombus formation is a well-known complication of severe mitral stenosis that can lead to systemic thromboembolism. The patient underwent mitral valve replacement, left atrial thrombus resection, and left atrial appendage closure that resulted in significant improvement in breathing.*
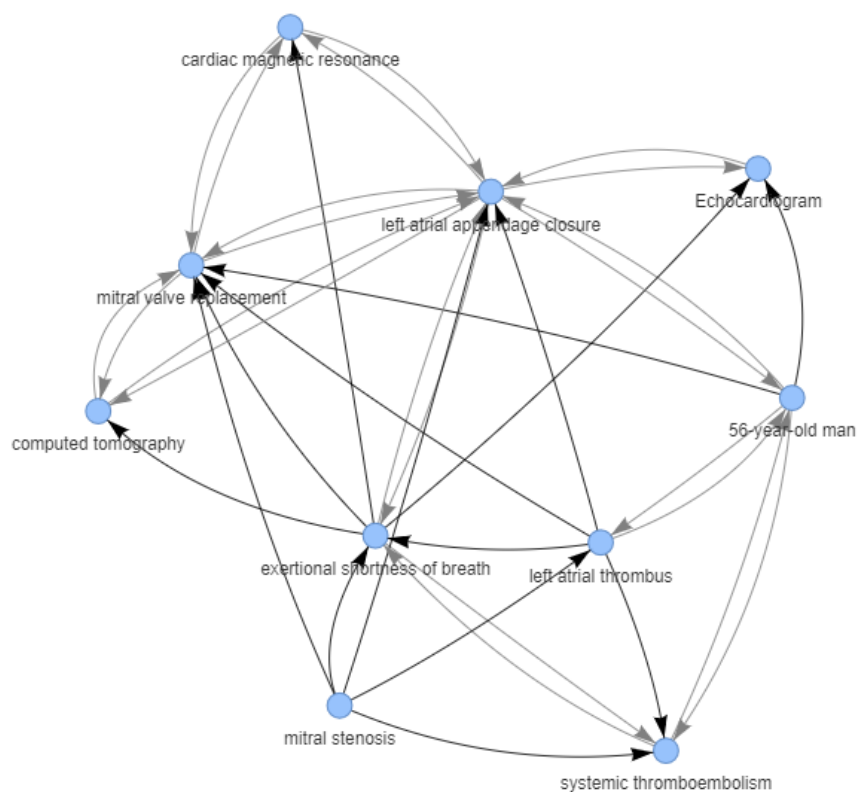
Figure E.9: "Left Atrial Thrombus in the Setting of Mitral Stenosis." graph.

# Bibliography

[1] S. Long et al. "Causal Discovery with Language Models as Imperfect Experts". In: *arXiv e-prints*, arXiv:2307.02390 (July 2023).

[2] J. Pearl, M. Glymour, and N. P. Jewell. *Causal inference in statistics: A primer*. John Wiley & Sons, 2016.

[3] S. Long et al. "Can large language models build causal graphs?" In: *arXiv e-prints*, arXiv:2303.05279 (Mar. 2023).

[4] J. Pearl. "Theoretical Impediments to Machine Learning With Seven Sparks from the Causal Revolution". In: (2018).

[5] J. Wei et al. "Emergent abilities of large language models". In: *arXiv preprint arXiv:2206.07682* (2022).

[6] E. Kıcıman et al. "Causal Reasoning and Large Language Models: Opening a New Frontier for Causality". In: *arXiv e-prints*, arXiv:2305.00050 (Apr. 2023).

[7] C. Glymour, K. Zhang, and P. Spirtes. "Review of Causal Discovery Methods Based on Graphical Models". In: *Frontiers in Genetics* 10 (2019). ISSN: 1664-8021.

[8] P. Spirtes and C. Glymour. "An Algorithm for Fast Recovery of Sparse Causal Graphs". In: *Social Science Computer Review* 9.1 (1991), pp. 62–72.

[9] T. Verma and J. Pearl. "Equivalence and Synthesis of Causal Models". In: *Probabilistic and Causal Inference* (1990).

[10] M. Willig et al. "Causal Parrots: Large Language Models May Talk Causality But Are Not Causal". In: *Submitted to Transactions on Machine Learning Research* (2023).

[11] C. Zhang et al. "Understanding Causality with Large Language Models: Feasibility and Opportunities". In: *arXiv e-prints*, arXiv:2304.05524 (Apr. 2023).

[12] A. Zanga and F. Stella. "A Survey on Causal Discovery: Theory and Practice". In: *arXiv e-prints*, arXiv:2305.10032 (May 2023).

[13] *A COMPLETE GUIDE TO Natural Language Processing*. URL: https://www.deeplearning.ai/resources/natural-language-processing/. (accessed: 23.08.2023).

[14] *OpenAI*. URL: https://openai.com/. (accessed: 29.08.2023).

[15]  *What is a REST API?* URL: https://www.redhat.com/en/topics/api/what-is-a-rest-api. (accessed: 29.08.2023).

[16]  *API Reference - OpenAI API*. URL: https://platform.openai.com/docs/api-reference/introduction?lang=python. (accessed: 23.08.2023).

[17]  *GPT - OpenAI API*. URL: https://platform.openai.com/docs/guides/gpt. (accessed: 29.06.2023).

[18]  *ChatGPT API Transition Guide*. URL: https://help.openai.com/en/articles/7042661-chatgpt-api-transition-guide. (accessed: 29.08.2023).

[19]  I. Fulford and A. Ng. *ChatGPT Prompt Engineering for Developers*. URL: https://www.deeplearning.ai/short-courses/chatgpt-prompt-engineering-for-developers/. (accessed: 10.08.2023).

[20]  *Prompt Engineering Guide*. URL: https://www.promptingguide.ai/. (accessed: 30.06.2023).

[21]  M. Hobbhahn and T. Lieberum. *Investigating causal understanding in LLMS*. URL: https://www.lesswrong.com/posts/yZb5eFvDoaqB337X5/investigating-causal-understanding-in-llms. (accessed: 02.06.2023).

[22]  The PyCoach. *OpenAI and Andrew Ng Just Released a FREE ChatGPT Prompt Engineering Course*. URL: https://artificialcorner.com/openai-and-andrew-ng-just-released-a-free-chatgpt-prompt-engineering-course-b0884c03e946. (accessed: 10.08.2023).

[23]  J. Shieh. *Best practices for prompt engineering with OpenAI API*. URL: https://help.openai.com/en/articles/6654000-best-practices-for-prompt-engineering-with-openai-api. (accessed: 10.08.2023).

[24]  *PubMed*. URL: https://pubmed.ncbi.nlm.nih.gov/. (accessed: 01.06.2023).

[25]  E. Sayers. *A General Introduction to the E-utilities*. URL: https://www.ncbi.nlm.nih.gov/books/NBK25497/. (accessed: 02.06.2023).

[26]  D. Johnson. "Finding All the Elementary Circuits of a Directed Graph". In: *SIAM Journal on Computing* 4.1 (1975), pp. 77–84.

[27]  A. A. Hagberg, D. A. Schult, and P. J. Swart. "Exploring Network Structure, Dynamics, and Function using NetworkX". In: *Proceedings of the 7th Python in Science Conference*. 2008, pp. 11–15.

[28]  T. Khuyen. *Pyvis: Visualize Interactive Network Graphs in Python*. URL: https://towardsdatascience.com/pyvis-visualize-interactive-network-graphs-in-python-77e059791f01. (accessed: 01.07.2023).

[29]  S. kaare Larsen. "Creating Large Language Model Resistant Exams: Guidelines and Strategies". In: *arXiv e-prints*, arXiv:2304.12203 (Apr. 2023).

STUDENTSUPSI

[30] K. Singhal et al. "Publisher Correction: Large language models encode clinical knowledge". In: 620.7973 (Aug. 2023), E19–E19. arXiv: 2212.13138.

[31] J. Liu et al. "Benchmarking Large Language Models on CMExam – A Comprehensive Chinese Medical Exam Dataset". In: *arXiv e-prints*, arXiv:2306.03030 (June 2023).

[32] R. Raimondi et al. "Comparative analysis of large language models in the Royal College of Ophthalmologists fellowship exams". In: *Eye (London, England)* (May 2023).

[33] D. Arora, H. Gaurav Singh, and Mausam. "Have LLMs Advanced Enough? A Challenging Problem Solving Benchmark For Large Language Models". In: *arXiv e-prints*, arXiv:2305.15074 (May 2023).

[34] J. Holmes et al. "Evaluating Large Language Models on a Highly-specialized Topic, Radiation Oncology Physics". In: *arXiv e-prints*, arXiv:2304.01938 (Apr. 2023).

[35] J. M. Mooij et al. "Distinguishing cause from effect using observational data: methods and benchmarks". In: *arXiv e-prints*, arXiv:1412.3773 (Dec. 2014).

[36] R. Tu, C. Ma, and C. Zhang. "Causal-Discovery Performance of ChatGPT in the context of Neuropathic Pain Diagnosis". In: *arXiv e-prints*, arXiv:2301.13819 (Jan. 2023).

[37] M. Scutari. *Bayesian Network Repository*. URL: https://www.bnlearn.com/bnrepository/. (accessed: 20.07.2023).

[38] A. Nogueira et al. "Methods and tools for causal discovery and causal inference". In: *WIREs Data Mining and Knowledge Discovery* 12.2 (2022).

[39] H. Pan. "Research progress on the protective mechanism of a novel soluble epoxide hydrolase inhibitor TPPU on ischemic stroke". In: *Frontiers in Neurology* 14 (2023). ISSN: 1664-2295.

[40] J. Tfelt-Hansen et al. "Risk stratification of sudden cardiac death: a review". In: *EP Europace* 25.8 (Aug. 2023), euad203. ISSN: 1099-5129.

[41] Y. Ye, J. Xu, and Y. Han. "Dexamethasone for delayed edema after intracerebral hemorrhageTo be or not to be?" In: *Heliyon* 9.7 (July 2023), e17621.

[42] T. Hosoya et al. "Papillary Glioneuronal Tumor Masquerading as Malignant Brain Tumors: A Case Report". In: *Yonago Acta Med* 66.3 (Aug. 2023), pp. 385–388.

[43] K. T. Alghamdi et al. "Gliosarcoma Invading the Temporal Bone, Temporalis Muscle, and Skull Base". In: *Cureus* 15.7 (July 2023), e42319.

[44] J. Jung, I. Shin, and Y. Choi. "A Rare Case of Nontuberculous Mycobacterial Abscess Mimicking Brain Tumor in an Immunocompetent Patient". In: *Brain Tumor Res Treat* ().

[45] L. A. McLendon et al. "Dramatic Response to Anti-IL-6 Receptor Therapy in Children With Life-Threatening Myelin Oligodendrocyte Glycoprotein-Associated Disease". In: *Neurol Neuroimmunol Neuroinflamm* 10.6 (Nov. 2023).

STUDENTSUPSI

[46]   K. Tanabe and J. Tanabe. "Role of echocardiography in patients treated with immune checkpoints inhibitors". In: *J Echocardiogr* (Aug. 2023).

[47]   T. Paul et al. "Advancing the science of management of arrhythmic disease in children and adult congenital heart disease patients within the last 25 years". In: *Europace* 25.8 (Aug. 2023).

[48]   O. Abdelkarim, Y. Saleh, and F. Nabi. "Left Atrial Thrombus in the Setting of Mitral Stenosis". In: *Methodist Debakey Cardiovasc J* 19.1 (2023), pp. 61–63.

STUDENTSUPSI