

# MScCBBi

MASTER IN  
**COMPUTATIONAL BIOLOGY  
& BIOINFORMATICS**

**SPECIALIZATION** MULTI-OMICS FOR LIFE AND HEALTH SCIENCES

**MARIA INÊS NUNES VILAR GOMES**

BSc in Aquatic Sciences

## **BENCHMARKING COMPUTATIONAL ALGORITHMS FOR ENHANCED DRUG DISCOVERY**

September, 2025



DEPARTMENT OF  
LIFE SCIENCES

# BENCHMARKING COMPUTATIONAL ALGORITHMS FOR ENHANCED DRUG DISCOVERY

INSIGHTS FROM CLARIVATE'S PRE-COMPETITIVE ALGORITHM BENCHMARKING  
CONSORTIUM

**MARIA INÊS NUNES VILAR GOMES**

BSc in Aquatic Sciences

**Advisers:** Dr. Filippo Ciceri

*Lead Data Scientist, Clarivate*

Dra. Cecilia Klein

*Director, Clarivate*

**Co-adviser:** Prof. Dra. Paula Maria Theriaga Mendes Bernardo Gonçalves

*Associate Professor, NOVA University Lisbon*

MASTER IN COMPUTATIONAL BIOLOGY AND BIOINFORMATICS  
SPECIALIZATION MULTI-OMICS FOR LIFE AND HEALTH SCIENCES

NOVA University Lisbon  
September, 2025

## **Benchmarking Computational Algorithms for Enhanced Drug Discovery Insights from Clarivate's pre-competitive Algorithm Benchmarking Consortium**

Copyright © Maria Inês Nunes Vilar Gomes, NOVA School of Science and Technology, NOVA University Lisbon.

The NOVA School of Science and Technology and the NOVA University Lisbon have the right, perpetual and without geographical boundaries, to file and publish this dissertation through printed copies reproduced on paper or on digital form, or by any other means known or that may be invented, and to disseminate through scientific repositories and admit its copying and distribution for non-commercial, educational or research purposes, as long as credit is given to the author and editor.

### **Declaration on the use of artificial intelligence tools**

This work used generative artificial intelligence tools, namely QuillBot [1] and Microsoft Copilot [2], to assist with text paraphrasing. These tools were used under the author's supervision, and all generated content was checked for accuracy. Authorship and content validation remained the author's responsibility, and the use of artificial intelligence complies with institutional standards of academic integrity.

## ACKNOWLEDGEMENTS

This thesis represents much more than a year dedicated to this project and its writing. It marks an important period of growth. A period in which I developed multidisciplinary skills by immersion in a fast-paced consulting environment of a multinational company. I would therefore like to express my sincere gratitude to all those who have supported and guided me throughout my journey.

First, I am deeply grateful to Cecilia, Filippo, Sarah for believing in me and for offering me the unique opportunity to join such a talented team. I would like to thank Alex who contributed substantially to the development and execution of this project, for giving me the privilege of being part of such an impactful and significant project, and for helping me at all times. A special thank you goes to my supervisor, Filippo, for sharing his knowledge and for making the entire journey enjoyable and rewarding, even though the most challenging moments. His guidance, encouragement, and constant support were essential. My sincere thanks also go to all the members of the Discovery and Translational Science Team. Each one, without exception, contributed significantly to my learning process, enriching my experience.

To Professor Paula, for all the support and willingness to help at any time, not only during the last year, but since the first day of the master's course.

Finally, I am infinitely grateful to my family and friends for their encouragement and unconditional support throughout all the challenges, as well as for the joy shared in all my achievements and successes.

## ABSTRACT

Drug development is a complex process that requires integrating information from many different fields of knowledge. During the preclinical phase, omics measurements are often used to clarify the phenotypical changes caused by exposure to a therapeutic agent, thereby enabling the reconstruction of the biochemical cascade responsible for the pharmacologic effect (often called Mechanism of Action (MoA) reconstruction). This step is crucial in characterizing a therapeutic agent, and researchers rely on numerous computational tools for this process. However, the abundance of available solutions and the lack of standardized assessment frameworks pose challenges in identifying reliable and efficient tools. To address this, Clarivate leads the pre-competitive subscription-based, Algorithm Benchmarking Consortium (ABC), which evaluates computational tools for various use cases. Within the ABC framework, 27 algorithms (13 topology-based, 6 connectivity-based, and 8 enrichment-based algorithms) for identifying key regulators and inferring dysregulated signaling proteins from transcriptomics data were benchmarked to provide an unbiased assessment of available methods. The algorithms were evaluated together with transcriptomic data derived from 7 databases and molecular networks from OmniPath and MetaBase<sup>TM</sup>. Performance metrics included accuracy in predicting causal regulators, computational scalability, and robustness to input data variations.

Causal reasoning, randomWalk, and non-causal, overconnectivity, algorithms showed superior results overall for direct target recovery. Among enrichment-based tools, wmean was the best performer. At the reference dataset level, MetaBase network lead to the best outcomes. At the query dataset level, tissue-derived signatures outperformed cell-line-derived ones. Through this use case, ABC emphasizes the importance of MoA reconstruction as a means to better understand the pharmacological effects of therapeutic agents. Ultimately, this is essential for characterizing safety profiles, identifying new indications (drug repurposing) and exploring potential combinatorial strategies with synergistic effects.

**Keywords:** Mechanism of Action, Transcriptomics, Algorithm benchmarking, Causal reasoning

## RESUMO

O desenvolvimento de medicamentos é um processo complexo que requer a integração de informações de muitas áreas diferentes do conhecimento. Durante a fase pré-clínica, dados ômicos são frequentemente utilizados para esclarecer as alterações fenotípicas causadas pela exposição a um agente terapêutico, permitindo assim a reconstrução da cascata bioquímica responsável pelo efeito farmacológico (frequentemente chamada de reconstrução do mecanismo de ação). Esta etapa é crucial para caracterizar um agente terapêutico, e os investigadores dependem de inúmeras ferramentas computacionais para este processo. No entanto, a abundância de soluções disponíveis e a falta de estruturas de avaliação padronizadas representam desafios na identificação de ferramentas confiáveis e eficientes. Com isso em mente, a Clarivate lidera o Algorithm Benchmarking Consortium (ABC), um consórcio baseado em subscrição para farmacêuticas que avalia ferramentas computacionais para vários casos de estudo. Dentro do ABC, 27 algoritmos (13 baseados em topologia, 6 baseados em conectividade e 8 baseados em enriquecimento) para identificar reguladores-chave e inferir proteínas de sinalização desreguladas a partir de dados transcriptômicos foram comparados para fornecer uma avaliação imparcial dos métodos disponíveis. Os algoritmos foram avaliados juntamente com dados transcriptômicos derivados de 7 bases de dados e redes moleculares da OmniPath e MetaBase<sup>TM</sup>. As métricas de desempenho incluíram precisão na previsão de reguladores causais, escalabilidade computacional e robustez em relação às variações dos dados utilizados pelos algoritmos.

Os algoritmos de raciocínio causal, randomWalk, e não causal, overconnectivity, apresentaram resultados superiores em geral para a recuperação direta de alvos. Entre as ferramentas baseadas em enriquecimento, o wmean foi o que apresentou melhor desempenho. A nível de dados de referência, a rede molecular MetaBase apresentou melhores resultados. E a nível de dados de consulta, as assinaturas derivadas de tecidos superaram as derivadas de linhas celulares. Através deste caso de uso, o ABC enfatiza a importância da reconstrução do mecanismo de ação de uma perturbação para obter uma melhor compreensão dos efeitos farmacológicos dos agentes terapêuticos. Em última análise, isso é essencial para caracterizar perfis de segurança, identificar novas indicações (reposicionamento de medicamentos) e explorar estratégias combinatórias potenciais com efeitos

sinérgicos.

**Palavras-chave:** Mecanismo de ação, Transcriptômica, Avaliação de algoritmos, Raciocínio causal

# CONTENTS

<b>List of Figures</b>	<b>viii</b>
<b>List of Tables</b>	<b>ix</b>
<b>Glossary</b>	<b>x</b>
<b>Acronyms</b>	<b>xii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Significance and Objectives . . . . .	1
1.2 Scope . . . . .	3
1.3 Other Contributions . . . . .	3
1.4 Structure . . . . .	4
<b>2 Literature Review</b>	<b>6</b>
2.1 Drug discovery: the importance of the compound’s mechanism of action	6
2.2 Transcriptomic data . . . . .	9
2.3 Prior knowledge Network . . . . .	19
2.4 Computational methods for MoA inference . . . . .	26
2.4.1 Connectivity-based tools for comparative analysis . . . . .	26
2.4.2 Enrichment-based tools for downstream analysis . . . . .	27
2.4.3 Topology-based methods for upstream analysis . . . . .	28
2.5 Benchmarking of computational methods for MoA inference . . . . .	30
<b>3 Materials and Methods</b>	<b>32</b>
3.1 Benchmarking architecture setup . . . . .	32
3.2 Data Description . . . . .	34
3.2.1 Gene expression data: Perturbation signatures . . . . .	34
3.2.2 Prior Knowledge: Interaction Networks . . . . .	37
3.3 Algorithms: implementation and wrapper function’s architecture . . . . .	39
3.3.1 Connectivity Mapping . . . . .	40



3.3.2	Pathway Enrichment . . . . .	41
3.3.3	Topology-based methods . . . . .	41
3.4	Evaluation . . . . .	50
<b>4</b>	<b>Results</b>	<b>52</b>
4.1	Algorithms . . . . .	52
4.2	Reference datasets . . . . .	56
4.3	Query datasets . . . . .	58
<b>5</b>	<b>Discussion</b>	<b>63</b>
	<b>References</b>	<b>69</b>
	<b>Annexes</b>	
<b>I</b>	<b>Supplementary Tables</b>	<b>80</b>

## LIST OF FIGURES

2.1	Visual representation of a compound-induced cellular signaling cascade. . .	8
2.2	Gene regulatory network based on regulon representation of human transcriptional interactions. . . . .	20
2.3	OmniPath molecular interactions database. . . . .	22
3.1	Schematic of the study architecture. . . . .	33
3.2	Flowchart representing the main steps for implementing connectivity mapping algorithms pre-built in the RCSM package. . . . .	42
3.3	Flowchart representing the main steps for implementing enrichment algorithms pre-built in the decoupleR package. . . . .	43
3.4	Flowchart representing the main steps for implementing the CARNIVAL algorithm. . . . .	44
3.5	Flowchart representing the main steps of ProTINA algorithm implementation. . . . .	45
3.6	Flowchart representing the main steps of CausalR algorithm implementation. . . . .	46
3.7	Flowchart representing the main steps of CIE algorithm implementation. . . . .	47
3.8	Flowchart representing the main steps of NicheNet algorithm implementation. . . . .	48
3.9	Flowchart representing the main steps of the wrapper function of both the CBDD baseline and CBDD causal reasoning algorithms implementation. . . . .	49
4.1	Runtime per signatures. . . . .	54
4.2	Mean scaled Area Under the Curve (AUC) values for each algorithm. . . . .	54
4.3	Mean win rate and WPAE across all evaluations for each algorithm. . . . .	56
4.4	Mean pathway enrichment values for each algorithm. . . . .	57
4.5	Final performance evaluation of 22 algorithms across five metrics. . . . .	58
4.6	Performance of nine reference datasets across four evaluation metrics. . . . .	59
4.7	Query dataset performance evaluation. . . . .	60
4.8	Comparison of algorithm performance between different signature types. . . . .	62

## LIST OF TABLES

2.1	List of resources that provide transcriptomic data driven from perturbagen.	13
2.2	List of resources that provide prior knowledge networks and pathways. . .	23
3.1	Summary of gene expression datasets used in this study. . . . .	36
3.2	Summary table of OmniPath and MetaBase prior knowledge resources. . . .	38
3.3	Summary of computational methods evaluated in the benchmarking study.	40
I.1	Summary of evaluated runs. . . . .	80
I.2	Algorithm execution success and run rate across all benchmarking runs. . .	81

## GLOSSARY

<b>Gene signature</b>	A gene signature is a specific collection of genes whose expression pattern is characteristic of a particular cellular state, disease outcome, or response to treatment. Usually, it includes the name/ID of the gene, together with a value representing their relative expression (fold-changes or p-values).
<b>Mechanism of Action (MoA)</b>	Molecular cascade by which a perturbation (such as a drug or genetic modification) produces its biological effect. The molecular cascade includes the interaction with the direct molecular target(s) and the immediate downstream cascade of events leading to a certain cellular outcome. The cascade can be reflected in changes in the gene expression.
<b>Molecular network</b>	A graph representation of molecular interactions where genes, proteins, and other molecular entities are represented by nodes, while their interactions are represented by edges. Networks can be directed (causality), weighted (interaction strength), or signed (effect, such as, activation/inhibition), or without any of these attributes. From a full network, it is possible to extract subsets of interactions (pathways). Pathways are cascades of molecular interactions, and some databases like KEGG or Reactome catalog those into specific types, such as metabolic, signaling, or regulatory pathways. Pathway enrichment analysis provides a high-level understanding of the biological processes by identifying coordinated network variations, instead of focusing on individual genes.

## **Transcriptomics**

The measurement of gene expression levels across the genome (or a subset of genes) obtained from a biological sample. One or more mRNA molecules are produced from the transcription of each gene. Transcriptomic technologies quantify individual mRNA molecules using various methods for measuring gene expression. Examples of some transcriptomic technologies are microarrays, L1000, and RNA-seq. Each one of these techniques can generate full gene expression profiles by capturing the expression of all the genes transcribed in a given sample. Transcriptomic analyses are usually performed by sampling at different time points, conditions, or treatments. The raw data can be further processed to identify DEGs, using metrics such as  $\log_2$  fold-change, z-scores, p-values, and q-values.

## ACRONYMS

<b>ABC</b>	Algorithm Benchmarking Consortium
<b>AUC</b>	Area Under the Curve
<b>AUPR</b>	Area Under the Precision-Recall Curve
<b>CARNIVAL</b>	CAusal Reasoning for Network identification using Integer VALue programming
<b>CBDD</b>	Computational Biology for Drug Discovery
<b>CCG</b>	Computational Causal Graph
<b>CDDI</b>	Cortellis Drug Discovery Intelligence
<b>CDS-DB</b>	Cancer Drug-induced Gene Expression Signature Database
<b>CIE</b>	Causal Inference Engine
<b>CMap</b>	Connectivity Mapping
<b>CRC</b>	Colorectal cancer
<b>CREEDS</b>	Crowd Extracted Expression of Differential Signatures
<b>CRISPR</b>	Clustered Regularly Interspaced Short Palindromic Repeats
<b>DD</b>	Drug discovery
<b>DEGs</b>	Differentially expressed genes
<b>DM2</b>	Diabetes Mellitus Type 2
<b>DNA</b>	Deoxyribonucleic acid
<b>DR</b>	Drug repositioning
<b>ES</b>	Enrichment score
<b>FDA</b>	Food and Drug Administration
<b>GEO</b>	Gene Expression Omnibus
<b>GRN</b>	Gene Regulatory Networks
<b>GSEA</b>	Gene Set Enrichment Analysis

<b>GWPS</b>	Genome-Wide Perturb-Seq
<b>HTS</b>	High-Throughput Screening
<b>KEGG</b>	Kyoto Encyclopedia of Genes and Genomes
<b>KS</b>	Kolmogorov-Smirnov
<b>LINCS</b>	Library of Integrated Network-Based Cellular Signatures
<b>MoA</b>	Mechanism of action
<b>mRNA</b>	Messenger RNA
<b>NCBI</b>	National Center for Biotechnology Information
<b>OE</b>	Overexpression
<b>ORA</b>	Over-Representation Analysis
<b>PGN</b>	Protein-Gene Regulatory Network
<b>PKN</b>	Prior Knowledge Network
<b>PPI</b>	Protein-Protein Interaction
<b>ProTINA</b>	Protein Target Inference by Network Analysis
<b>QS</b>	Quaternary scoring
<b>R</b>	R Programming Language
<b>RCSM</b>	Recommended Connectivity-map Scoring Methods
<b>RNA</b>	Ribonucleic acid
<b>RNA-seq</b>	RNA sequencing
<b>scRNA-seq</b>	single-cell RNA sequencing
<b>shRNA</b>	short hairpin RNA
<b>TF</b>	Transcription factor
<b>TS</b>	Ternary scoring
<b>WPAE</b>	win percentage above expected
<b>xCos</b>	explainable cosine
<b>XSum</b>	eXtreme Sum

# INTRODUCTION

*This section summarizes the study’s underlying motivation, rationale, and goals, emphasizing its significance in the field. It provides context by giving some background on the supporting company and the initiative, along with other contributions. Furthermore, it outlines a reading guide for this thesis.*

## 1.1 Significance and Objectives

Drug discovery (DD) and development is a time-consuming, resource-intensive, multidisciplinary effort that can be challenging from many points of view. Over half of clinical trial failures are attributed to lack of efficacy, underscoring the importance of identifying and validating pharmacological targets, and highlighting the lack of knowledge of the drug’s Mechanism of Action (MoA) as one of the major barriers to clinical efficacy [4–6]. As thorough understanding of the Mechanism of action (MoA) is such a critical step, computational methods can accelerate the identification of pharmacologically active agents by providing more efficient and cost-effective alternatives to traditional approaches. Computational methods can accurately identify a new target or propose new indications for a known molecule, thereby reducing the time dedicated to *in vitro* target validation [7].

A key to understanding a compound’s MoA lies in transcriptomic profiling, which captures the changes in gene expression triggered by a perturbagen. While traditional RNA-seq methods remain too costly for large-scale expression signatures, recent High-Throughput Screening (HTS) advances, such as the L1000 assay [8], enable affordable generation and analysis of large-scale omics datasets. Several existing databases provide public access to transcriptomic data from experiments exposing cell lines to range of chemical and genetic perturbagens. These datasets can be leveraged, using various computational tools, to establish the causal chain of gene expression changes triggered by a specific compound. Three primary approaches have emerged: causal reasoning, connectivity mapping, and enrichment tools [9].

Causal reasoning is a topology-based method, that determines potential causes for



an observed gene expression profile, starting from a perturbation signature and a biological interaction network. The Molecular network is defined as a signed and directed graph describing relations between nodes (e.g., proteins or genes). Efforts to compile causal molecular relation networks have increased, resulting in several publicly accessible databases (such as OmniPath), which offers curated Prior Knowledge Network (PKN). Among the networks commercially available [10] we can find MetaBase<sup>TM</sup>, developed and curated by Clarivate.

The Connectivity Mapping (CMap) method is instead focused on collecting and analyzing perturbation signatures. In this case, a similarity score is used to compare a set of known MoA/compound reference signatures, with a query gene expression profile from a perturbagen of interest [9, 11]. The principle behind CMap is that the higher the similarity between the query and the reference signature, the more likely it is that the underlying mechanism is the same.

On the other hand, enrichment tools take perturbation signatures as query input and utilize prior knowledge, such as a regulon network or collections of perturbation-induced DEGs, as a reference. The purpose of these tools is to determine whether specific gene sets or regulons (genes interacting with Transcription factor (TF)) are enriched in the perturbed data [12].

With a systematic benchmarking approach, this study guides the selection of the most appropriate algorithms and inputs for evaluating the plethora of current solutions for MoA reconstruction. The evaluation process relies on three inputs:

- Gene signatures derived from chemical or genetic perturbation experiments.
- A PKN of the molecular interactions within the system.
- A gold standard dataset for validation of the results.

These data are used to feed the evaluated methods, serving as the reference, query, and gold standard datasets, respectively. This study analyzes the three types of tools depicted below:

**Topology-based tools** Eight causal reasoning algorithms (CAusal Reasoning for Network identification using Integer VALue programming (CARNIVAL), CausalR, Protein Target Inference by Network Analysis (ProTINA), Causal Inference Engine (CIE), NicheNet, plus causalReasoning, SigNet, quaternaryProd from the Computational Biology for Drug Discovery (CBDD) R Programming Language (R) package [13] , and five node prioritization algorithms (networkPropagation, randomWalk, over-connectivity, hiddenNodes, interconnectivity - from CBDD).

**Similarity-based tools** Six algorithms from the Recommended Connectivity-map Scoring Methods (RCSM) R package.

**Enrichment-based tools** Eight algorithms from the decoupleR package [12].

Performance is assessed by comparing the results obtained against gold standard datasets, along with the robustness, and the computational efficiency. With this approach, the project aims to identify the most appropriate tools to contextualize gene expression data and if the choice of the inputs can also have an impact on the outcomes.

## 1.2 Scope

This project was conducted within the framework of the Algorithm Benchmarking Consortium (ABC), a subscription-based initiative for pharmaceutical companies led by Clarivate. ABC is dedicated to evaluating a wide range of computational tools for a variety of applications in the life sciences and healthcare field. The topic for this thesis is the development of ABC's tenth use case - Causal Regulation - which focuses on benchmarking tools designed to identify key regulators from transcriptomic data and PKN.

## 1.3 Other Contributions

This study expands the state of the art in causal reasoning using gene expression data and molecular interactions, by presenting a robust framework and a systematic algorithm benchmarking approach. The study was presented during the following poster communication:

**XIV Edition of Bioinformatics Open Days** M. Gomes, A. Ishkin, F. Ciceri, C. Klein. Benchmarking Computational Algorithms for enhanced Drug Discovery: Insights from Clarivate's pre-competitive Algorithm Benchmarking Consortium. XIV Edition of Bioinformatics Open Days, 26 March 2025, Braga, Portugal.

Beyond the topic of this thesis, during the MSc industry placement, I was also involved as developer in other activities aimed at identifying reliable solutions for several different external stakeholders in the pharmaceutical business. Although not related with the topic of this thesis, these experiences enriched my consultant experience, allowing for an expansion in the expertise across various domains of computational biology and data science. These projects included:

**Skin Microbiome Atlas** This project involved the curation and re-analysis of publicly available skin microbiome datasets. My role in this project included conducting an in-depth scientific literature review, compiling relevant datasets, and pre-processing raw sequencing data to ensure consistency and comparability across studies.

**Natural Language Processing** A pipeline was set up for automatic text classification of epidemiology abstracts, leveraging different versions of the BERT foundational model. Model fine-tuning on a minimal dataset was attempted by generating synthetic data using paraphrasing techniques.

**ABC - Spatial Niche** As part of an internal case study for ABC, I implemented several Python wrappers for selected algorithms relevant to spatial niche analysis. My role in this use case mainly included the management of conda environments, the development of functions to run those algorithms and their integration with an R-based pipeline.

**Google Data Extraction Tool** In collaboration with another team in the company, I developed an automated script for retrieving pharmacological information from different web sources. The pipeline involved querying URLs and keywords and extracting structured data through Google search API calls.

**Transcriptomic comparative analysis** I was involved in a project to perform a comparative analysis of transcriptomic profiles from three types of cancer. The workflow included exploratory data analysis, identification of DEGs and hub genes, pathway enrichment analysis, causal reasoning to infer upstream regulators, and a survival analysis. The analysis also included the comparison of both the Human Papillomavirus status and tumor location.

**Proteomics analysis** In a proteomics-focused project, I was responsible for carrying out exploratory data analysis and functional enrichment analysis to uncover biological insights from protein-level data.

**Single-cell atlas** For the single-cell atlas project, I performed some downstream analysis, including enrichment analysis on a already built atlas. At this moment, I am involved in the construction of a new single-cell atlas, with the integration, pre-processing and quality control of datasets.

**Indication-Prioritization** This project involved the development of a pipeline to integrate disease and target data annotation from several sources, and prioritize disease indications based on customer's preference criteria.

## 1.4 Structure

This study is organized in five chapters. After this introductory chapter, Chapter 2 starts by describing what MoA is and its importance in DD, followed by a review of the data and algorithms that can be used for MoA reconstruction. This chapter also characterizes the role of a benchmarking study when there is a wide diversity of data and methods available. Chapter 3 describes the methodology used in this study, including the selection of gene expression datasets, Molecular networks, and algorithms for benchmarking, the design of the experimental framework, and the evaluation metrics used to assess algorithm performance. Chapter 4 presents the results of the benchmarking, including the analysis of algorithm performance across different levels of the data. In Chapter 5 the results are interpreted, putting them in the context of the current literature in the field, and

summarizes the key findings of the study and their significance, offering recommendations for future work. At the end, Annex I includes the supplementary tables.

## LITERATURE REVIEW

*This section presents a summary of the relevant research related to the topic of the project. First, it begins by highlighting the importance of understanding the MoA in DD, followed by a description of two fundamental components for MoA reconstruction: transcriptomic data and biological networks. Then, the computational methods used to apply various scoring algorithms (topology-, similarity-, and enrichment-based algorithms) are summarized, laying the groundwork for a systematic evaluation of tools designed to infer compound MoA. Finally, this chapter outlines best practices for conducting a benchmarking study.*

### **2.1 Drug discovery: the importance of the compound's mechanism of action**

Developing new drugs is an extraordinarily complex process. The high prevalence of complex diseases, which collectively account for 70% of all the deaths in Europe and affect around 25% of the population, is one of the challenges faced by this industry [14]. In addition, statistics show that *de novo* DD has become an extensive and costly process, taking on average 13 years and 2\$ billion to develop a new drug, with most of clinical trials lasting 95 months and non-clinical development 31 months [15–17]. These challenges have led to fewer drug approvals by regulatory bodies, resulting in a significant gap between therapeutic demand and available treatments. Hence, as the current treatments become less effective, there is a strong interest in finding alternatives to optimize critical steps in the drug development pipeline and developing more advanced therapeutic methods [18]. Efforts to address these challenges are reflected in the growing number of studies, both in industry and academia.

Drug repositioning (DR) became an attractive strategy to tackle the constraints faced by traditional DD by reducing the initial cost to 1/3 and the duration to 3-9 years, and it continues to gain increasing attention, as nearly 30% of the drugs approved by the Food and Drug Administration (FDA) are identified using this approach [19, 20]. The fundamental goal of DR is to broaden the indication of known, safe, and previously approved drugs. From multiple points of view, this is a particularly interesting approach.

It allows to investigate therapeutic agents that have been put on hold because of failed clinical trials [21], and also it enables to identify treatment for conditions with unmet clinical needs. This is the case of rare diseases, which are not providing sufficient returns to pharmaceutical companies to justify a conventional DD pipeline. Many studies have demonstrated the success of establishing new drug-disease relationships [22]. Sildenafil is a well-known example. It was first discovered in the 1980s as a possible treatment for angina pectoris. The FDA approved it in 1998 to treat erectile dysfunction and again in 2005 to treat pulmonary arterial hypertension [19, 20, 23]. Another classic example is Thalidomide, originally used for sedation and morning sickness, and afterward repurposed for multiple myeloma, leprosy [19, 20], and to minimize the hippocampal neuronal loss [24]. Moreover, the low success rate (5%) for phase I clinical studies of cancer treatments led to increased attention in DR for oncology, resulting in several promising findings [20, 25]. Noteworthy cases include the schizophrenia drug Spiperone, which has been studied for inducing apoptosis in Colorectal cancer (CRC) cells [26], and Raloxifene, indicated for osteoporosis, which proved to be effective in reducing breast cancer risk in postmenopausal women [20, 27].

Understanding how cellular signaling (Figure 2.1) is modulated upon a stimulus is essential for identifying potential drug targets and finding new indications for an existing drug. When a drug enters a biological system, it typically interacts directly or indirectly with cellular targets, regulating the activity of signaling networks and pathways. This is commonly referred as the MoA [28, 29]. These interactions are relevant across the whole DD process, from initial investigation to clinical trials. The reconstruction of a drug's MoA allows to uncover drug-exposure biomarkers, anticipate early adverse effects, and even synergistic effects resulting from drug combinations. Nevertheless, FDA approval can be obtained without knowing the drug's MoA if the drug exhibits sufficient safety and efficacy [9, 30]. Yet, not knowing the mechanisms of the compounds can be extremely disadvantageous. This is demonstrated by the case of Dimebon. Originally developed as an antihistamine drug, it later entered clinical trials (with the MoA still unknown) for the treatment of Alzheimer's disease, failing to show meaningful clinical efficacy in phase 3 studies. Later, it was discovered that it was the activity on the histamine and serotonin receptors that caused the initial observed cognitive efficacy, instead of the stabilization of mitochondria (as first hypothesized) [9, 31].

Although we refer to the target(s) of a compound as direct, this is often not the case. From a chronological point of view, there are a series of interactions that result in modulation of biological processes, and what is detected at a given moment does not always linearly reflect what happened previously. Indeed, the basic definition of MoA is just the tip of the iceberg, given the chain of biochemical reactions forming part of the cell signaling cascade. This process is characterized by the signaling pathways leading to a certain cellular response. These pathways can also interact with each other through crosstalk [23], forming a complex network of interconnected and distinct nodes. The impact of a certain compound in the complex cell signaling cascade can be defined and observed on a system

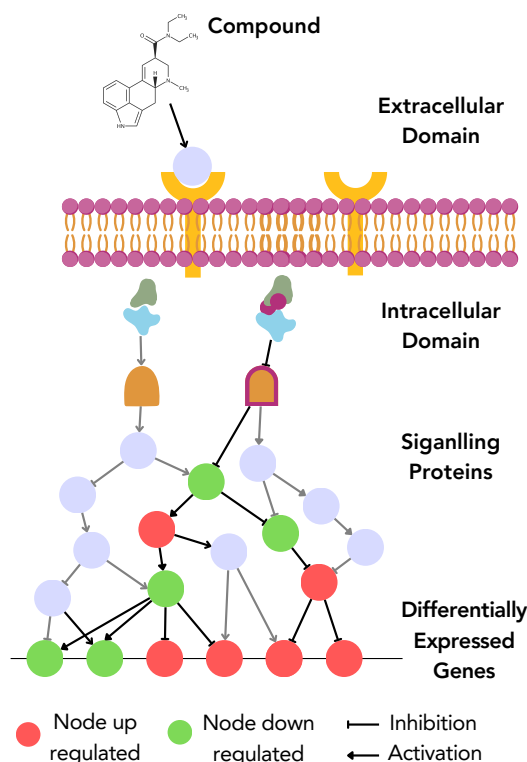


Figure 2.1: Visual representation of a compound-induced cellular signaling cascade, where binding of the perturbing compound to its extracellular receptor domain triggers downstream intracellular signal transduction events through various signaling proteins, ultimately altering gene expression levels in the nucleus. The red nodes represent upregulated genes, the green nodes represent downregulated genes, arrows denote activation events, and T-bar edges indicate inhibitory interactions. Adapted from Trapotsi *et al.* (2022) [9].

level by high throughput approaches such as Genomics, Transcriptomics, Proteomics, Metabolomics, and even Phenomics. Each of them provides a different perspective on the compound's activity. Despite these technological progresses, experimental identification of targets, signaling proteins, and biological pathways (MoA) modulated by an uncharacterized compound can be extremely difficult. To facilitate this process HTS *in silico* methods have become an attractive option, given the affordability, speed and the amount of data provided. These methods act as a screening process, helping to identify and produce mechanistic hypotheses for further experimental validation [9]. Many accessible computational resources integrate omics data with prior knowledge graphs, such as Gene Regulatory Networks (GRN), to enhance and pinpoint the drug's potential cellular targets. However, choosing the most suitable data and computational tool for each situation is not always simple, requiring precise identification of the scientific question that needs to be addressed. By doing this, researchers can choose the most appropriate data type and bioinformatics tools to efficiently study a compound's MoA.

## 2.2 Transcriptomic data

Transcriptomic data provides a comprehensive view of gene expression changes in response to a compound. Following the compound's perturbation, these data will reflect the differential mRNA expression, by capturing modulated signaling and TF activity changes triggered by the perturbagen. For this reason, after a cellular disturbance, changes in gene expression levels give rise to a perturbed gene expression signature [32]. In 2000, the first reference resource was built from a compilation of *Saccharomyces cerevisiae* gene expression signatures derived from pharmacological and genetic perturbations. However, the same authors recognized that it would require less expensive technologies to create a larger database of gene expression responses, because the existing methods were too costly to apply at scale [33–35]. Since then, driven by growing interest in DD optimization, several databases (Table 2.1) have emerged to aggregate and publicly provide perturbagen transcriptomic signatures.

DrugMatrix was the first larger molecular toxicology database, created in 2006 by Iconix Pharmaceuticals (later acquired by the National Institute of Environmental Health Sciences) and publicly available since 2011 [33, 36, 37]. The database comprises the gene expression responses, detected via microarray, to more than 600 perturbagens in rat tissues. Additionally, it provides information on chemical treatments related to histopathology, hematology, and clinical chemistry, enabling the investigation of specific types of toxicity [38]. However, studies *in vivo* limit the number of disturbances that can be studied, due to excessive costs that make it impractical to generate data on a large scale, as well as all the associated ethical implications [11]. In addition, transcriptional changes are usually cell-specific, so a precise analysis of transcriptomic changes should be carried out considering the cell type, further increasing the complexity of this effort [39].

CMap is a resource that systematically establishes connections between the MoA of chemical compounds, diseases, and biological processes through pattern matching between signatures derived from cell lines exposed to different perturbagens. The first version of CMap (CMap 1.0) generated 453 signatures derived from 164 distinct small molecules applied at two-time points (over a range of concentrations), to four human cell lines (MCF7, PC3, HL60, and SKMEL5) [11]. Although the goal could be achieved using signatures from various omic layers, this resource focused only on mRNA expression data through Deoxyribonucleic acid (DNA) microarrays. While the concept behind this database has become extensively utilized, the data generated by this pilot project was too small for the potential applications in the DD area. The experimental conditions tested were few and undiversified in terms of perturbagens and biological systems. The recent advances in high-throughput technology allowed large data acquisition, as in the case with the L1000 assay. The Library of Integrated Network-Based Cellular Signatures (LINCS) program extended the CMap to a second version (CMap 2.0 or LINCS L1000) that measured the expression of 978 landmark genes (representatives of various biological processes), in cells exposed to a different set of stimuli [8]. Using the Luminex L1000



platform, expression of these landmark genes is directly measured, and computational methods and *in silico* imputation then infer the expression of 11350 additional genes, to provide wider reconstruction of the genomic profiles [8]. In a preliminary phase, LINCS released over 6000 signatures from around 1300 small molecules, many of which were FDA-approved [33]. As of today, LINCS includes more than one million gene expression profiles from over 20000 chemical and genetic perturbagens, tested at multiple time points and doses across various human cell lines. Currently, the dataset is more than a thousand times larger than the CMap pilot dataset. The perturbagens include pharmacologic (small molecules) and genetic modulators, such as gene knockdowns using short hairpin RNA (shRNA) and/or CRISPR and induced Overexpression (OE). LINCS L1000 provides data on five levels. Levels 1 to 4 contain data at different pre-processing stages, and Level 5 contains the final signatures, where replicates (usually three per treatment) are combined into a single differential expression vector. This level is recommended by the authors to use for downstream analyses [8]. The data provided by LINCS L1000 has increased the understanding of the association between changes in gene expression and certain disorders, facilitating DR and contributing to the generation of testable hypotheses about the MoA of less characterized compounds. However, the presence of gene expression changes imputed and not measured directly can lead to some inaccuracies. In addition, the inherent complexity of cellular responses must be considered, since gene expression snapshots may not fully capture the dynamic nature of biological processes and do not always correlate with protein expression due to post-translational modifications [9]. Despite these limitations, several computational methods have been developed to apply these data and facilitate the DD process. These methods will be described later in section 2.4.

The Cancer Drug-induced Gene Expression Signature Database (CDS-DB) [40] is an user-friendly resource, released in September 2023, aiming to provide gene expression data from patient samples, following exposure to anti-cancer therapies. It compiles gene expression profiles from 78 patient-derived paired pre- and post-treatment datasets, (from Gene Expression Omnibus (GEO) and ArrayExpress databases), with manually curated clinical information. These sources have been organized into 219 CDS-DB datasets, composed of paired pre- and post-treatment gene expression profiles from human samples. Pairing has been performed considering a wide range of factors (therapeutic regimen, administration dosage, cancer subtype, sampling location, time, and drug response status). From those, datasets containing at least two patients have been used to generate differential expression analyses, resulting in 181 gene perturbation signatures. In addition, 2012 patient-level perturbation signatures have been derived by comparing pre- and post-treatment profiles from individual patients (e.g., baseline vs. 14-day; baseline vs. three months). All transcriptomic data in CDS-DB were uniformly re-processed from raw files (microarray or RNA-seq), the metadata was manually curated, and the terminologies for drugs, cancers, and genes were harmonized. This database is a valuable resource for MoA elucidation studies as it provides well-curated gene expression data for various

cancer types and treatments, including pre- and post-treatment samples [40]. Nonetheless, LINCS mainly focuses on cancer cells, and CDS-DB has only signatures derived from cancer patients, so they are not ideal for answering questions related to transcriptional changes in non-transformed cells [39].

ChemPert [39], emerged as a manually curated resource that maps the relationships between chemical perturbations, their protein targets, and downstream transcriptional signatures in non-transformed cells. It provides a user-friendly interface that includes two sections: the database, and a web analysis tool. The database has three main components: (1) Direct signaling protein targets of chemical perturbagens, curated from Drug Repurposing Hub, DrugBank, and STITCH v5.0; (2) Initial gene expression profiles of untreated non-transformed cells, extracted from GEO, ArrayExpress, and LINCS L1000 (Chemical perturbation datasets of non-cancer cells from LINCS L1000 at Level 3); (3) Transcriptional responses after perturbation, categorized as upregulated or downregulated. ChemPert database encompasses over 82000 transcriptional signatures from the exposure to 2566 chemical compounds across 167 different non-transformed cells and tissues. It includes target information for 57818 chemical compounds, capturing activation, inhibition, or unknown effects. Additionally, ChemPert also offers two built-in analysis tools: (1) Given a perturbagen and the initial gene expression profile, the users can predict how TF will respond; (2) Given a specific transcriptional response the users can identify the potential perturbagens based on that input data.

Bulk transcriptomic databases average gene expression across cells, potentially masking important heterogeneity in the biological responses, such as the existence of rare cell subpopulations surviving chemotherapy [41]. To capture these variations, single-cell perturbation sequencing methods have emerged. Techniques like Sci-Plex (for chemical perturbations) and Perturb-seq (for genetic perturbations) leverage mass screening technologies in combination with single-cell resolution, to provide a more detailed view of cellular responses [42]. Although traditional single-cell RNA sequencing (scRNA-seq) is essential for analyzing tissue heterogeneity, its high cost remains a barrier.

Sci-Plex [41] was introduced to overcome this limitation, by combining two techniques: nuclear hashing and combinatorial indexing-based RNA sequencing (sciRNA-seq), to assess global transcriptional responses at single-cell resolution [43, 44]. Nuclear hashing labels cell nuclei with unique DNA barcodes before pooling, allowing multiple treatment conditions to be multiplexed in one experiment. sciRNA-seq uses successive rounds of combinatorial indexing to uniquely tagged transcripts from individual cells, enabling high-throughput scRNA-seq at a much lower cost. Sci-Plex was applied to A549 (lung adenocarcinoma), K562 (chronic myelogenous leukemia), and MCF7 (mammary adenocarcinoma) cell lines, exposed to 188 different compounds at four doses [41]. The integration of the two techniques allowed to simultaneously profile thousands of single-cell transcriptomes across nearly 5000 samples.

Genome-Wide Perturb-Seq (GWPS) [45] is a large-scale database describing how genetic changes affect cell behavior, by using a single-cell genetic perturbation sequencing

(Perturb-seq), a Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR) interference (CRISPRi) screening technique combined with scRNA-seq [45]. With this approach, every expressed gene was silenced, to capture detailed transcriptional responses. This massive dataset allowed to assign functions to previously uncharacterized genes, and to identify new regulators involved in key cellular processes [45].

Table 2.1: List of resources that provide transcriptomic data driven from perturbagen.

Database	Description	Ref.
ArrayExpress	ArrayExpress is a curated repository with microarray and HTS data from various perturbation studies (not involving toxic compounds). Established in 2003, it maintains high-quality standards through manual curation, offering structured metadata and accessioned datasets for diverse biological conditions, including compound treatments and diseases. It provides access to data from other sources, such as DrugMatrix [46] and Open TG-GATES [38].	[47]
CDS-DB	Cancer Drug-induced gene expression Signature DataBase (CDS-DB) is a patient-derived cancer drug response database that compiles pre- and post-treatment microarray and RNA-seq data. It provides data for studying treatment effects, drug responses, and resistance mechanisms. It includes curated metadata from GEO [48] and ArrayExpress [47] and supports data browsing, searching, and signature analysis.	[40]
CEBS	Chemical Effects in Biological Systems (CEBS) is a publicly resource for toxicogenomic data, established in 2008. It integrates multiple types of experimental data, including detailed study designs, histopathological, microarray and proteomics data, collected from studies investigating the effects of both exposure to specific compounds and genetic alterations [49].	[50]
ChemPert	ChemPert is particularly useful for comprehending the molecular impacts of chemicals in non-cancer cellular contexts. It provides 82270 manually curated transcriptional signatures, derived from 167 non-cancer cell types incubated with 2566 chemical perturbagens. It also includes the protein targets for 57818 chemical compounds, and the effects (activation, inhibition, or unknown) of those on the targets.	[39]

*Continued on next page*

Table 2.1 – *Continued from previous page*

Database	Description	Ref.
CMap / LINCS	Connectivity mapping compiled gene expression responses from human cell lines both of epithelial (MCF7 -breast cancer- and PC3 -prostate cancer) and non-epithelial origin (HL60 -leukemia- and SKMEL5 -melanoma), treated with 164 small molecules prior gene expression profiling using Affymetrix GeneChip. Building on this, the LINCS (clue.io) uses Luminex bead arrays to measure 978 reference genes and infer the expression of 11,350 additional transcripts, resulting in a database with over one million gene expression profiles, from 20000+ chemical and genetic perturbagens. LINCS data are available in Phase 1 (GEO access: GSE92742) and Phase 2 (lincsportal; GEO access: GSE70138). iLINCS [32] provides an interactive platform to explore these datasets and the relationships between compounds, gene expression changes, and disease.	[8, 11]
CREEDS	Crowd Extracted Expression of Differential Signatures (CREEDS) is a manually validated collection with GEO data, assembled from crowdsourcing [48]. The manual curation resulted in 2176 genetics, 875 chemicals, and 828 disease perturbation signatures, from both human and mice. It was later expanded with 8620 genetic, 4295 chemical, and 1430 disease perturbation signatures automatically extracted from 2543 GEO studies.	[51]
DrugMatrix	Public toxicogenomic database (drug matrix; GEO Access: GSE59927) with microarray-based gene expression profiles from rat tissues exposed to 600+ compounds and with 5000+ signatures available. The database includes repeated and single-dose studies at 6 h, 24 h, 3 days, and 5 days.	[36, 37]

*Continued on next page*

Table 2.1 – *Continued from previous page*

Database	Description	Ref.
DRUG-seq	Digital RNA with perturbation of Genes (DRUG-seq) is a high-throughput RNA sequencing platform designed for cost-effective transcriptomic profiling. It is suitable for testing multiple treatments in parallel, since it uses in-well cell lysis in 384-well plates, enabling large-scale screening of compounds. Initially, DRUG-seq was applied to osteosarcoma cells treated with 433 drugs at 8 dosages. A follow-up study further validated the platform with extensive testing, and introduced an open-source analysis pipeline. Novartis has deposited two datasets in GEO using this technique. The first contains bulk RNA-seq data from osteosarcoma cells treated with 7 small molecules at 3 concentrations for 12 hours. The second treats the same cell line with 14 compounds, each tested across an 8-point dose-response range with 3 replicates per dose. (GEO access: (1) GSE120222; (2) GSE176150; Data analysis pipeline: DRUG-seq).	[52, 53]
GEO	Gene Expression Omnibus (GEO) is a public repository for user-submitted transcriptomics data spanning a wide range of perturbation studies, diseases, and experimental conditions across various organisms and platforms. GEO is an appropriate resource for data mining and large-scale transcriptome analysis since it is updated frequently with a vast collection of datasets. GEO provides tools for depositing, querying, and retrieving gene expression and molecular abundance data. It serves as a foundation for more specialized databases.	[48]
GWPS	Genome-Wide Perturb-Seq (GWPS) is a public resource that provides single-cell genetic perturbation data generated by CRISPR interference (CRISPRi) in over 2.5 million human cells (cell lines: K562, chronic myeloid leukemia and RPE1, retinal pigment epithelial). It includes 1946 signatures, each representing a loss-of-function perturbation that triggers a strong transcriptional response.	[45]

*Continued on next page*

Table 2.1 – *Continued from previous page*

Database	Description	Ref.
Open TG-GATEs	Toxicogenomic Project - Genomics Assisted Toxicity Evaluation System (Open TG-GATEs) is a public resource that contains 1483 unique signatures from microarray-based gene expression profiles of human and rat liver tissues exposed to 170 compounds. It is built around repeated concentration time-course studies, enabling to evaluate the long-term consequences of exposure to toxic compounds. In addition to transcriptomic profiles, the database also contains histopathology, biochemistry, and hematology data.	[54]
PertOrg	PertOrg 1.0 is a public database with curated gene expression data from genetically modified organisms. It curates non-human high-throughput gene expression and phenotypic data from <i>in vivo</i> genetic perturbation experiments in eight model organisms. The perturbation includes gene knockdown, knockout, and overexpression. It currently aggregates 58707 transcriptome profiles and 10116 comparison datasets, including 122 scRNA-seq datasets, with a total of over 8.6 million Differentially Expressed Gene signatures, retrieved from GEO and ArrayExpress. This tool not only enables the retrieval of the curated data, but it also provides a platform to search and browse various genetic perturbations and to compare gene lists against their signatures, linking perturbations to pathways, cell types, and phenotypic outcomes.	[55]
PerturBase	PerturBase is a public database for single-cell perturbation data. It curated 122 datasets from 46 publicly available studies, covering 24254 genetic and 230 chemical perturbations from approximately 5 million cells, across 31 human and murine cell types. PerturBase is organized into two main modules: the Dataset and the Perturbation modules. The former is for exploring individual studies, whereas the latter for comparing perturbation effects. This resource enables downstream analysis across various cellular contexts, with a direct download option.	[42]

*Continued on next page*

Table 2.1 – *Continued from previous page*

Database	Description	Ref.
PerturbAtlas	PerturbAtlas is a resource that re-analyzes publicly available perturbed genetic gene expression signatures, including knockdown, knockout, knock in, OE, mutations, and multi-condition experiments. Currently, it provides a vast curated collection of 122801 RNA-seq libraries from 7778 studies across 13 species, sourced from ENCODE, GEO, ArrayExpress, and SRA.	[56]
Sci-Plex	Sci-Plex is a method that combines nuclear hashing with combinatorial indexing-based RNA sequencing (sci-RNA-seq), to quantify transcriptional responses to thousands of perturbations at single-cell resolution and in a single experiment. Sci-Plex screened 3 cancer cell lines (A549, K562, MCF7), exposed to 188 compounds (GEO Access: GSE139944).	[41]
STARGEO	The Search Tag Analyze Resource for GEO (STARGEO) is an open resource constructed using GEO's publicly available functional genomics data [48]. STARGEO platform provides 3031859 reliable and annotated samples of gene signatures from humans, mice, and rats.	[57]
ToxicoDB	ToxicoDB integrates and harmonizes diverse <i>in vitro</i> datasets from three sources (Open TG-GATEs [54], DrugMatrix [46], and ArrayExpress [47]). The aim is to easily perform queries and to summarize the relationships between gene expression and toxicant effects [58]. Currently, ToxicoDB encompasses curated datasets from liver tissue and three cell lines (Hep-G2, HepaRG, and Hepatocytes) in humans and rats, covering a total of 234 compounds.	[58]



GEO [48] is a public repository managed by the National Center for Biotechnology Information (NCBI), that archives microarray and next-generation sequencing data from various organisms, cell lines, etc [45]. The vast array of high-throughput experimental data stored frequently serves as the foundation for more specialized databases, making it an essential building block for several other bioinformatic databases. Additionally, querying GEO often requires metadata manual verification [33]. This results in a proportion of the publicly available data not usable, because of lack of associated metadata. To overcome this, databases using data originally from GEO often perform curation and quality checks to ensure that all the metadata are correctly provided. The Crowd Extracted Expression of Differential Signatures (CREEDS) is a database that resulted from a crowdsourcing project, to improve the annotation and reanalysis of data from the GEO [51]. Several datasets with gene, drug, and disease perturbation signatures were manually curated. Then, those signatures were used as a training set for models that automatically search GEO for more signatures with accurate metadata, allowing the database to grow efficiently. CREEDS tackles the key challenge of metadata inconsistencies and lack of standardized annotation by using both manual curation and computational methods [51]. The platform enables the download of only human-validated gene expression signatures, excluding the signatures from automated signature extraction tools.

PertOrg 1.0 [55] is another database built on extracted data from GEO [48] and ArrayExpress [47], allowing the analysis and download of curated gene expression and phenotypic data from genetically modified organisms. This extensive database contains induced *in vivo* genetic perturbations across eight diverse model organisms, including mammals (mouse and rat), non-mammalian vertebrates (zebrafish), invertebrates (nematode worm and fruit fly), microorganisms (bacteria and yeast), and plant (thale cress). The database covers various types of genetic modifications, such as gene knockout (complete removal or inactivation of a specific gene), gene knockdown (partial suppression of gene expression using Ribonucleic acid (RNA)i or similar techniques), gene OE (increased expression of a target gene through vector-based methods), mutations and other genetic modifications (specific point mutations, insertions, or deletions introduced to study gene function and disease models). PertOrg has identified over 8.6 million DEGs associated with genetic modifications, derived from microarray, RNA-seq and scRNA-seq. The database has two main built-in tools: the differential gene overlapping analysis, which investigates perturbation datasets significantly enriched in the user-provided gene set via a hypergeometric test, and dataset enrichment analysis for perturbation datasets where the user-provided gene set is over-represented [55]. The huge collection of curated non-human data and all the functionalities provided make this a valuable resource, bridging the gap between genetic disorders and their phenotypic outcomes.

Perturbation signatures provide important information about the compound's effect, through gene expression changes. However, finding the key targets and pathways involved in the compound's action can be achieved by integrating transcriptomic data with PKN. Molecular networks provide additional information and context that explain the

underlying mechanisms for the gene expression changes.

## 2.3 Prior knowledge Network

To understand a drug's MoA, we should consider the molecular interactions within a system. Computational methods help the integration of molecular entities's interactions with omics data [9]. These interactions can be represented with more or less complexity and can be included in the analysis as supplementary data sources. A PKN is a collection of interactions where nodes represent molecular entities (such as proteins, genes, or metabolites) and edges illustrate their relationships. Understanding causal graphs is key to modeling and interpreting these networks, as they depict cause-and-effect relationships. In such graphs, nodes represent variables, while directed edges represent causality, indicating that a change in one variable affects another [59]. Furthermore, edges can be signed, indicating whether a causal node employs a positive or negative effect on the second variable, and weighted, to show the connection strength [59]. In causal graphs that model biological networks, multi-edge connections are common, with two or more edges linked to the same node.

Molecular networks can be classified based on interaction types and node characteristics. Protein-Protein Interaction (PPI) networks show direct interactions between proteins. GRN (Figure 2.2) illustrate how TFs influence gene expression [60]. Signal transduction networks describe how cells process external signals. Metabolic networks display relationships between enzymes and metabolites. Furthermore, networks are not always composed of molecular entities, as is the case with the disease networks, which links diseases using genes and mutations as connections. These networks fit experimental data to predictions from causal graphs describing the system. The choice of PKN should match the data type. For instance, for transcriptomic data, integrating a PKN may be beneficial, while for metabolomic data, metabolic networks are more suitable.

Researchers have made significant efforts to construct regulatory networks. A primordial example was the functional characterization of yeast genes through PPI analysis. This study leveraged the guilt-by-association principle, by inferring an unknown protein's function by looking at its interactions with nearby entities [59, 61]. The guilt-by-association principle is a foundational concept in biological networks. It suggests that genes with similar functions often interact with the same proteins or have similar expression patterns [62]. This principle also applies to drugs that cause similar transcriptional responses and may have comparable MoA [20].

Biological interactions can be described with different levels of complexity, such as full networks, pathways, and regulons. A network is an intricate representation of the global interactome, linking all entities in the system. These interactions are initially characterized in published experimental research with a varying amount of associated information. Based on supporting studies, each interaction may include details about direction, signal, and confidence level. This helps filter data, creating a network with more

reliable relationships. Still, networks can be noisy and incomplete, with false positives and false negatives, and a tendency for well-researched entities to become overrepresented [9, 63, 64]. In contrast, a pathway is a simpler version of a network. It illustrates a series of molecular interactions that begin with one entity and follow a specific signaling cascade. This arrangement helps classify entities by their common biological roles. Yet, pathways often do not capture crosstalk among them, providing a static view of what is a dynamic process [9]. The entities' overrepresentation issue also applies to pathways. Another way to show interactions is through regulons. Regulons are groups of co-regulated genes controlled by a common TF and are usually represented as GRN (Figure 2.2). The choice of network type should be tailored on the scientific question and the type of data with which the network is used. The type of interactions between molecular entities and its complexity should be considered. For example, if pathways are used instead of a full network, they might miss some interactions. On the other hand, if a study targets a specific cell type or tissue, it's important to use tissue-specific networks. Databases like TissueNet [65] provide molecular interactions specific to a certain cellular context.

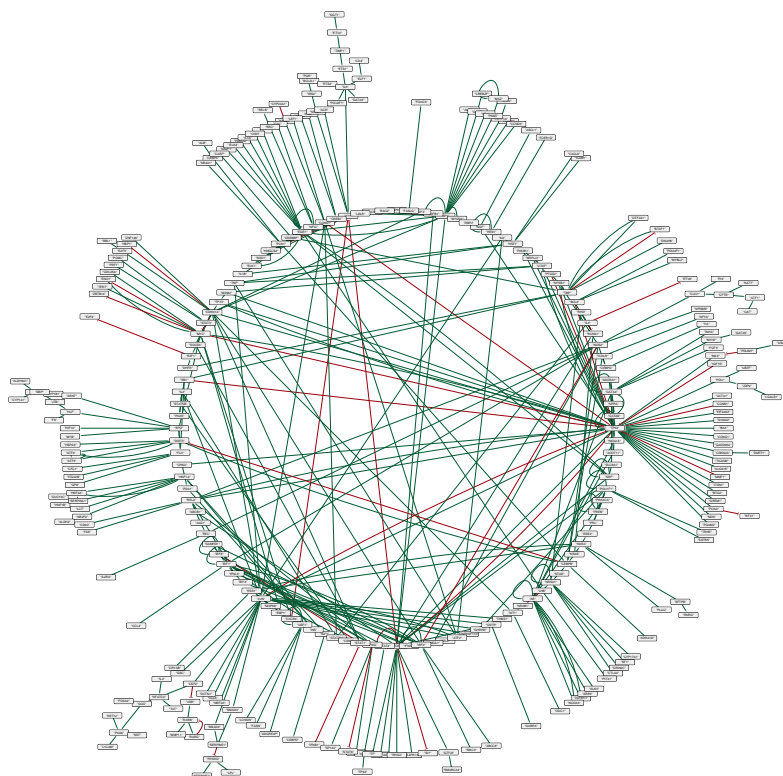


Figure 2.2: GRN based on regulon representation of human transcriptional interactions. CollecTRI-derived regulons were extracted from the decoupleR (v. 2.12.0) package, which provided 43,178 interactions. CollecTRI collection [60] is a comprehensive, curated resource of TFs and their target genes, expanding on DoRothEA. This figure illustrates a subset of those 1,000 interactions, with edge color indicating the mode of regulation (green for activation and red for inhibition) using RCy3 (v. 2.26.0) R package for visualization.

The use of causal networks with transcriptomic data was first introduced by Pollard *et al.* [63, 66]. This study aimed to infer the molecular causes of the changes in oxidative

phosphorylation gene expression in skeletal muscle from type Diabetes Mellitus Type 2 (DM2) patients. For this purpose, the gene expression data were integrated with a large-scale model created from over 210,000 molecular relationships based on DM2 literature. Computer-aided causal reasoning on these complementary data identified that the observed changes are linked to decreased glucose transport, impaired insulin signaling, and increased risk of post-transplant diabetes [63].

Given the good results obtained from supplementing the studies with PKN, the identification of interactions began to receive more attention. The development of HTS technologies [67] allowed the construction of PPI. Which in turn have begun to be deposited in databases that provide molecular interaction data. Nowadays, there are several public and commercial network and pathway resources, some of them summarized in Table 2.2. Two resources that provide composite public and commercial networks are OmniPath and MetaBase, respectively. OmniPath [68] is a freely available resource of prior knowledge in molecular biology. It combines data from over 100 resources and builds five integrated databases with different types of data: Interactions (several molecular interactions organized into sub-networks), Post-Translational Modifications (enzyme-substrate reactions), Complexes (35,000+ protein complexes), Annotations (proteins and complexes annotations, such as the function, localization, tissue, etc.) and Intercell (inter-cellular signaling) [68]. The interactions database is a composite signaling network that offers several manually curated subnetworks, encompassing a total of 282,504 unique interactions. Each subnetwork has different types of interactions, including post-translational interactions, transcriptional interactions, post-transcriptional interactions, and other interactions involving small molecules. The number of interactions per subnetwork is described in Figure 2.3. One of the GRNs that is provided by this database is the CollecTRI-derived regulons [60] (Figure 2.2). This collection contains high-confidence signed TF - target gene interactions compiled from 12 resources. OmniPath data can be accessed through the OmnipathR R package [69].

MetaBase<sup>TM</sup> [70] is a proprietary, commercial database from Clarivate that offers one of the most comprehensive, manually curated systems biology datasets available. It contains over 4.2 million molecular interactions, including protein-protein, protein-RNA, compound-protein, compound-compound interactions, and transport reactions, with details on directionality, mechanisms, and effects. In addition, MetaBase provides more than 1,500 pathway maps that cover regulatory, disease, metabolic, and toxicity characteristics, alongside over 10,000 disease-related networks and 1,000+ validated networks. Each interaction is assigned a trust score that reflects its reliability, helping users distinguish well-established interactions from those obtained via HTS. MetaBase can be accessed via metabaseR package in R, which simplifies visualization, functional analysis, and network manipulation. Furthermore, the CBDD R package offers 73 advanced algorithm implementations for analyzing and extracting insights from these networks.



Figure 2.3: OmniPath molecular interactions database. This figure presents a hierarchical visualization of human molecular interactions curated by OmniPath, which compiles data from diverse resources. The graph displays the numbers of distinct interaction types, including post-translational interactions (PPI), transcriptional interactions (GRN), post-transcriptional interactions (such as miRNA-mRNA, TF-miRNA, and lncRNA-mRNA interactions), and other interactions involving small molecules (encompassing drug-target, ligand-receptor, enzyme-metabolite, among others).

Table 2.2: List of resources that provide prior knowledge networks and pathways.

Database	Type of data	Description	Ref.
BioSNAP	Network	The Biological Stanford Network Analysis Project (BioSNAP) includes multiple biological networks encompassing different entities and relationships, such as a disease-drug association (1,334,088 edges); side effects-drug relationships (4,649,441 edges); tissue-specific protein-protein interactions (70,338 edges), physical protein-protein interactions identified in human (342,353 edges) and many more. BioSNAP also provides a tool, Mambo, for the construction, representation, and analysis of large multimodal networks.	[71]
IntAct	Network	IntAct is a public molecular interaction database with over 1,600,000 interactions, derived from literature curation (23,000+ publications) and user submissions (75,000+ experiments) for multiple species.	[72]
KEGG	Pathways	Kyoto Encyclopedia of Genes and Genomes (KEGG) primarily features metabolic pathways but also includes signal transduction and disease-specific data. KEGG Pathway consists of manually curated pathway maps, illustrating molecular interactions, reactions, and networks across various domains, such as metabolism, cellular processes, human diseases, drug development, etc. KEGG covers 578 human pathways, 12,629 drugs (including 2,499 drug groups), and 2,894 human diseases.	[73]

*Continued on next page*

Table 2.2 – *Continued from previous page*

Database	Type of data	Description	Ref.
MetaBase	Network	MetaBase (MetaBase) provides commercially available manually curated networks, larger (higher number of nodes) and denser (higher number of connections) than other publicly available databases. It contains 4.2 M+ molecular interactions with directionality, mechanism, and effect. Although this is primarily described as a network resource, MetaBase also includes curated pathway maps. It covers human, rat, and mouse genes. metabaseR is a software package that facilitates access to MetaBase content in R.	[70]
OmniPath	Network	OmniPath is a freely accessible resource that integrates molecular biology data into five main databases: Interactions, Post-Translational Modifications, Complexes, Annotations, and Intercellular. The interactions database contains high-confidence data from different sources encompassing 282,504 unique interactions, organized into various subnetworks. Software tools are also provided for R (OmnipathR) and Python (pypath), together with a plug-in for network visualization (OmniPath Cytoscape) [69, 74].	[68]
Pathway commons	Pathways	Pathway Commons is a public repository of pathway and interaction data from 23 databases formatted in the BioPAX standard. It currently includes 6,692 pathways and 3,579,336 interactions from sources like KEGG, IntAct, BioGRID, Reactome, etc. It provides an R package [75], and the CyPath2 plugin for pathway visualization in Cytoscape.	[76]

*Continued on next page*

Table 2.2 – Continued from previous page

Database	Type of data	Description	Ref.
Reactome	Pathways	Reactome is a public and reviewed pathway database that is manually curated. It includes 2,769 human pathways, 23,911 interactions, 11,574 proteins, and 1,057 drugs, with over 40,000 literature references. The pathways are organized under 29 categories, such as cell-cell communication, signal transduction, etc. Reactome offers a Pathway Browser for visualizing and interacting with the data, as well as an integrated Analysis Tool [77] for pathway identifier mapping, over-representation, and expression analysis. It also provides two R packages, ReactomePA [78] for pathway analysis and ReactomeGSA [79] for Multi-Omics Comparative Pathway Analysis, and the ReactomeFIPlugIn for pathway visualization in Cytoscape.	[77]
STRING	Network	STRING is a public database that includes known and predicted PPI. It currently covers over 59.3 million proteins from 12,535 organisms, providing a comprehensive PPI network with confidence scores based on various data sources.	[80]
WikiPathways	Pathways	WikiPathways is another open-source curated pathways database. It provides pathways for 27 organisms and includes 959 curated human pathways. The platform supports collaborative annotation and refinement of pathways. It also provides several software tools and workflows in R (rWikiPathways) and Python (pywikipathways), along with Cytoscape and PathVisio plug-ins for pathway visualization, among others [81]	[82]



Integrating biological knowledge with experimental data is key to understanding how cellular regulation impacts gene expression. Known interaction networks are usually used to predict the results of regulatory events, but they can also be used in the opposite direction, to find upstream regulators that cause expression changes [63]. Computational tools play a crucial role here. They combine high-throughput omics data with established cellular interactions, like PPI and signaling pathways, to give a broader context. While network data show the complete interactome of molecular interactions, pathway data arrange these interactions into cascades. Each of these data sources forms prior knowledge. When combined with experimental results helps to create mechanistic hypotheses about, for instance, how a perturbation works in a system. This integration of experimental and interaction data sets the stage for some of the computational methods covered in the next section. These methods aim to uncover the mechanisms that cause an observed transcriptomic change.

## 2.4 Computational methods for MoA inference

Due to technological advances, large-scale transcriptomic datasets can now be generated affordably for many perturbations. However, extracting biological insights from this data can be complicated, requiring dedicated computational tools for the analysis. These methods include *in silico* experiments that combine experimental data with prior knowledge. They fall into three main categories: topology-, similarity-, and enrichment-based methods [16, 83]. Choosing the appropriate tool depends on several factors, including the type of input data, runtime, computational complexity, and the specific scientific questions being addressed, along with the inherent strengths and limitations of each method [9]. For example, in Hill *et al.*'s study, [59] three types of topology-based algorithms were evaluated. Node prioritization algorithms rank the nodes in the network based on the distance from start nodes and the connection between them; causal regulation algorithms infer and rank upstream nodes in the network by combining the direction of edges, and subnetwork identification algorithms extract regions of the input network that are enriched for perturbed nodes. These approaches help generating mechanistic hypotheses about the cellular targets and pathways affected by a perturbagen, more efficiently and accurately [84].

### 2.4.1 Connectivity-based tools for comparative analysis

Similarity-based approaches for CMap focus on matching gene expression signatures from query to reference signatures. The approach emerged from the need to connect changes caused by perturbagens, with those observed in diseased or other biological states. As stated by Lamb *et al.* [11] the resource was named CMap, due to its potential to connect drugs, genes, and diseases, with the foundational idea that if a compound induces a transcriptional signature similar to a known gene expression pattern, it likely shares

the same MoA or therapeutic effect. The two main components of this approach are: a query signature, a ranked list of genes up- or downregulated in a condition of interest and a reference database. Using a pattern-matching algorithm, the query signature is systematically compared against a reference database of perturbation-induced expression profiles [85]. The seminal study that proposed the first connectivity-based approach also introduced the first CMap database, described in section 2.2. As a strategy to score the similarity between data, the authors employed a nonparametric, rank-based Kolmogorov-Smirnov (KS) test [11, 86]. This approach has been adapted and extended by several methods, including those based on weighted KS statistics, such as the variations of the Gene Set Enrichment Analysis (GSEA), and alternative metrics such as the eXtreme Sum (XSum) score and signed rank-based methods like the ZhangScore [86].

For each reference profile, the algorithm evaluates whether query upregulated genes cluster at the top of the ranked list and downregulated genes at the bottom, indicating a positive connection, or a negative connection. This yields a connectivity score between 1 (strong positive correlation) and -1 (strong negative correlation), with scores near zero indicating no significant association [11]. In this way, these methods not only generate the direction of correlation, but also provide information on the strength of the link. The similarity scores can be directly linked to biological interpretations. A high positive score may indicate that a compound could enhance a biological condition, while a high negative score suggests a potential inhibitory effect [46].

CMap has become a powerful tool for DR and mechanistic studies because it leverages large-scale gene expression data, to connect drugs, diseases, and gene perturbations. By comparing a query signature with a database of perturbagen signatures, candidate drugs that may reverse or have similar transcriptional changes patterns can be identified [46]. The CMap concept not only facilitated the creation of several drug-induced molecular perturbation signature databases [40], but also supported studies applying similarity-based algorithms in DR and MoA elucidation [39].

### 2.4.2 Enrichment-based tools for downstream analysis

As the number of gene expression studies increases, it becomes harder to extract relevant information from it. Pathway analysis helps to frame large changes in gene expression that, if isolated, lack biological context. These methods convert gene lists into a meaningful and interpretable biological process, by linking expression data to specific biological pathways.

Biological pathways are identified by characterizing the cascade of interactions occurring in the system in response to a given perturbation. These interactions can be observed in transcriptomic data and translated into significantly enriched pathways, to capture how changes in gene expression are preferentially affecting some pathways more than others. However, pathway analysis does not trace the causal paths leading to observed gene expression. One commonly used method to detect these relevant pathways is GSEA. The

input for GSEA is a ranked list of genes, based on metrics like fold-change and a predefined gene set. The algorithm then produces an Enrichment score (ES) by determining whether each gene from the gene sets clusters at the top or bottom of the ranking list. The statistical significance is evaluated using a KS test [87]. Since it is a rank-based approach, there is no need to use thresholds (for example, for fold-change), a full gene expression profile can be used, not only DEGs. This reduces bias and increases sensitivity to pathways where individual gene changes are small but consistent. However, this thoroughness comes with a high computational cost.

Over-Representation Analysis (ORA), on the other hand, is a simpler and faster method for pathway enrichment. It is ideal for large-scale screening and initial data exploration. This analysis uses a hypergeometric test or Fisher's exact test, a background gene set and list of genes of interest, meeting a specific fold-change or statistical threshold (such as DEGs). The statistical test compares how each pathway for each gene in the list overlaps with the input gene list. Since ORA only needs a non-ranked gene list, it can be more straightforward than GSEA. However, by relying on a cutoff for the input list the pathway detection is affected by a certain degree of subjectivity. Additionally, ORA does not consider the magnitude or direction of gene expression changes and, by assuming independence among genes and pathways, it may overlook complex interactions.

Overall, enrichment-based methods help provide a broader and better view of biological processes affected downstream of transcription. However, they lack the context of mechanisms behind observed changes. One way to understand the how and why those changes are observed in the gene expression data is through topology-based algorithms. These algorithms can combine transcriptomic data with network information to predict drug or disease targets.

### **2.4.3 Topology-based methods for upstream analysis**

While pathway enrichment methods compare gene expression data with the respective encoded proteins and the signaling pathways, on the other hand, topology-based methods treat a gene expression pattern as the result of a specific perturbation. Several algorithms fall into the category of topology-based methods, as they take advantage of a topology network as prior knowledge. A subset of them is categorized as causal reasoning algorithms. Causal reasoning can be described as the process of looking at what happened, the effects, and trying to infer the upstream causes. In this context, change in a process in response to specific stimulus is considered as a proxy for causality, in line with the following axiom A causes B if a change in A leads to a difference in B, assuming everything else stays the same [88]. This thinking has deep philosophical roots, and, in systems biology, it has a more testable and specific meaning. In early medical applications of causal reasoning, A represented the treatment, and B was the outcome.

In DD studies, this concept can be applied by overlaying high-throughput measurements, such as gene expression data, onto a network. The algorithm itself is a sophisticated

causal node prioritization tool that makes predictions based on network topology. It requires direct interactions, meaning directed cause-effect relationships, connecting each pair of biological entities. Additionally, it uses edge effects, to know (at least for some edges) whether it represents activation or inhibition. For DEGs, it doesn't necessarily require the exact fold-change values; instead, it is sufficient to specify which genes are up- and down-regulated. At its core, a causal reasoning algorithm works backward. It starts with a node in the network and follows it downstream for a few steps, predicting what would happen if this node were activated (or inhibited) in the biological system of interest. For example, by analyzing the network structure, it can determine if the activation of a specific regulatory protein will also activate a certain TF, directly affecting the expression of target genes in a positive or negative way (respectively by activating or repressing gene expression). The algorithm makes these predictions and then compares them with what is observed in the data. This helps identify which upstream regulators best explain the downstream changes that are seen [10]. However, not every edge in a prior network is active during a specific experiment, and not all predicted downstream effects occur. This is the reason why statistical scoring is crucial.

In causal-network inference, all evaluation methods compare each regulator predicted downstream effects against the observed gene expression changes, but the exact scoring metric depends on the type of network used [89]. The simplest approach just counts how many predictions are correct and incorrect. A more refined ES uses Fisher's exact test to evaluate whether the targets of a regulator are enriched among the DEGs while ignoring whether each edge is activated or inhibited. To capture that extra layer of information, the Quaternary scoring (QS) computes a z-score to quantify how well the predicted up/down directions match the observed up/down changes in the data [90]. Early work by Pollard *et al.* [66] combined the two methods, computing an overlap and concordance p-values for regulators in type-2 diabetes expression data [90]. However, Fisher's test can not leverage signed edges even when those signs are known. To address this, Chindelevitch *et al.* [63, 91] developed a Ternary scoring (TS) method that models the signed network and observations as a dot-product distribution, allowing exact p-values for both the activation and inhibition hypotheses of each regulator. This scoring method requires a fully annotated causal graph, i.e., every edge must have a known direction and effect information. More recently, Fakhry *et al.* [92] introduced a QS method that extends TS to networks containing both signed and unsigned interactions, preserving directional inference whenever possible while still accommodating unannotated edges [89]. Early methods used simpler calculations, while newer methods used more precise statistical tests to ensure accuracy. However, simpler z-score and enrichment-based methods are still popular because they are faster to compute and easier to use and understand.

Causal reasoning algorithms are widely used but have many nuances depending on the network's structure and how far one can trace from regulators to downstream genes. Improving its accuracy in predicting regulators is a key focus. Also, they are complex and computationally demanding, especially with the increase in input data.

There are evidences that topology-based methods may be more useful in DD than pathway enrichment methods. Topology-based causal reasoning takes advantage of directed and signed Molecular networks to infer the most plausible upstream regulators from transcriptomic endpoints. By inferring key signaling nodes whose activity can directly explain the changes observed in the experimental data, these methods pinpoint candidate drug targets and generate mechanistic hypotheses for further validation. On the other hand, pathway enrichment methods only infer which biological processes are affected by changes in gene expression. These methods do not consider that after transcription, events such as translation and post-translational modifications will also affect protein activity [29].

## 2.5 Benchmarking of computational methods for MoA inference

The use of computational methods to elucidate MoA is becoming increasingly indispensable for integrating and interpreting multidimensional biological data. Given the plethora of existing computational tools, choosing the appropriate data and methods to answer specific scientific questions can be challenging. When a new tool is developed and published, it is usually benchmarked against popular existing methods. Without a deep expertise in the area, distinguishing the benefits of novel tools can be difficult.

A comprehensive benchmarking study is crucial for evaluating available methods in a standardized way, providing sufficient information to accurately choose the best tools and data for a given study [93]. A key component of a benchmarking study is the use of gold-standard datasets, against which results are compared. By this process, often formalized via well-defined good practises, it is possible to evaluate performance metrics and statistical analyses and to compare different algorithms on specific types of data.

Benchmarking studies can be carried out by the authors that implemented the tool, independent groups, or as organized challenges, such as those by the Dialogue on Reverse Engineering Assessment and Methods (DREAM) Consortium [94]. If the evaluation is performed by the authors, the aim is usually to demonstrate the advantages and performance improvements over other techniques. In other contexts, it is important to define the benchmark's scope and purpose. The selection of the methods should reflect the relevance of the study's objective and include publicly available implementations to ensure accessibility. Parameter optimization can significantly affect a tool's behavior, including runtime, yet finding the optimal values is not always straightforward. Thus, balancing default settings with computational efficiency is important. Regarding datasets, in MoA studies, it is crucial to include diverse data sources, to ensure representativity and a credible assessment of performance. For instance, transcriptomic data should ideally include both bulk RNA-seq and scRNA-seq data to broaden the options and use two widely known types of data. Since there are no perfect, fully curated datasets, it is necessary to ensure quality by choosing, for example, resources that provide manually curated data, to avoid biasing the results and performance of the tools [94]. The same applies to the gold

standard datasets, which serve as the ground truth, and they represent the essence of a benchmarking study.

Other similar studies arise from an effort to contextualize gene expression data with several computational algorithms. Hosseini-Gerami *et al.* [29] evaluated the performance of different causal reasoning algorithms to find the targets of small molecules and associated signaling pathways using gene expression data. The study compared four causal reasoning algorithms against networks from two diverse sources and transcriptomic data from one database. Hill *et al.* [59] conducted a study that provided a more comprehensive framework, by analyzing a diverse range of algorithms, networks, and datasets, to assess how well network-based algorithms prioritize and connect gene lists derived from transcriptomic data. This study integrated 17 algorithms, categorized into three main groups: (1) Node prioritization algorithms, which rank network nodes based on their connectivity, (2) Causal regulator algorithms, which identify upstream regulators of transcriptomic changes, and (3) Subnetwork identification algorithms, which extract subnetworks linking input genes. The algorithms were applied to three PPI networks, each with different structures and levels of curation, using hundreds of datasets from four sources. The first network combined data from various databases, resulting in a mix of signed/unsigned and directed/undirected interactions. The second network included only signed, directed, and high-confidence interactions, while the third was a large-scale, undirected PPI network. This study exemplifies good practices when comparatively analyzing different algorithms and integrating different types of input data, although it does not explore the algorithms' parameter landscape (considered beyond the scope of the work).

Typically, a benchmarking analysis ranks the algorithms in terms of the most appropriate use for distinct applications [59]. By providing a robust assessment of the capabilities of existing algorithms, these studies provide guidelines to choose the most appropriate tool for the scientific question of interest [93].

## MATERIALS AND METHODS

*The following section describes the methodologies and workflow of the benchmarking study. It begins by describing the input data used, both transcriptomic data and prior knowledge data. In addition, it details the implementation of the selected algorithms and their execution. Finally, the metrics used to evaluate their performance are described.*

### 3.1 Benchmarking architecture setup

Several tools and algorithms are available for most research tasks in computational biology, and new algorithms and tools are published every week. At the same time, finding the right computational tool for a given research question is essential. Systematic benchmarking of tools is a time- and resource-consuming endeavor, while a lack of benchmarking carries several potential risks. Researchers usually carry out published benchmarking to demonstrate that their tool performs better than others. ABC is a consortium created in 2021 to help members reduce R&D risks, saving time and resources by distributing the effort of benchmarking computational biology algorithms. ABC maintains the same workflow regardless of the case study. It consists of three main steps: (1) Voting, (2) Curation, and (3) Coding. The consortium members suggest and vote on the use case (1). Once the use case is determined, the curation phase begins, where Clarivate collects the most appropriate datasets and algorithms according to the voted use case (2). Again, the members vote on the final selection of datasets and algorithms (1). Finally, the last phases - implementation, execution, and reporting - are conducted by Clarivate (3). The project description will fall within the third phase of the workflow, specifically concerning the algorithm's implementation and execution, where I actively participated since all the data was already collected and voted on when I started the project. A visual representation of the study workflow is provided 3.1 and will be explained in detail in the following sections. The implementation of ABC causal regulation use case was carried out by Dr. Alexandr Ishkin and me. The entire workflow was conducted in R Statistical Software (v4.4.1) [95], and run on a server with an AMD EPYC 7R13 processor featuring 64 logical CPUs and 246 GB of RAM.

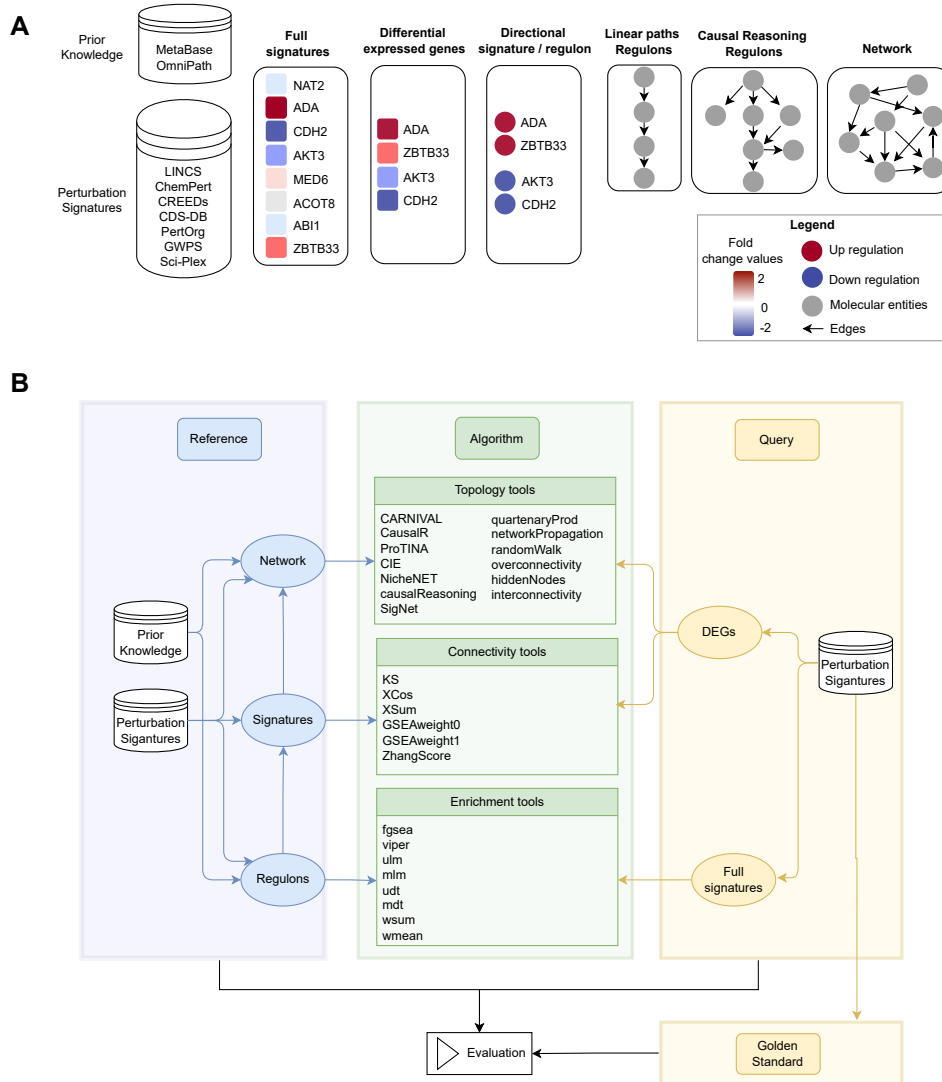


Figure 3.1: Schematic of the study architecture. A. Perturbation signatures collected from seven public sources are used in the benchmarking framework either as reference, query, and gold standard (known targets) datasets. Prior knowledge networks, used as reference, were derived from two sources: OmniPath (public) and MetaBase (commercial). From OmniPath, a global network, and regulons were used as references. From MetaBase, it was also used a full network, regulons, and linear pathways. B. Three classes of computational methods were evaluated: topology-based, connectivity-based, and enrichment-based, comprising a total of 27 algorithms. Depending on the method, input data may consist of a global interactome (network), curated signaling pathways, or perturbation signatures (typically directional gene sets or full transcriptomic profiles, which can be reduced to gene sets if needed). These input types are often interrelated, and the arrows in the diagram indicate the required data transformations specific to each algorithm. The output of each method is systematically compared to the gold standard targets for evaluation.



## 3.2 Data Description

As represented in Figure 3.1, each algorithm should receive two types of inputs: query and reference dataset. The query dataset refers to the data derived from perturbed signatures (Full profiles or DEGs lists). The reference dataset can be derived either from perturbed signatures (Full profile, DEGs, or directional/regulons) or from prior knowledge (Networks or pathways/gene sets). The databases and datasets used as perturbed signatures and as prior knowledge are described below.

### 3.2.1 Gene expression data: Perturbation signatures

Currently, there are several publicly available perturbation-driven gene expression datasets. This study comprehends transcriptomic datasets from seven different public sources, summarized in Table 3.1. Chemical and genetic perturbagens were included, analyzed by bulk microarray, bulk RNA-seq, and scRNA-seq assays. Each dataset contains more than hundreds of perturbation signatures. For each collection, the perturbagen type, the total number of unique perturbagens profiled, and the subset for which a gold standard target annotation is available were recorded. The gold standard is necessary for the evaluation and it consists of a set of known targets (drug-protein interactions or genes deliberately perturbed), used to assess the ability of the algorithms to recover true upstream regulators from observed expression changes. The LINCS expands upon the original CMap by leveraging the cost-effective L1000 platform, which directly measures 978 landmark transcripts and imputes the remaining transcriptome to reconstruct genome-wide expression profiles. LINCS comprises several distinct collections of perturbations in human cell lines: over 30,000 unique small-molecule treatments, CRISPR knockouts targeting 5,156 genes, cDNA OE of 3,780 genes, and shRNA knockdowns of 4,854 genes. The level 5 data were retrieved from the CLUE platform (available at CLUE). This level already contains the differential expression signatures with z-scores aggregated across biological replicates without p-values. Since each perturbagen usually appears under multiple conditions (different doses, time points, and cell lines), these were condensed into a single consensus signature per perturbagen by extracting every available gene's z-score and then using the median value across signatures. For the gold standard, directional effects were assigned as follows: for chemical perturbations Cortellis Drug Discovery Intelligence (CDDI) annotations were used to identify molecular targets inhibited by specific compounds; for CRISPR and shRNA datasets, the target genes were assigned with inhibition effect, and for OE, each target was assigned with activation effect.

ChemPert is a manually curated compendium of 82,256 gene expression signatures derived from non-cancer cell compound perturbation experiments. Most signatures originate from bulk expression studies in various cell lines, and each is represented as a list of DEGs indicating only up- or down-regulation (no fold-change values or p-values). From the total number of signatures, only 2,587 have distinct compounds. A

set of consensus DEGs lists were derived to reduce redundancy and runtime. For each compound, only genes appearing as DEGs in at least two signatures and with the same regulation direction were kept. As well as only signatures with at least 50 consensus DEGs. This resulted in 1,304 signatures which was the dataset used instead of the original ChemPert.

CDS-DB contains 78 cancer patient-derived, paired pre- and post-treatment transcriptomic datasets, all with associated metadata such as drug dosages, sampling times, and locations. 181 study-level gene perturbation signatures (85 therapeutic regimens across 39 cancer subtypes) were extracted. The perturbagen consists of drugs, and the expression is measured by microarray or RNA-seq (including fold change and p-values).

The sci-Plex dataset is based on a single-cell transcriptomics method that uses nuclear hashing. Sci-Plex dataset profiled three cancer cell lines treated with 188 small-molecule compounds. The data contains full transcriptomic signatures with around 11,000 genes each, containing dose-response effect estimates and associated p-values. Only signatures linked to compounds with known CDDI targets were kept. For each of the 135 perturbagen with a target, the gene expression responses were measured across 3 cell lines, resulting in 405 signatures.

CREEDS is a crowd-sourced, manually curated collection of perturbation signatures from GEO. Includes both small-molecule and genetic perturbations in mouse and human with the expression from different bulk gene expression platforms. These signatures are represented as DEGs lists indicating the regulation direction without fold change values. Only perturbations with CDDI target annotations were retained, and all mouse data were mapped to human orthologs, using the metabaser package.

PertOrg is a curated collection of *in vivo* genetic perturbation (such as knockdown, knockout, and overexpression) signatures across 8 model organisms. Only mouse signatures with more than 5,000 genes were kept and mapped to human orthologs. For the gold standard, perturbation effects were considered as activation in the case of knock-in, overexpression and activation, and inhibition for all the remaining ones. Since PertOrg originally contained 7,398 signatures but only 2,321 distinct target genes, a filtering criteria was applied. Each signature should have at least 50 DEGs, and the target gene's fold change should be ranked in the top 5% by absolute value among all measured genes within that signature. Then, the selected signature was the one with the highest number of DEGs per target and perturbation type combination. This resulted in 951 signatures used as PertOrg dataset.

The GWPS dataset represents a large-scale effort for single-cell CRISPRi profiling across more than 2.5 million human cells. It targets 9,866 genes and was generated using the 10x Genomics platform. The dataset includes 1,946 perturbation signatures corresponding to gene knockdowns. Each signature consists of full transcriptomic profiles by z-scores without p-values. Although the DEGs per signature were also provided by the authors, only the full signatures were used in the analysis.

The concept of causal inference can be described as the ability of algorithms to find

Table 3.1: Summary of gene expression datasets used in this study. Each dataset includes transcriptomic signatures derived from chemical or genetic perturbations. The Perturbagen column specifies the type of compound or gene perturbation applied, while the Type column labels it as either chemical or genetic. The Number of perturbagens refers to the unique compounds or genes perturbed in the dataset. The Signatures with the gold standard column indicate how many signatures have associated targets that can be used for benchmarking. The Signature type describes the format and content of the signature, such as full profiles or DEGs lists.

Data set	Perturbagen	Type	Number of perturbagens	Signatures with gold standard	Signature type	Ref.
LINCS compounds	Compound (small molecules)	Chemical	33,627	3,540	Full	[8]
ChemPert	Compound (small molecules, ligands, drugs)	Chemical	2,508	1,304	DEGs (up/down gene sets)	[39]
CDS-DB	Compound (small molecules) Patient-derived	Chemical	181	181	DEGs (Full)	[40]
Sci-Plex	Compound (Single cell; Different doses)	Chemical	189	405	Full (scRNA-seq)	[41]
CREEDS	Disease, small molecules, single gene perturbations	Chemical Genetic	3051 (875 drugs, 2176 genes)	2,642	DEGs (up/down gene sets)	[51]
LINCS CRISPR	CRISPR KO	Genetic	5,156	5,156	Full	[8]
LINCS OE	cDNA over-expression	Genetic	3,780	3,780	Full	[8]
LINCS shRNA	shRNA interference	Genetic	4,854	4,854	Full	[8]
PertOrg	shRNA interference; CRISPR knockdown; Over-expression	Genetic	7,398	951	DEGs (up/down gene sets)	[55]
GWPS	CRISPR interference	Genetic	1,946	1,946	Full (scRNA-seq)	[45]

the target candidates of a perturbation, based on gene expression data generated from a specific experimental study. Each dataset described in this section feeds into the benchmarking workflow as the query (or reference dataset) and as a gold standard (signature associated with the set of known targets). Golden standard target annotations are mandatory, not for running the algorithms, but for the evaluation step. During the evaluation, the performance of each algorithm will be assessed based on how well the targets were identified. When using signatures derived from drug perturbation, it can be hard to identify the exact compound used only from gene expression. Instead, it's easier and more meaningful to infer the target(s) of the compound (i.e. the biologically active protein that the drug binds to). Although MetaBase also contains compound information, most networks do not, but they do include gene or protein targets that can be used as proxy instead. Even for connectivity scoring methods, knowing the drug targets helps when querying compound perturbations versus gene perturbation references (or vice versa). Five chemical perturbation datasets (LINCS compounds, ChemPert, CDS-DB, Sci-Plex, and CREEDS) were subjected to this mapping through three approaches. (1) The authors' target information was extracted from the dataset/database whenever possible, including all target gene symbols. (2) Small molecules were mapped against the drugs in

the CDDI database, then, depending on the annotation level, one or more of the following information were included for downstream analyses: target drug annotation, names, synonyms or structural information. (3) The target lists provided by the authors and the one from CDDI were then merged to form the final set of targets for each therapeutic agents.

### 3.2.2 Prior Knowledge: Interaction Networks

One of inputs that can serve as a reference is the prior knowledge data, required for contextualizing gene expression signatures. The benchmarking framework depends on three complementary types of this data: PKN (global interaction networks), regulons (regulator-target gene sets), and pathway-derived linear maps (Table 3.2). Although these resources vary in their coverage, they are interconnected, as illustrated in Figure 3.1. Including sources of different sizes and densities is particularly important for understanding how the performance of topology-based algorithms is affected by the type of the input. Additionally, an increase in network size can also introduce noise that may disturb the extraction of biologically relevant information.

The interactions are obtained from two databases, OmniPath [68] and MetaBase [70]. OmniPath is a public database with protein-protein, transcriptional, and RNA-related interactions. MetaBase is a manually curated systems-biology database, provided by Clarivate, containing over 4 million directional molecular interactions, such as PPI, protein-RNA, compound-protein, etc. From each of these two sources, PKN and regulons were obtained and used as input in the benchmarking process. Canonical linear pathways were extracted only from MetaBase and annotated according to four main concepts: directionality, effect, mechanism, and weight. Directionality indicates the intended flow of signal, from the source to the target node. The effect (or edge type) denotes whether the interaction is inhibition (-1), unknown (0), or activation (1). Mechanism distinguishes generic molecular interactions (interactions from receptors upstream to the transcription factors downstream, coded as 0) from transcriptional regulation edges (TFs with their target genes, coded as 1). Finally, weight determines interaction confidence based, among the others, on literature support. Regardless of the database source used, the mandatory annotation for each interaction is directionality information, whereas any other information will not be used by algorithms.

OmniPath PKN was constructed by integrating signaling and TF-target interactions using OmniPathR (v. 3.14.0) R package. Signaling interactions were extracted using the `import_omnipath_interactions()` function and with assigned `mechanism = 0`, whereas transcriptional regulatory interactions imported using `import_transcriptional_interactions()` were annotated as `mechanism = 1`. These were combined into a single network, and interactions with `effect = 0` were kept only if `mechanism = 0`. Nodes in OmniPath are proteins or protein complexes (UniProt IDs), with the corresponding gene symbol(s).

Table 3.2: Summary table of OmniPath and MetaBase prior knowledge resources. Number of nodes and edges are displayed for each resource. Network refers to the full interaction network, while regulons and linear path regulons are downstream derived, regulators subsets. The Regulator and Target rows correspond to the number of source nodes and target nodes, respectively. Gene space, counts how many of those nodes correspond specifically to genes (not proteins or others). The three edge type rows indicate the number of Activation, Inhibition, or Transcriptional regulation interactions, with the Total number of interactions for each resource also provided.

Network components		OmniPath		MetaBase		
		Network	Regulons	Network	Regulons	Linear Path Regulons
Nodes	Regulator	6,166	4,442	33,927	11,739	2,922
	Target	6,723	6,723	15,229	10,476	9,465
	Gene space	7,809	5,622	17,693	9,988	3,185
Edges	Activation	119,113	5,842,390	81,866	23,844,526	3,493,007
	Inhibition	13,680	4,270,032	61,214	21,469,352	1,361,149
	Transcriptional regulation	64,367	10,112,422	101,752	45,313,878	4,854,156
	Total	145,896	10,112,422	657,746	45,313,878	4,854,156

For using MetaBase as another source of PKN, the global network was extracted using the `networkFromMetabase()` function, via the `metabaser` (v. 5.1.0) and `CBDD` (version 20.0.3) R packages. Unlike OmniPath, it already includes both signaling interactions (`mechanism = 0`) and transcription regulation interactions (`mechanism = 1`). Only high-confidence interactions with defined effect (activation or inhibition) were kept. Originally, the network contained specific MetaBase network objects, that were processed to add only the corresponding gene symbols to the network (using the function `metabaser::annotate.nwobj2gene()`).

Another way of representing interactions that can be used as reference data for both topology-based and enrichment-based tools are the regulons. The regulons were extracted using the causal reasoning algorithm, through the `CBDD::hypothesisGeneration()` function, providing as parameters the downstream depth for the search (in this case, 4 steps) and the position in the pathway where transcription regulation links may appear (set to anywhere in the path). Bearing in mind that both networks used as input contain the directionality of the signal, this function will then predict which targets are influenced (by activation or inhibition) by each specific regulator. Finally, all possible interactions were filtered to retain only those where a node and all downstream activated or repressed genes are present. The final number of nodes and interactions of the regulons is also detailed in Table 3.2. In addition to the network and the regulons, canonical linear pathways, available from MetaBase were also included. Canonical linear pathways are sequences of biological entities and interactions between them. They are automatically generated from pathway maps and represent highly curated canonical signaling paths, starting from an important signaling molecule and ending, usually with a transcription regulation interaction.

### 3.3 Algorithms: implementation and wrapper function's architecture

To carry out a systematic and robust comparative evaluation of inference algorithms, wrapper functions were developed to build a common framework and to standardize the input data and output, to ensure compatibility between each algorithm data requirements and processing methods. A wrapper is a function that serves as an intermediary layer. These are required to handle data type conversions, parameter standardization, and result formatting, allowing diverse algorithms to be executed consistently regardless of their underlying implementation differences. Here there are two types of wrappers: shared and individual ones. The shared wrapper architecture incorporates an already established package that bundles several algorithms inside, unlike the individual ones that incorporate single algorithms. The connectivity mapping from the RSCM package, enrichment methods from decoupleR, and topology-based algorithms from CBDD were implemented in shared wrappers. On the other hand, causal reasoning CARNIVAL, CausalR, ProTINA, CIE, and NicheNet were incorporated in individual wrappers. Table 3.3 provides a complete list of algorithms with their annotations.

Some supporting helper functions were also implemented to facilitate essential data conversions across all wrappers. Those functions include mapping identifiers between transcriptomic datasets and network nodes, to ensure matching IDs and converting the input data to the required format, if necessary. For the query input data, the tool may need a full signature or DEGs. When DEGs are required, the full signature can be filtered using a fold change and p-value threshold, or by simply taking the top threshold for DEGs by fold change magnitude. By default, signatures are converted to DEGs using the top 500 genes with the strongest changes, taking half from the most upregulated and half from the most downregulated genes. The default parameters were also used to run each algorithm, given the complexity of the study, evaluate other parameters was out of the scope. However, the wrapper function is prepared to accept custom parameters for each algorithm, if needed, including for converting to DEGs using fold change and p-value thresholds, instead of using the top genes. As reference datasets, the workflow can start with PKN or full signatures. Topology, enrichment, and CMap algorithms require, respectively, networks, regulons/gene sets, and full signatures. To use this large variety of input data and tools, some conversions between these data are required. All the conversions are represented by the arrows in Figure 3.1 B.

As with input data, output data must also have a similar shared format so that the performance of each algorithm can be evaluated systematically. For that reason, at the end of each run, all algorithm wrappers return a table with all prioritized regulators identified without any significance filtering applied. The output contains a score column, with a larger score reflecting greater confidence in the causal regulator for the observed differential expression patterns. Score may also be signed if the tools can predict the directionality of the perturbation. In that case, regulators are ranked by absolute value of

Table 3.3: Summary of computational methods evaluated in the benchmarking study. A total of 27 algorithms categorized into the following methodological approaches: 1) Enrichment, 2) Connectivity Mapping, 3) Topology with Causal Reasoning, and 4) Topology with Node Prioritization. For each algorithm, the corresponding R package implementation and version used are reported. The Reference and Query columns indicate the required input data types to run the algorithm. The Output column specifies whether algorithms produce node rankings (prioritized lists of potential regulators) or subnetworks (connected molecular entities involved in the MoA).

Method	Algorithm(s)	R Package	Reference	Query	Output	Ref.
Enrichment	fgsea	decoupleR (v. 2.12.0)	Regulons	Full signatures	Node ranking	[12]
	viper					
	ulm					
	mlm					
	udt					
	mdt					
	wsum					
Connectivity Mapping	wmean	RCSM (v. 0.3.0)	Full Signatures	DEGs	Node Ranking	[86]
	KS					
	explainable cosine (xCos)					
	XSum					
	GSEAweight0					
Topology (Causal Reasoning)	GSEAweight1	CARNIVAL (v. 2.16.0) CausalR (v. 1.38.0) Protina (v. 0.1.0) CIE (v. 1.0.0) nichenetr (v. 2.2.0) CBDD (v. 21.0.0)	Network	DEGs	Subnetwork	[96]
	ZhangScore					
	CARNIVAL					
	CausalR					
	ProTINA					
	CIE					
	NicheNet					
Topology (Node Prioritization)	causalReasoning	CBDD (v. 21.0.0)	Network	DEGs	Node ranking	[91] [99] [92] [100] [101] [102] [103] [104]
	SigNet					
	quaternaryProd					
	networkPropagation					
	randomWalk					
	overconnectivity					
	hiddenNodes					
	interconnectivity					

score, and activation/repression status is stored in separate column effect (coded as -1/1 respectively).

### 3.3.1 Connectivity Mapping

Figure 3.2 represents the wrapper function framework for running CMap algorithms from the RCSM R package [86]. This package provides uniform implementations of several CMap scoring methods including KS, and GSEA-based approaches. The function is designed to receive filtered DEGs lists as query input and full perturbation signatures as reference data. If full signatures are used as the query, they are internally converted

to DEGs using the filtering parameters previously described. These algorithms are designed to quantify the similarity between query and reference perturbation signatures. Some of the similarity metrics used include the KS statistic, implemented in the original CMap [11]; xCos, a cosine similarity metric between query and reference fold-changes; Xsum connectivity map statistic based on the sum of reference fold-change values of query genes; GSEAweight0 is a GSEA weighted KS ES with parameter  $p = 0$ , which ignore the fold-change magnitude for computation; GSEAweight1 with parameter  $p = 1$ , where fold-change magnitude contributes linearly to the final score; and Zhang, a CMap score first suggested in S. D. Zhang and T. W. Gant (2008) [105]. The wrapper function handles the different algorithm requirements by preparing either separate up- and down-regulated gene lists or simple gene vectors for xCos, also including optional regulator filtering for TF mode. The output is formatted to return regulator rankings with similarity scores, directional effects, and optional statistical significance measures. The results are sorted by absolute score magnitude to prioritize the most relevant regulatory relationships regardless of similarity direction. For these algorithms, the regulator score measures the similarity of the query versus the reference perturbation signature.

### 3.3.2 Pathway Enrichment

For running the enrichment-based algorithms, the decoupleR package [12] was used. The package was initially used to benchmark approaches for TF activity inference. It contains 12 algorithms already implemented to expect prior knowledge resources (gene sets or regulons) to derive biological processes from omics data. Some of them take directionality into account (i.e., can work with regulon-gene set with activated and repressed genes). Of all the algorithms already implemented in decoupleR, only GSEA and the others that respect directionality were considered for this benchmarking effort. As for CMap algorithms, a shared wrapper function (Figure 3.3) was built to prepare the input and output data, designed to accept full signatures as query and regulon table or a gene regulatory network as reference data. If the reference is a list of signatures or DEGs, it is converted to directed regulatory networks using the common filtering parameters described above. The implementation supports TF-mode by filtering the network to keep only transcription-regulation edges. The query signatures are converted to a fold-change matrix and ID space conversions can also be performed, if required to match network node identifiers. For algorithms that support directional information, the wrapper uses the edge type from the network. The output of the wrapper function consists of regulator rankings with scores, effects, and p-values (if returned by the algorithm) organized by signature name and sorted by absolute score magnitude.

### 3.3.3 Topology-based methods

CARNIVAL [96] is prepared to integrate several types of prior knowledge (signed and directed PPI, TF-targets, and pathway signatures) to yield a causal subnetwork



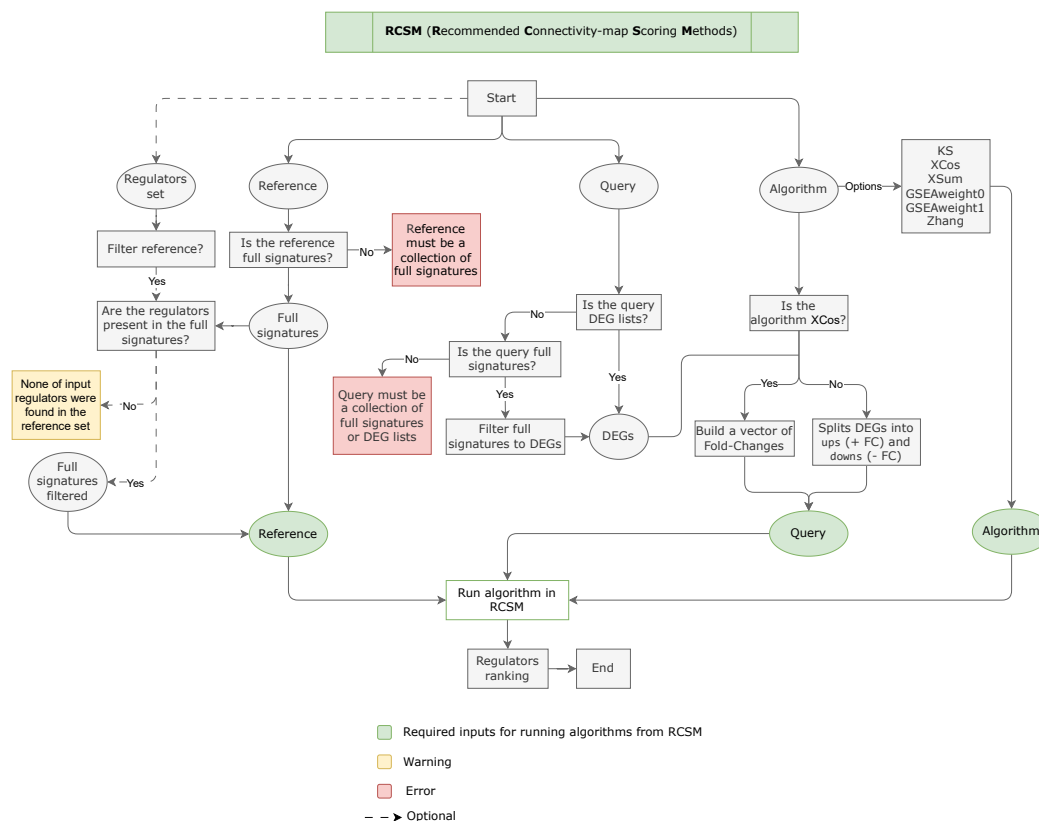


Figure 3.2: Flowchart representing the main steps for implementing connectivity mapping algorithms pre-built in the RCSM package. The general computational pipeline for executing connectivity-based methods, showing the main input requirements, data pre-processing steps, algorithm execution, and output generation. Green indicates required inputs, while red highlights potential errors.

explaining the MoA behind the observed omics data. This algorithm expects query DEGs as input and a network as reference. CARNIVAL wrapper (Figure 3.4) begins by validating essential inputs, including the CBC solver path required for the optimization engine. It determines the execution mode based on whether perturbation targets and pathway weights are provided, enabling either the standard CARNIVAL or the inverse CARNIVAL algorithms. For reference network preparation, the wrapper handles multiple input formats by converting, if needed, signature collections or DEGs lists into networks (with source node - interaction - target node). The full signatures are also converted into matrix format. CARNIVAL performs TF activity inference with either a network-derived approach, if effect and mechanism are present, or using DoRothEA regulons from decoupleR. As an option, depending on the execution mode selected, a pathway activity score can be also calculated using PROGENy from decoupleR. The network is filtered to contain a subset of the interactions, keeping only nodes reachable from relevant TF. When CARNIVAL algorithm is executed, the results are provided with regulator scores reflecting the frequency where each node appears in the causal subnetwork.

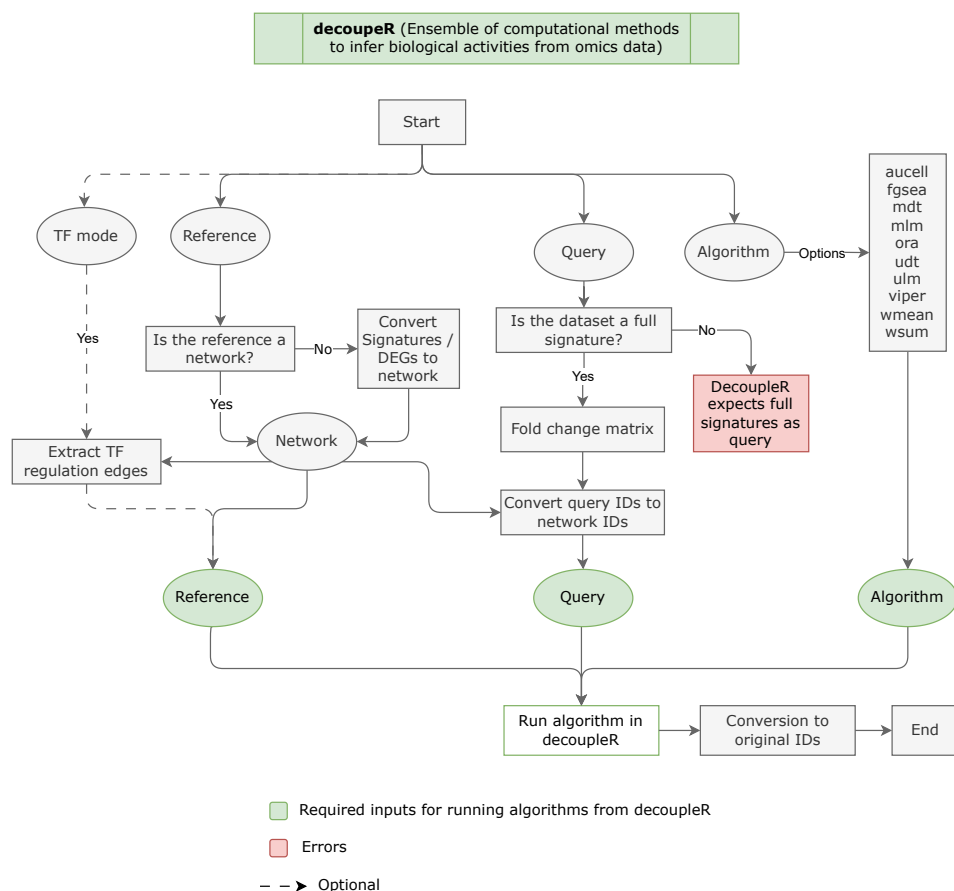


Figure 3.3: Flowchart representing the main steps for implementing enrichment algorithms pre-built in the decoupleR package. The general computational pipeline for executing enrichment-based methods, showing the main input requirements, data preprocessing. Green indicates required inputs, while red highlights potential errors.

ProTINA [97] algorithm generates protein activity scores based on gene expression changes through network perturbation analysis. The wrapper (Figure 3.5) starts by checking if the query is a collection of full signatures and if the reference data is a network. If the reference is a signature or a DEGs list, it is then converted to a network. Then, query data is filtered to remove entire signatures without any DEGs, as ProTINA cannot process full profiles with only zero fold-change expression. After that, the full signatures are converted to a matrix of fold changes, to remove genes that have zero standard deviation across all signatures. ProTINA was originally designed to handle experiments with multiple replicates, time points, or dosages for each perturbation. Since none of the data used for benchmarking have measurements in different conditions, the function only creates a vector with consecutive integers starting from one to the number of signatures, representing the number of groups. Reference data is processed to construct a Protein-Gene Regulatory Network (PGN), which combines PPI from the input network with TF-target gene relationships. From the PGN it is calculated an adjacency matrix

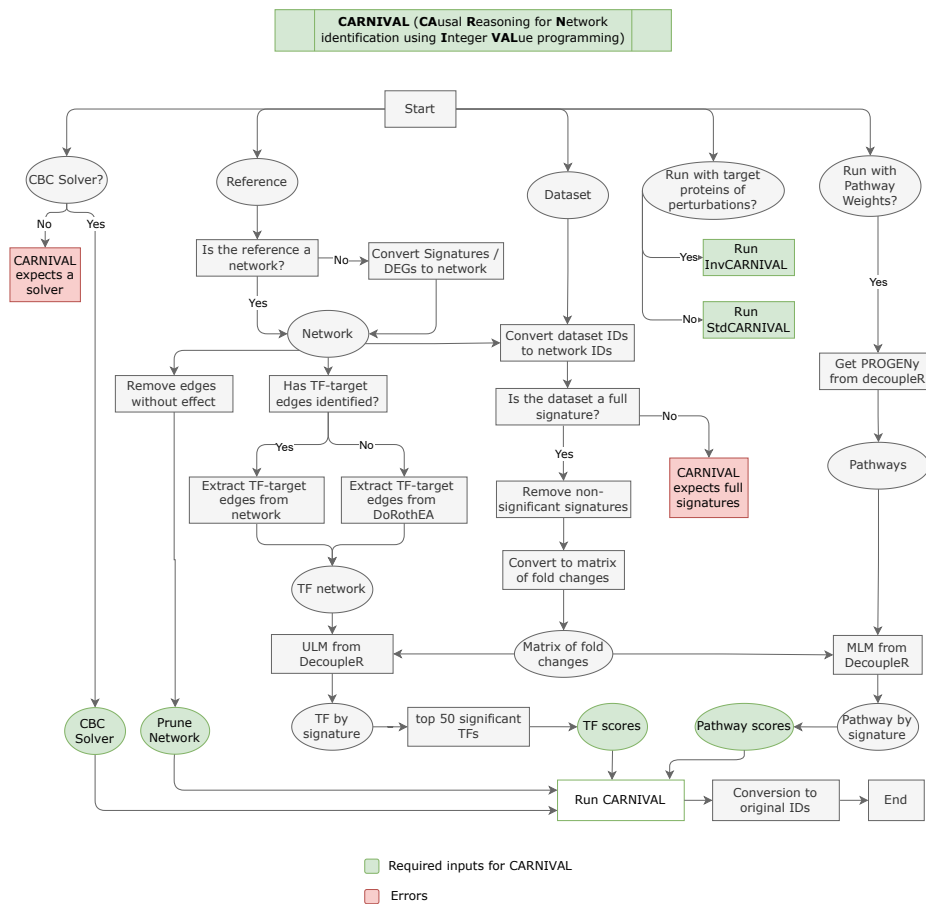


Figure 3.4: Flowchart representing the main steps for implementing the CARNIVAL algorithm. The general computational pipeline for executing this topology-based method, showing the main input requirements, data preprocessing steps, algorithm execution, and output generation. Green indicates required inputs, while red highlights potential errors.

where rows represent proteins, columns represent genes, and matrix elements indicate the regulatory relationships. Finally, the algorithm returns a matrix of protein (regulator) activity scores (signed Z-scores) for each perturbation group. The wrapper converts the results back to gene symbols and ranks regulators by their activity scores, providing both magnitude and directionality of predicted protein activities.

To identify DEGs, the CausalR [10] wrapper (Figure 3.6) starts by processing the query signatures using the common filtering parameters described above. The resulting data is then converted to the required CausalR format where genes are assigned with regulation status values (1 for upregulation, -1 for downregulation, 0 for unchanged). Non-DEGs within the signature are marked as unchanged, rather than excluded. The final format for the query dataset is a matrix with two columns for each experiment, reporting gene name and regulation status (1, -1, or 0) instead of fold-change values. As for ProTINA, signatures that do not have any DEGs are removed. If reference data are not already provided in the correct format, the function can convert signature collections or DEGs

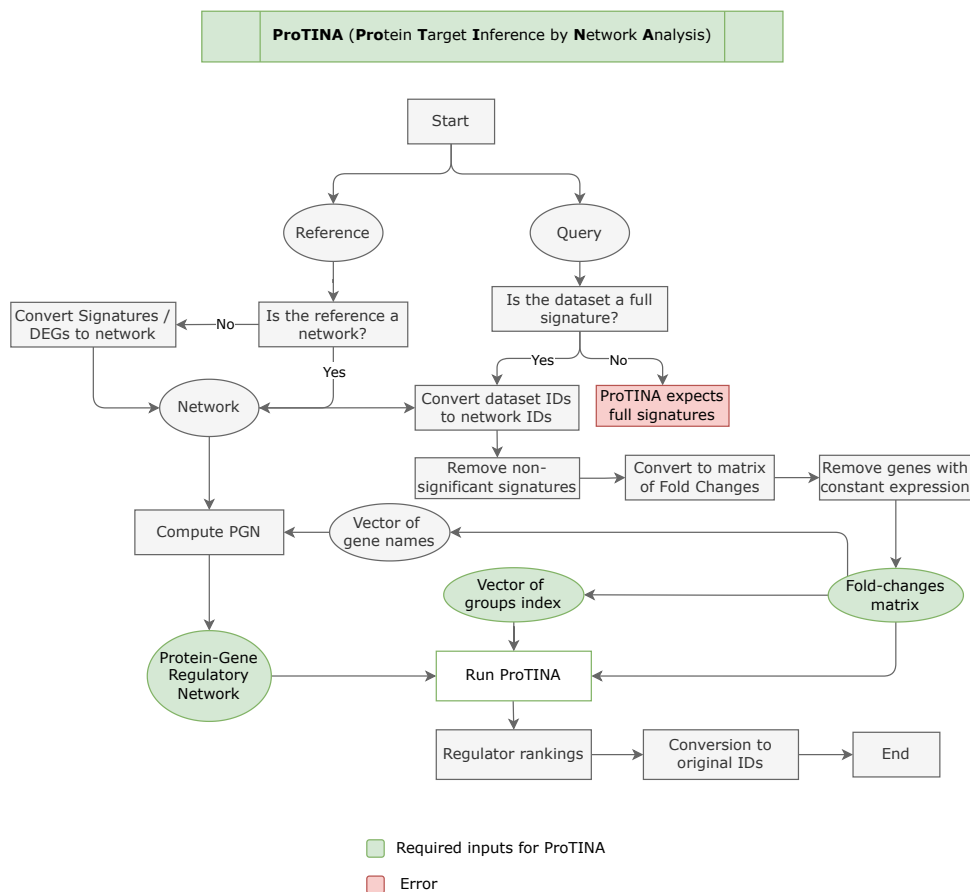


Figure 3.5: Flowchart representing the main steps of ProTINA algorithm implementation. The general computational pipeline for executing this topology-based method, showing the main input requirements, data preprocessing steps, algorithm execution, and output generation. Green indicates required inputs, while red highlights potential errors.

lists into directed causal networks. Edges connect the regulators (signature names) to their targets (DEGs) with edge effects corresponding to DEGs fold-change signs. After having a network as reference, CausalR requires the construction of a Computational Causal Graph (CCG), which contains twice the number of nodes and edges as each regulator is reported both as up and down/regulated. CausalR allows two options of algorithms, RankTheHypotheses and runSCANR. RankTheHypotheses algorithm uses the configurable path-length parameter delta, to control how many network edges can be traversed from regulator hypotheses to the observed gene expression data, enabling from a direct transcriptional regulation ( $\delta = 1$ ) analysis to a multi-step causal cascade ( $\delta > 1$ ). Finally, the algorithm returns, per each signature, a regulator ranking with the regulator name, and the corresponding score (difference between correctly and incorrectly predicted DEGs), p-value, and the predicted regulatory effect.

The CIE [89] wrapper (Figure 3.7) starts by converting (if needed) full signatures to DEGs, then transforming reference data into a network format. A key aspect of this

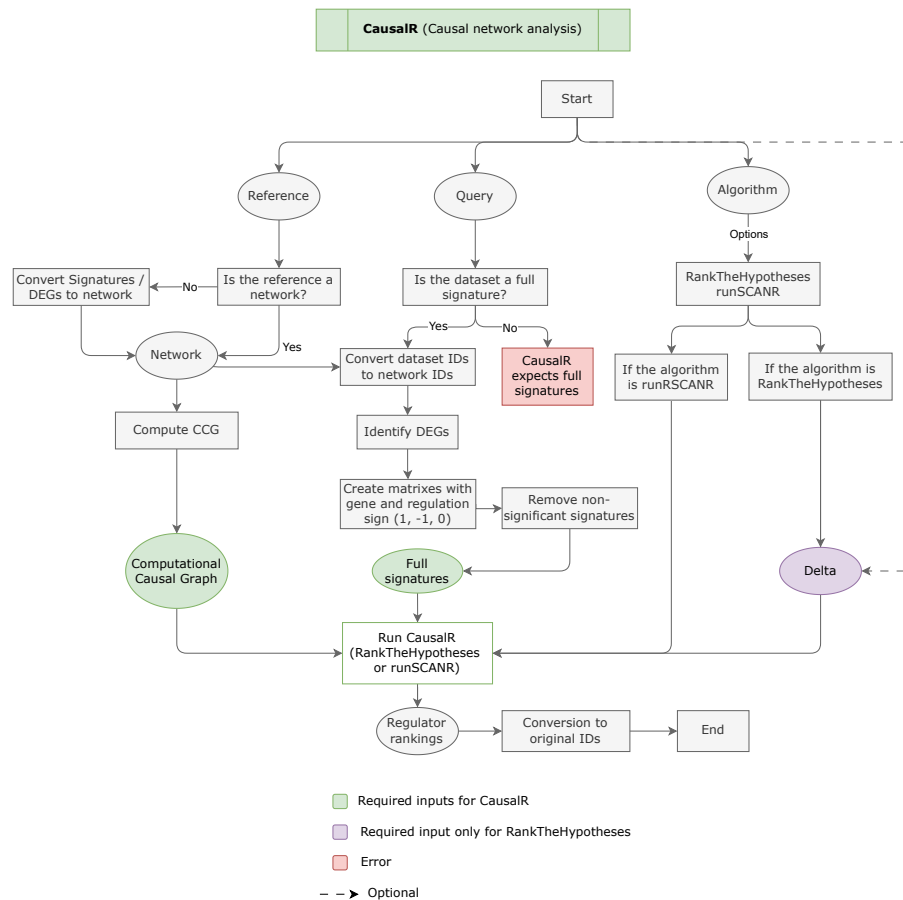


Figure 3.6: Flowchart representing the main steps of CausalR algorithm implementation. The general computational pipeline for executing this topology-based method, showing the main input requirements, data preprocessing steps, algorithm execution, and output generation. Green indicates required inputs, while red highlights potential errors.

implementation involves preparing the network data structure to create two essential inputs: network entities (nodes) and network interactions (edges). The entities data frame contains unique gene IDs classified as mRNA or Protein based on the targets of the transcriptional regulation edges (mechanism = 1). The relations data frame links source and target nodes using their distinct identifiers and maps the network edges with their regulatory effects. The function supports three CIE algorithms: Fisher for unsigned networks, Ternary for completely signed networks, and Quaternary for partially signed networks (default). The results from CIE are the regulators and the corresponding causal reasoning scores, regulatory effects, and a p-value ranking.

As the other wrappers, the one implemented for NicheNet [98] (Figure 3.8) also ensures that the input formats are correct by converting the input signatures to DEGs and reference data to networks. From the reference network, three subnetworks are created based on the regulatory mechanism: one ligand-receptor network and one signaling network from the signaling interactions (mechanism = 0), and one transcriptional

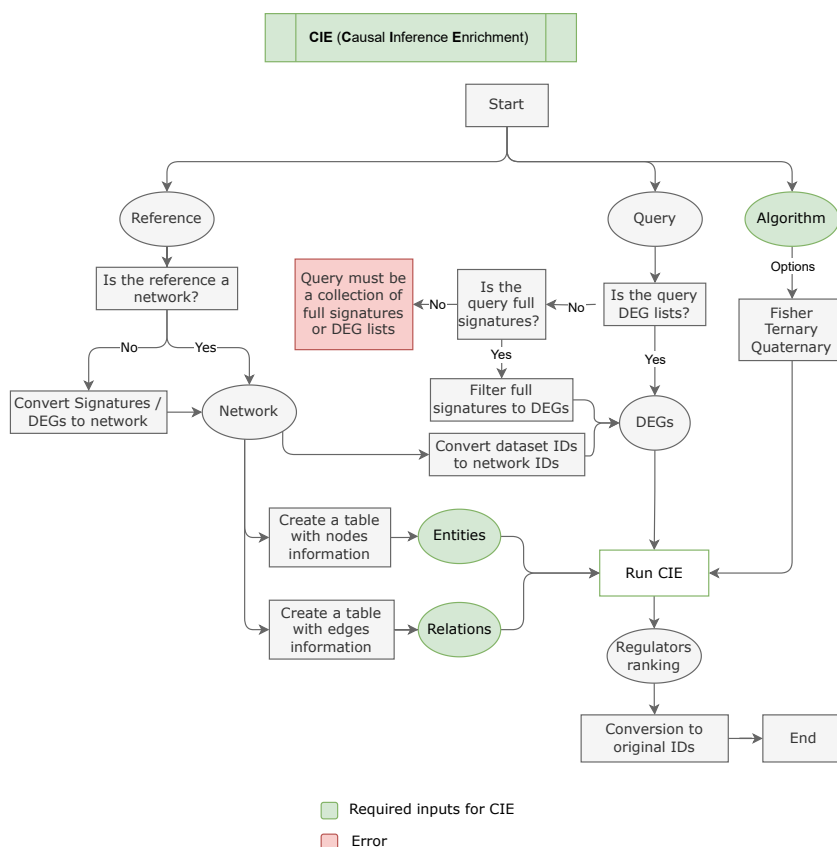


Figure 3.7: Flowchart representing the main steps of CIE algorithm implementation. The general computational pipeline for executing this topology-based method, showing the main input requirements, data preprocessing steps, algorithm execution, and output generation. Green indicates required inputs, while red highlights potential errors.

network from gene regulatory interactions (mechanism = 1). The ligand-receptor and signaling networks are constructed using an optional vector of regulators, which provides a list of source nodes to filter signaling interactions. If this optional vector is not provided, only unique source nodes are used to create the set of regulators. The ligand-receptor network includes only edges where the source node is a regulator, whereas the remaining edges are presented in the signaling network, with the downstream signaling cascades down to the target genes. In the next step, a weighted network is constructed (`construct_weighted_networks()` function from `nichenetr` R package [98]). This function merges the three subnetworks into one graph, and assigns weights to edges based on their origin. By default, equal weights are assigned to all edge types (ligand-receptor, signaling, and gene regulatory). To down-weight the influence of highly connected nodes, hub correction is applied to the network (`apply_hub_corrections()`). With the weighted network, a ligand-target matrix is built (`construct_ligand_target_matrix()`) containing target genes as rows, ligands as columns and scores for each entry (inferred signaling

strength), from each ligand to each gene. At this point, the inputs required to run NicheNet core function (`predict_ligand_activities()`) are ready: a vector of DEGs, the full set of genes present in the ligand-target matrix (used as the background), the ligand-target matrix of regulatory scores and the list of ligands (regulators) to prioritize. The output from NicheNet is ranked by corrected Area Under the Precision-Recall Curve (AUPR) scores and identifiers are converted back if needed.

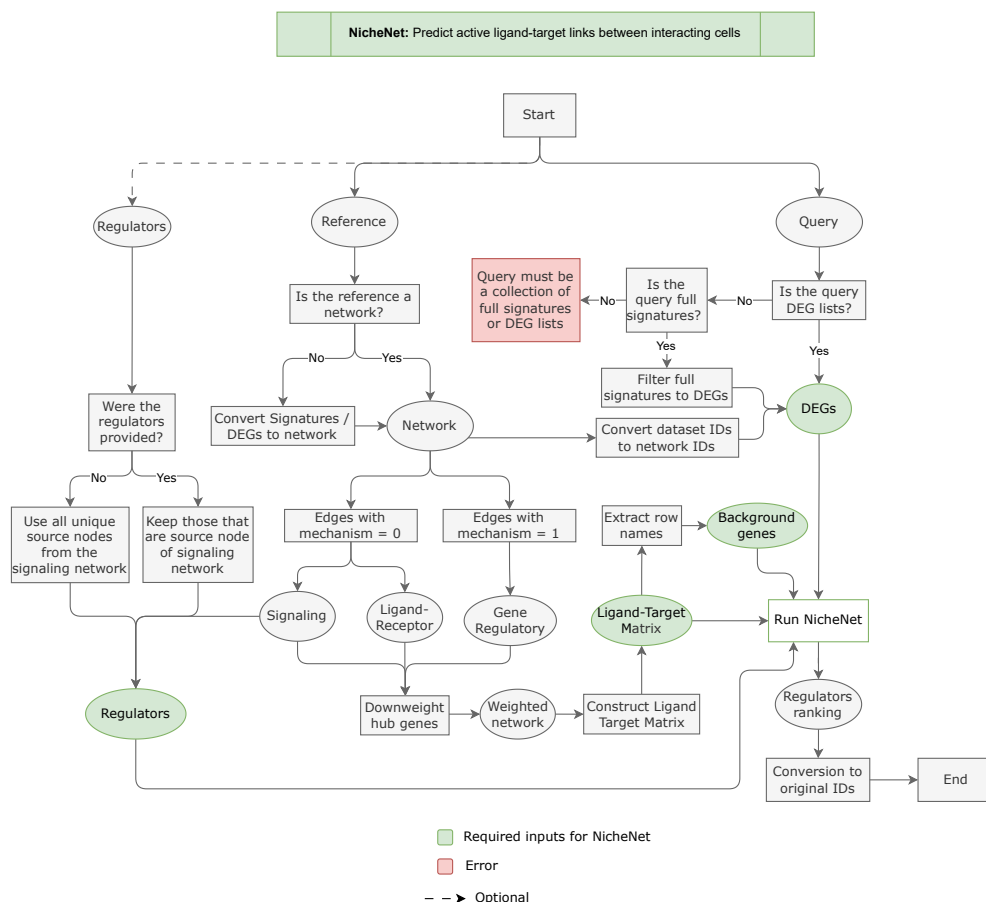


Figure 3.8: Flowchart representing the main steps of NicheNet algorithm implementation. The general computational pipeline for executing this topology-based method, showing the main input requirements, data preprocessing steps, algorithm execution, and output generation. Green indicates required inputs, while red highlights potential errors.

The implementation of the CBDD wrapper functions (CBDD baseline and CBDD causal reasoning) follows an identical structure, with both functions sharing the same core preprocessing pipeline and output handling. For this reason, they are both represented in a single schema, with the differences highlighted (Figure 3.9). Both start with validation and preprocessing of input data, converting query datasets to DEGs lists, handling reference networks as described previously and performing ID conversion to map query genes/regulators to network-specific identifiers. The wrapper functions then execute their

respective algorithms, and return the results. Despite the similarities, there are differences between the implementations that reflect their distinct computational requirements. The CBDD baseline function supports randomWalk, overconnectivity, interconnectivity, networkPropagation, and hiddenNodes algorithms, while CBDD causal reasoning runs causalReasoning, SigNet, and quaternaryProd. The causal reasoning wrapper requires network attributes such as the effect and mechanism, reflecting the need for directed, signed edges for causal inference. Moreover, the CBDD baseline converts networks to igraph objects, while CBDD causal reasoning creates a causal graph. Essentially, both wrappers return regulator scores, but since CBDD causal reasoning considers directionality, the results contain additional information about the effect of the relationships (predicted activation or inhibition).

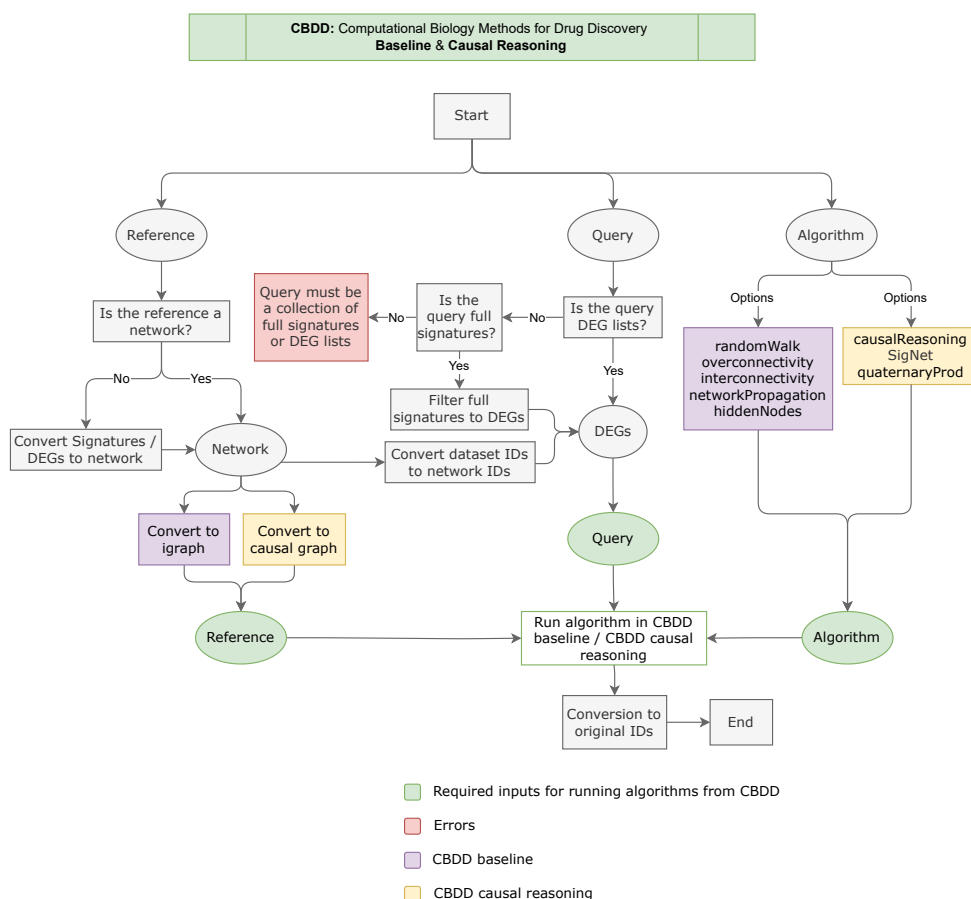


Figure 3.9: Flowchart representing the main steps of the wrapper function of both the CBDD baseline and CBDD causal reasoning algorithms implementation. General computation pipeline and color coding as previously.



### 3.4 Evaluation

The final part of the study pipeline involves evaluating the algorithms by comparing the outputs from the algorithms' runs with the gold standard data. This evaluation process involves three main steps: configuration, running the algorithm, and evaluating performance against the standards. The configuration step includes selecting the input data (both the query and reference), as well as the standard with the verified list of the targets. Then, the execution requires defining which algorithms to run for the selected input data. Finally, the metrics are defined to assess performance of each algorithm, considering the known targets provided by the gold standard data.

The evaluation was performed at the level of each signature, aggregating results to obtain an average assessment by dataset. In general, the algorithms return all the ranked and scored regulators, so the threshold for considering which regulators are significant was also defined in two ways. If an algorithm provides p-values for the regulators, the threshold of 0.05 is used to consider that regulator as significant (if p-value < 0.05). If this is not the case, the top 1% of regulators with the highest scores are selected. In the case of CARNIVAL, which returns a subnetwork, all the regulators in that network are already considered significant.

The performance was measured at a target and pathway level. The former considers target recovery as the percentage of gold standard targets among all significant regulators, top 1% and top 100 regulators. Target enrichment was calculated using hypergeometric test enrichment ( $-\log_{10}(\text{p-value})$ ), for the intersection of significant regulators and the gold standard targets. It was also analyzed the rank of the top predicted gold standard target and target effect correctness, when directional information was available. The target recall was measured at multiple percentage cutoffs (0.001% to 1% in 0.001% steps) to enable the calculation of the AUPR. This led to three general classifications:

1. Overall recovery: average of target recovery percentages, enrichment scores, and scaled AUPR (for the top 1%) across all signatures.
2. Recoverable signatures only: target recovery percentage calculated only on the subset of signatures whose true targets appear in the reference network.
3. Win rate: to assess each algorithm's ability to identify actual target genes, the win rate represents how frequently a given output ranking contains a true target higher than all other competing algorithms. The win percentage above expected (WPAE) is obtained subtracting the expected win rate under random chance (1 divided by the number of algorithms), as the number of competing algorithms can change between runs.

The assessment of performance at the pathway level is equally important to understand whether the algorithm not only recovers the target itself, but also the biological pathways to which these targets are connected. From the selected significant regulators, an enrichment

analysis of the biological pathways was done using the CBDD `enr` function, with MetaBase mapping gene sets to the pathways that are most associated with the targets predicted by the algorithms. The statistical significance of the overlap between biological pathways identified from the targets and those manually annotated via MetaBase was calculated with a hypergeometric test (`hh_pvalue()` function from CBDD). The result was recorded as a pathway enrichment score ( $-\log_{10}(\text{p-value})$ ), where a higher score indicates that the algorithm is able to identify pathways directly affected by the perturbation. For comparing the results, the mean pathway enrichment score for each algorithm was computed by averaging those signatures whose true targets appear in at least one MetaBase pathway.

## RESULTS

*In the following section, the results of the benchmarking study are presented across three main levels: the algorithms, the reference and the query dataset. For each level, the metrics to assess the performance are analyzed to capture both computational feasibility and biological relevance. A final ranking combines all the scores, aiming to identify the best algorithm providing real-world applicability at each dimension of data.*

A full evaluation was conducted with 86 runs (Table I.1), resulting from the combination of 10 query datasets and 9 reference datasets (excluding self-comparisons). Given the extensive number of individual runs, the results will be presented as performance summaries of the four key dimensions: algorithms, datasets, reference, and target identification performance. Table I.1 presents the runs that were evaluated, detailing which combination of query and reference datasets was used, and some characteristics of each query dataset.

### 4.1 Algorithms

To evaluate the performance of the 27 algorithms, several metrics were used to describe the output generated from each package. One is the practical scalability of each tool, and the other is the biological accuracy of the results. To address the first aspect, the runtime and the reliability were measured. To evaluate the biological relevance of the predictions, the AUC, the win rate, and the pathway enrichment scores were measured.

Execution time and algorithm failures can prevent the use of algorithms in research settings. For this reason, evaluating those metrics is important for understanding the feasibility of each computational approach. The runtime represents the time in seconds that each tool took to execute, and it is directly correlated with query and reference datasets size. In this context, a mean runtime captures how long on average each tool takes to generate the desired output. Whilst this value can be informative, a better approach to present these data is to consider mean runtime per signature (Figure 4.1) or a value normalized considering the size of the reference dataset (the final runtime score in Figure 4.5). The final runtime score was calculated as follows. First, within each

reference, each tool's per-signature times were rescaled from 0 (slower) to 1 (faster). Then, those values were averaged across all references and inverted, so that faster algorithms get higher scores. As shown in Figure 4.5, the majority of topology and enrichment methods topped the rankings with scores greater than 0.90, except for the enrichment methods *udt* and *viper*, which had lower scores (0.63 and 0.39, respectively). The topology methods, *causal reasoning*, *interconnectivity*, and *network propagation*, along with the enrichment methods *ulm*, *wmean*, and *wsum*, topped the rank with scores near 1, providing final output on average in 1-2s per signature. Connectivity-mapping tools stayed in the middle of the ranking with scores ranging from 0.80 to 0.90, and mean runtimes ~13~24s per signature. In contrast, methods like *NicheNet* and *ProTINA* (respectively 416s and 126s per signature) were markedly slower, thus obtaining a score close to 0.

To complement the runtime evaluation, a reliability score was also computed to assess how often each algorithm failed. The final reliability score is an average of a success rate (runs that returned results) and a run rate (runs launched) from a scale of 0 to 1, where scores close to 1 indicate more reliable algorithms. Details of the rates can be found in Table I.2, and the final reliability scores are represented in Figure 4.5. At the top of the rank, some of the topology and enrichment methods, with *causalReasoning*, *randomWalk*, and *ulm* topped scoring close to 0.99 (i.e. returning results in 98% of the launched runs). Next in the ranking it can be found *overconnectivity* (score of 0.98) and *wmean/wsum* (scores of 0.91). Consistent with the runtime, the connectivity methods are in the middle with scores between 0.80 and 0.82, with *XCos* further down with 0.75. At the bottom of the ranking, *NicheNet*, which although it achieved total success when launched, only ran in 2 out of 86 runs (run rate of 0.02), resulting in a reliability score of 0.51. The lowest-performing algorithms in this category were *CARNIVAL* (score of 0.06) and *CIE* (score of 0.05), which were skipped or failed in every single run.

The reliability score reflects the capability of providing the desired output. Thus, only the tools that were launched on almost every run and that did not fail too often, obtained scores close to 1. From a practical point of view, beyond the runtime, reliability is also a critical criterion. Based on these two metrics, *CARNIVAL*, *CIE*, *CausalR*, *hiddenNodes*, and *quaternaryProd* did not complete any benchmarking runs, reflecting their impracticality at this scale. For this reason, these five algorithms were removed from the final evaluation, since they did not behave as expected for a fair comparison.

In the research field, only a small number of predictions can be experimentally validated. Hence, the choice of the algorithms to identify potential target genes should prioritize those that are most likely to include the true target. To evaluate that, the metric AUC was used to assess performance in recovering the true target among the top 1% of the predicted regulators. The AUC scaled mean was first calculated by taking the AUC mean across all query signatures. Then, for each reference dataset, these means were scaled from 0 to 1 and averaged across all runs for each algorithm. The final AUC scores with values close to 1 mean better performance. The mean scaled AUC values obtained for each algorithm can be found in Figure 4.2, and the final AUC scores in Figure 4.5.

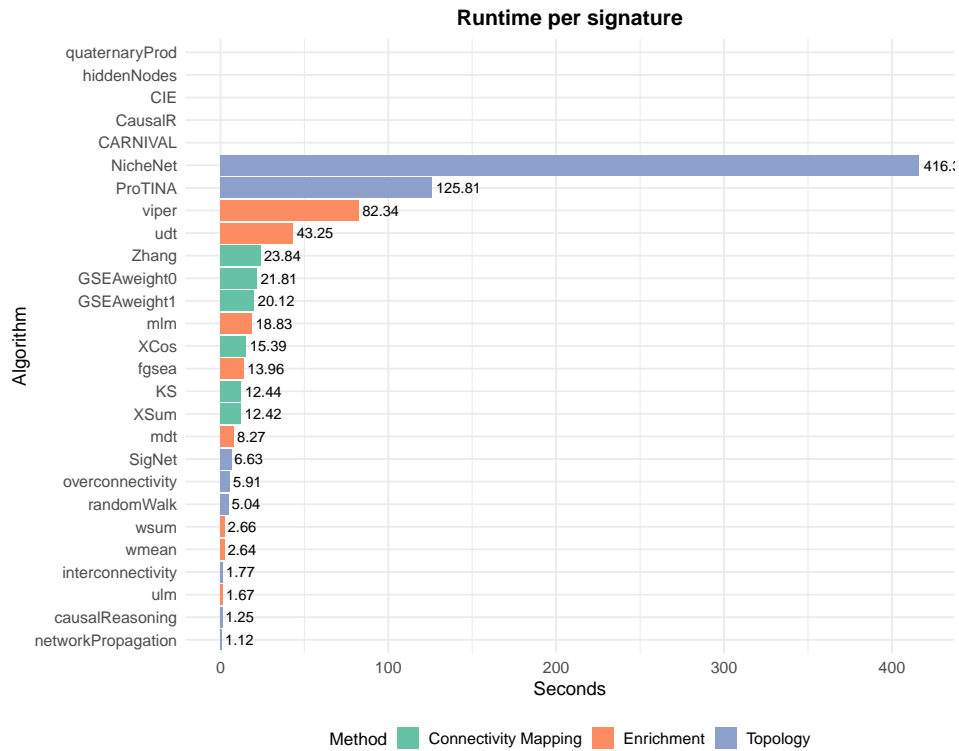


Figure 4.1: Average runtime per signature in seconds for each algorithm, with values displayed on bars. Algorithms are colored by category.

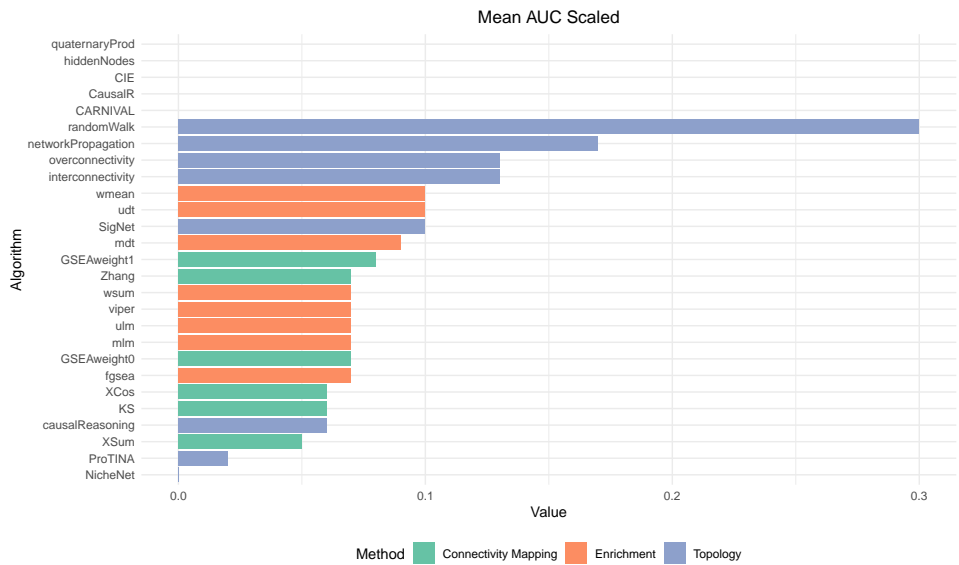


Figure 4.2: Mean scaled AUC values for each algorithm, ordered by performance. Algorithms are colored by category.

The topology-based algorithms with node prioritization occupied the top of the rank, with randomWalk showing the highest mean scaled AUC of 0.30 (score of 1), achieving 30% recovery within the top 1% of predictions. The others (networkPropagation, interconnectivity and overconnectivity), although at the top, only achieved half the performance of the randomWalk (scaled AUC of 0.17, 0.13 and 0.13, respectively). The enrichment-based methods displayed a uniform performance, with scaled AUC values clustering between 0.07 and 0.10. Connectivity mapping approaches do not seem suitable for this task, with all five methods achieving scaled AUC values below 0.08. The XSum method performed worst overall (scaled AUC: 0.05, score: 0.16), while GSEAweight1 represented the best among connectivity tools but still achieved only 28% of randomWalk's performance. These results suggest that methods designed for perturbation signature matching lack the mechanistic reasoning necessary for an accurate target identification.

The win rate and the WPAE were two metrics used to evaluate how often an algorithm ranks the true targets and the performance over a random selection of target genes. A win rate of 0.5 means that out of 25 different query datasets, the true target is ranked among the top regulators in 25 of them. The WPAE is a normalization, to account for a different number of competing tools in different runs.

A positive WPAE means that the performance was better than chance, and a negative indicates an underperformance of the tool compared to a random selection of target genes. The results of these metrics are shown in Figure 4.3. The randomWalk algorithm stood out compared with other algorithms with a win rate of 0.29 and a WPAE of 0.14. The win rate of 0.29 indicates that randomWalk ranked the true target highest in only 29% of the runs. A WPAE of 0.14 indicates that the algorithm performed 14% better than random chance in ranking the true target highest. Although it seems a modest result, this is a significant achievement given the complexity of the task and the number of competing algorithms. Only two other topology-based methods achieved positive WPAE values, overconnectivity (win rate: 0.17, WPAE: 0.02) and interconnectivity (win rate: 0.14%, WPAE: 0.02). Most algorithms (19 out of 22) showed negative WPAE values, indicating they win less frequently than random chance would predict, revealing that they rarely identify the true target. networkPropagation, despite showing a good AUC performance, achieved a slightly negative WPAE (-0.05), and won only 8% of comparisons. This discrepancy suggests the algorithm excels at enriching true targets within its top predictions, but rarely ranks them at the absolute top position. The win rate score (Figure 4.5) close to 1 for randomWalk (calculated as mean WPAE rescaled between 0 and 1), reinforce the superiority of this approach against the other comparators considered.

The mean pathway enrichment (Figure 4.4) assesses whether the regulators identified by each algorithm are biologically meaningful. This metric is not evaluating direct target recovery, but rather the ability to identify regulators that are functionally related to the target. Hence, an algorithm with a high pathway enrichment performance should identify a high number of regulators that are enriched in the same pathways as the true

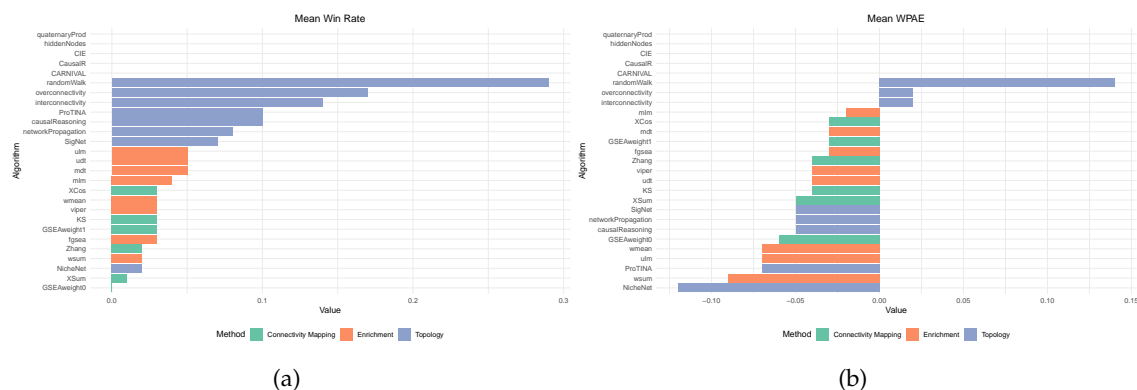


Figure 4.3: A. Mean win rate across all evaluations for each algorithm, showing the frequency of ranking the true target highest among competing methods. B. Mean WPAE for each algorithm, showing performance relative to random chance. Algorithms are colored by category.

target. The pathway enrichment scores (Figure 4.5) are a rescaled mean pathway enrichment metric across all runs, with values closer to 1 indicating higher levels of recovery of target-related pathways. Pathway enrichment performance displayed interesting patterns when compared with direct target recovery. Enrichment-based methods, which showed moderate performance in direct target recovery, demonstrated an improvement in pathway-level performance. The wmean and wsum methods achieved the best pathway scores (1.00), with average enrichment values of 2.15. This is confirmed by a low hypergeometric p-value, showing a strong overlap between the target pathways and the pathways enriched in the significant regulators predicted by the algorithm. The ulm follows with a mean enrichment of 1.98 (score of 0.92). In contrast to these are viper, mlm, fgsea, mdt, and udt, which performed poorly with values between 1 and 0.24 (scores between 0.45 and 0.08), respectively. Within topology-based methods, causalReasoning achieved the second-highest overall performance (enrichment: 1.97, score: 0.91), followed by overconnectivity (enrichment: 1.72, score: 0.79). Similar moderate performance was shown by networkPropagation, SigNet, and randomWalk (scores in the 0.54 - 0.55 range). Connectivity-mapping approaches uniformly failed at pathway enrichment, with all five methods achieving scores below 0.11.

The outperformance of some of the algorithms in this benchmark should be kept in mind when selecting the best computation approach to understand broad biological context rather than simply identifying single causal nodes.

## 4.2 Reference datasets

Assessing the influence of other components on algorithm performance is also important. The type of prior knowledge chosen to execute the algorithm is just as critical as the algorithm itself to provide meaningful results. For this reason, the nine reference datasets

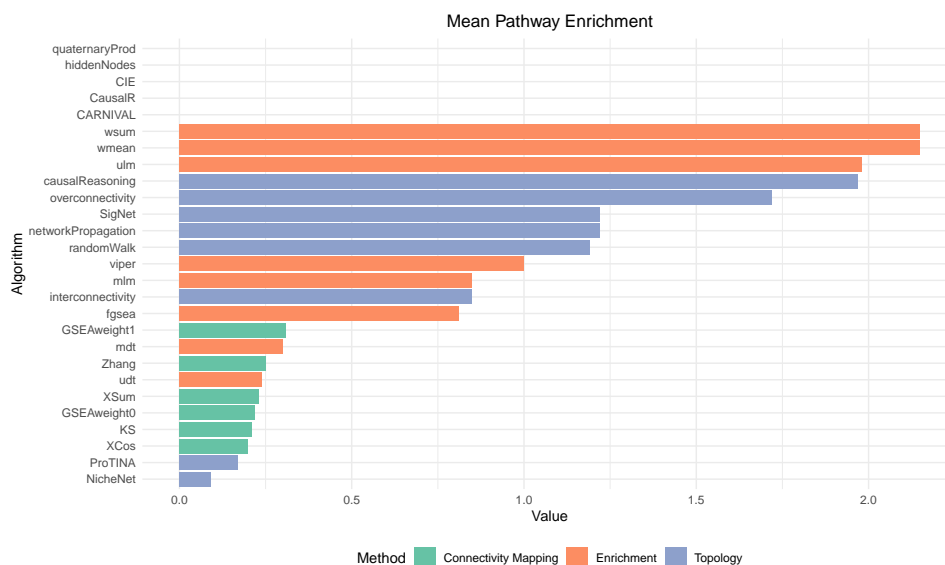


Figure 4.4: Mean pathway enrichment values for each algorithm, indicating the ability to identify biologically relevant pathways containing the true targets. Algorithms are colored by category.

were evaluated to understand which allows more effective recovery of targets identified from the gold standard. As before, four metrics were used to evaluate performance: coverage, win rate, target recovery, and pathway enrichment. For each reference dataset, the metric values and the final rank, with the final scores (values rescaled from 0 to 1) are represented in Figure 4.6. Evaluating the degree of coverage is crucial because if targets are not present in the prior knowledge provided to the algorithm, then they will never be returned in the results. Looking at the broader picture, Figure 4.6 shows the weight of the coverage in the final rank, which again highlights the importance of having well annotated and detailed prior knowledge. Overall, the MetaBase network was the best reference dataset. MetaBase network achieved the top rank in both coverage (92% of targets are present) and target recovery (AUC score of 1). OmniPath network ranked immediately next, given the substantial differences in coverage (63%). However, it has the highest win rate (and a WPAE of 0.05) and pathway enrichment scores, suggesting that the quality of annotation can balance the possible lack of specific regulators. For MetaBase network, the coverage is in line with the number of regulatory interactions present in this database, ensuring that no causal regulators are missed because of a lack in the prior knowledge dataset. What stands out most when considering performance of different reference datasets is the fact that topological network references generally score higher than signatures. Also, the clear contrast between MetaBase Network's 92% coverage and GWPS's 27% coverage directly translates to the performance differences across all other metrics evaluated. GWPS achieved a score of zero in all metrics, showing that even sophisticated algorithms cannot overcome insufficient molecular coverage in a prior knowledge reference data set. LINCS CRISPR performed best among perturbation references, but scored poorly in pathway





Figure 4.5: Final performance evaluation of 22 algorithms across five metrics. Heatmap showing algorithm performance rankings, where higher values indicate better performance. Algorithms (rows) are ordered by their final aggregate rank from best (top) to worst (bottom). Each cell displays the normalized performance score (0-1 scale) for the corresponding algorithm-metric combination, with cell colors representing the rank (1-24). The color gradient ranges from red (low rank/poor performance) through yellow (medium performance) to green (high rank/excellent performance).

enrichment (0.50), indicating that perturbation signatures may lack pathway-level context. Finally, it is important to mention the performance of regulon-based references. MetaBase Linear Path Regulons achieved the best pathway enrichment (score of 1) in contrast with the poor coverage (48%) and target recovery (AUC: 0.18). This can indicate that although this dataset captures well the biological context, it probably oversimplifies the complexity needed for precise target identification.

### 4.3 Query datasets

It is essential to comprehend the performance of all components used in computational algorithms, and the query dataset is no exception. An assessment at the individual and group level was performed for the query datasets using the AUC and pathway enrichment metrics. The goal was to evaluate which perturbation approach enables more effective recovery of the gold standard perturbed genes or targets. At the individual dataset level

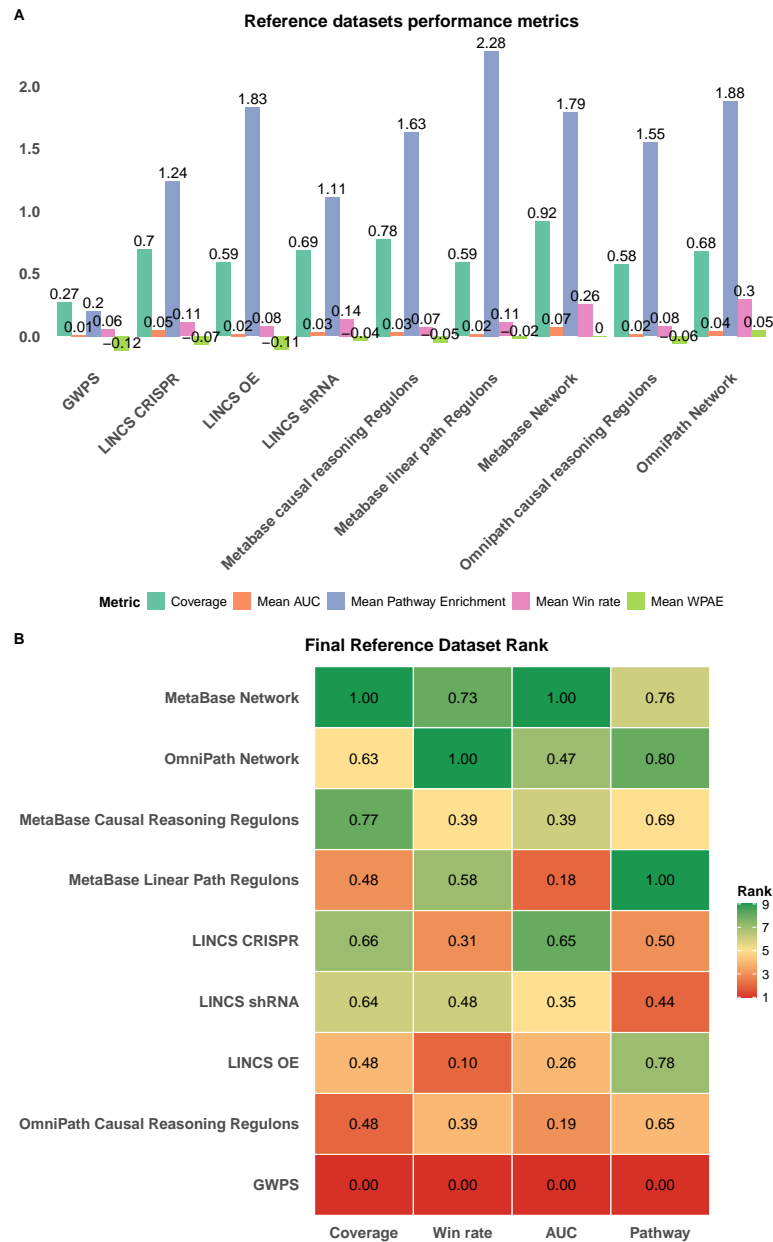


Figure 4.6: Performance of nine reference datasets across four evaluation metrics. A. Performance mean values. For each reference dataset, the performance is colored by metric. The mean pathway enrichment values represent the significance of overlap between true target pathways and pathways enriched with significant regulators ( $-\log_{10}(\text{p-value})$  of the hypergeometric test) averaged across the signatures where true targets are present in one or more MetaBase pathways. B. Performance ranking with final scores. Heatmap showing reference dataset performance rankings, where higher ranks indicate better performance. Reference datasets (rows) are ordered by their final aggregate rank from best (top) to worst (bottom). Each cell displays the scaled performance score for the corresponding dataset-metric combination, with cell colors representing the rank. The evaluation metrics include Coverage, Win rate, AUC, and Pathway Enrichment. The color gradient ranges from red (low rank/poor performance) through yellow (medium performance) to green (high rank/excellent performance).

(Figure 4.7), each query was evaluated separately. Then, to systematically investigate performance patterns, algorithms were evaluated across grouped dataset categories: perturbation types (drug vs genetic), sample sources (tissue vs cell-line), and sequencing approaches (bulk RNA-seq vs scRNA-seq) (Figure 4.8).

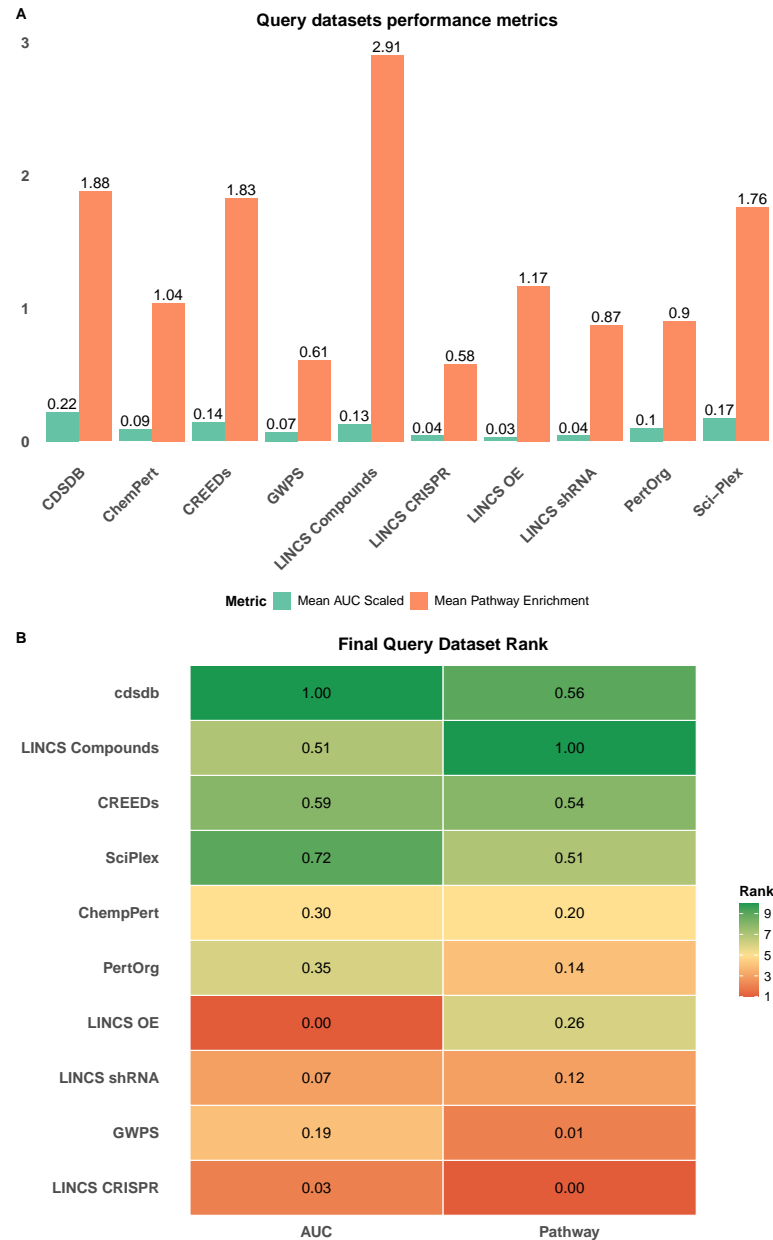


Figure 4.7: Query dataset performance evaluation. A. Bar plot displaying performance metrics (AUC and pathway enrichment scores) for each query dataset, with values shown above bars. B. Heatmap visualization of query dataset rankings based on AUC and pathway enrichment performance, ordered by final aggregate rank from lowest (bottom) to highest (top). Cell values represent normalized scores (0-1 scale) while color intensity indicates rank position, with green representing better performance. Query datasets include both drug and genetic perturbation experiments from various sources including cell lines, tissues, and single-cell platforms.

CDS-DB, a drug perturbation dataset derived entirely from cancer patient samples, obtained the highest overall performance with AUC and pathway enrichment scores of 1.00 and 0.56, respectively. CDS-DB, CREEDS, and PertOrg, some of the top-performing datasets, represent tissue or *in vivo* samples. This pattern is also visible when comparing by grouping tissue versus cell-line (Figure 4.8). The tissue versus cell-line comparison showed better AUC for tissue-derived signatures, aligning with the individual dataset observations. The drug versus genetic perturbation comparison also revealed a pattern. Algorithms consistently performed better on drug perturbation datasets. This distinction is evident from most algorithms falling below the diagonal line for AUC. Drug perturbations apparently are more reliable signatures for target recovery. This is even more pronounced individually (Figure 4.7) for LINCS, which comes from the same platform. LINCS Compounds ranked second overall with the best mean pathway enrichment value. On the other hand, genetic perturbations from the same LINCS platform (CRISPR, shRNA, and OE) showed poor target recovery with AUC values between 0.03 and 0.04, with LINCS CRISPR obtaining the last position in the rank. The bulk versus single-cell comparison revealed no specific differences, with algorithms grouping evenly around the diagonal in Figure 4.8. Looking at individual single-cell datasets, Sci-Plex (AUC score: 0.72; Pathway enrichment score: 0.51) performed far better than GWPS (AUC score: 0.19; Pathway enrichment score: 0.01). The inconsistent performance of single-cell and bulk query datasets may suggest that the sequencing technology by itself is not the factor most responsible for the variability of algorithm performance.

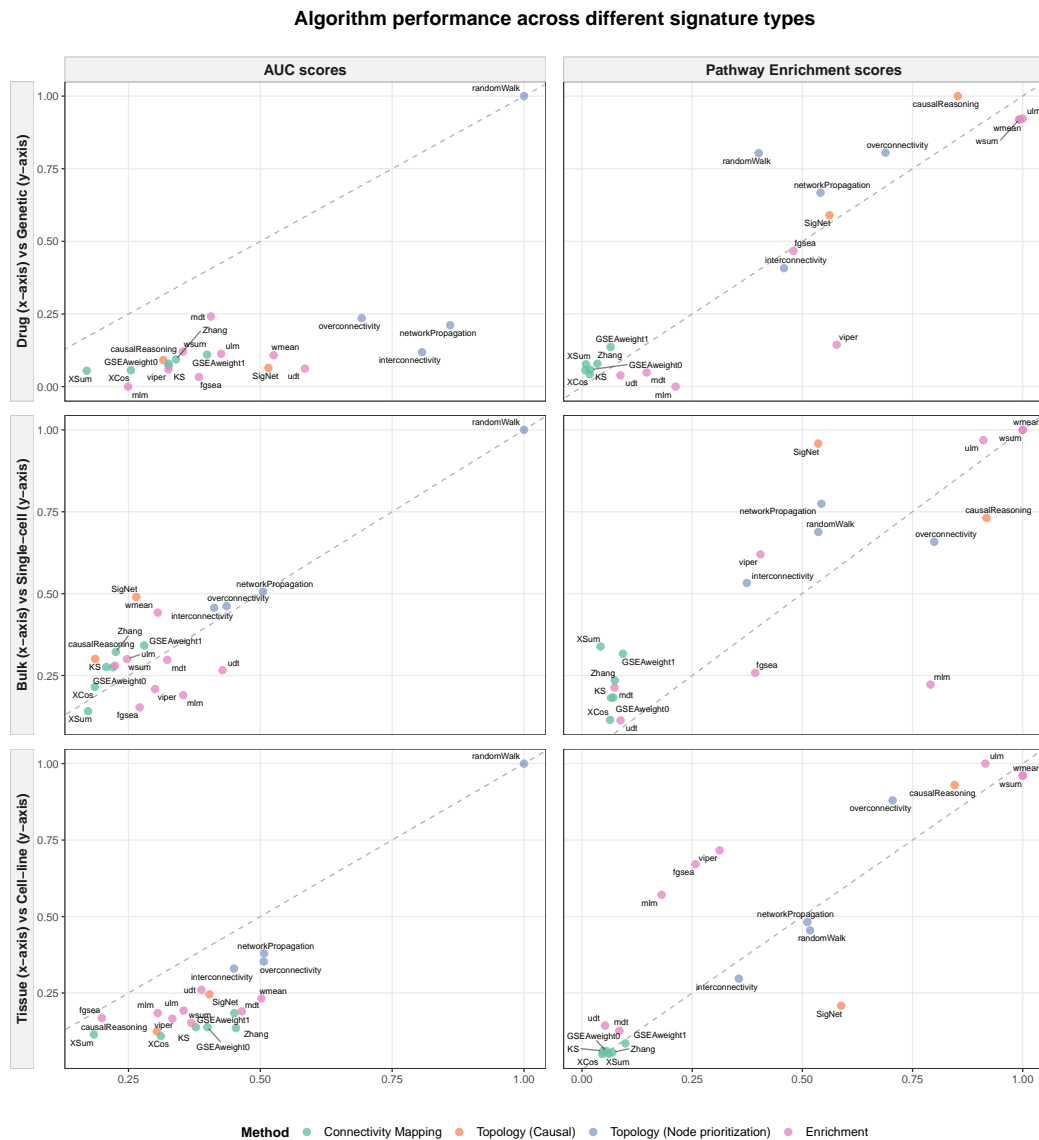


Figure 4.8: Comparison of algorithm performance between different signature types, across two evaluation metrics (AUC scores and Pathway Enrichment scores). Different comparison: drug perturbation vs genetic perturbation signatures (top), bulk RNA-seq vs scRNA-seq data (middle), and tissue-derived vs cell-line-derived signatures (bottom). Each dot represents an algorithm, colored by category. The diagonal dashed line represents equal performance between the two signature types being compared. Dots above the line indicate better performance on the y-axis signature type, while dots below the line indicate better performance on the x-axis signature type.

## DISCUSSION

*This chapter discusses the results by highlighting their implications and comparing them with the existing literature. Finally, the key findings are summarized with conclusions and suggestions for some future perspectives.*

The MoA of a small molecule involves its interaction with specific molecular targets. When these primary targets are activated or inhibited, they trigger changes in downstream signaling pathways, modulating the activity of TFs and therefore the expression of individual genes. All these changes will ultimately produce the desired drug effect together with any additional undesired side effects. As with many topics in the scientific community, there are different views on the real need for a complete understanding of the molecular target and the MoA of a drug [106]. Some defend and prioritize clinical benefit, given the high number of drugs on the market for which the MoA is currently unknown. On the other hand, some argue that understanding the effects of compounds is essential in the early stages of DD. Knowing the MoA of a drug not only helps to guide the process but also increases predictability and can even help to understand potential side effects. While understanding these mechanisms isn't strictly necessary for successful drug development, it plays a crucial role in improving efficiency. For this reason, advances in system biology are focused on reducing the process for MoA reconstruction, keeping reliable results. For example, in 2020, the CMap Institute created a challenge on the Kaggle platform [107] to develop an algorithm to predict the MoA of a new drug. Along with these collective efforts, several bioinformatics tools have emerged to infer causal regulators. Authors who develop a new algorithm typically perform comparative benchmarking to highlight the novelty of their discovery [94]. However, these studies often have some limitations. The main one is that simulated data, such as the LINCS dataset, is used instead of real experimental data. This approach may fail to capture the variability and complexity of biological systems, leading to a biased evaluation of the computation tool. Also, the datasets used are often small and curated, which may not reflect the challenges of real-world applications. Finally, the evaluations often focus on a single aspect of performance, such as accuracy, without considering other important factors like scalability, robustness, and computational efficiency. Benchmarking studies are therefore essential because they

aim to test various components in an unbiased and realistic way, to assess how well different algorithm classes perform.

This study represents one of the most comprehensive evaluations of MoA recovery methods published to date to our knowledge. While previous studies have focused on specific categories of methods, limited network sources, or small selected datasets, this study covered 27 algorithms, including topology-based, enrichment, and CMap methods, incorporating public and commercial reference datasets, and experimental and simulated data from seven distinct sources. This extends previous benchmarking efforts and shows critical differences between theory and practical feasibility.

Topology-based methods consistently achieved the best performance across all metrics in direct target retrieval. Consistent with other studies, randomWalk was at the top of all the rankings with the exception of runtime and pathway analysis. In Hill *et al.* [59], node prioritization methods, including randomWalk followed by networkPropagation and GeneMANIA, were also identified as the top performers, using a combination of interactions from STRING, MetaBase, and BioPlex as reference dataset. Our study confirmed the superiority of randomWalk, together with evidence that performance depends on the choice of reference network, with the Metabase network providing the best results, probably due to its greater molecular coverage.

The similar performance of randomWalk with different reference datasets demonstrates the accuracy of the algorithm. The overconnectivity and networkPropagation algorithms also performed well in our study. However, each one had some limitations in terms of accuracy, reliability, or run time. For example, networkPropagation obtained results similar to randomWalk in terms of AUC, but with reduced reliability in large-scale executions, suggesting possible concerns with scalability.

Our study reveals an important trade-off between direct recovery and the biological context, not previously characterized in other studies. Methods that achieved the highest pathway enrichment scores, such as wmean and wsum, also obtained negative WPAE values, indicating that they do not outperform direct target classification methods. On the other hand, although randomWalk performs well in target identification, it seems to compromise pathway enrichment performance, by obtaining half the scores of wmean and wsum. This inverse relationship suggests that algorithms designed to capture biological context end up sacrificing accuracy in identifying individual causal nodes. As a consequence, methods selection goes beyond simple performance rankings, as it should also account for the need for a deeper understanding of the biological context surrounding a specific target. If understanding the broader biological context is more important than identifying single targets, these enrichment methods can provide better results, despite lower accuracy. In contrast, when experimental validation resources limit testing to a few candidates, the accuracy of randomWalk seems to be more appropriate.

A study by Lin, K., *et al.* [86] and others [108] identified ZhangScore as the best algorithm among CMap methods. However, in our study, this algorithm performed poorly. In general, all CMap algorithm achieved low AUC values ( $< 0.08$ ), win rate values ( $< 0.03$ ),

and pathway enrichment scores (0.20 0.31). Among the CMap algorithms, GSEAweight1 showed the best results, followed by KS, with Zhang algorithm distinguished itself only by having the worst computational efficiency (highest runtime). This huge difference compared to published results can be interpreted considering the task being evaluated. While the studies mentioned above were assessing the similarity between signatures, in this study we are evaluating the ability to identify the true target responsible for the perturbation signatures. These results suggest that similarity-based algorithm fails to capture the causal regulatory relationships underlying perturbation responses, meaning that gene expression pattern relationship does not necessarily reflects shared regulatory mechanism.

Consistent with Hill *et al.* [59], we observed that causal reasoning methods performance is better at pathway-level enrichment, compared with precise target identification. For example, causalReasoning achieved perfect reliability, but modest AUC, in line with the observation that these methods struggle to pinpoint exact drug targets. This pattern was observed for SigNet and CausalR [27] and it can reflect that many regulatory targets are not TFs or are not directly observable in expression data. The most interesting and surprising finding from this study is the issue of scalability for causal reasoning algorithms, which resulted in complete failure of the most sophisticated among them (CARNIVAL). In previous benchmarking studies, SigNet and CARNIVAL had the best performance [27], with OmniPath network as reference data. However, despite CARNIVAL's theoretical ability for capturing complex regulatory interactions, its computational requirements make it unusable for realistic dataset sizes. This emphasizes that algorithm benchmarking must consider not just accuracy but also scalability and runtimes, that are often neglected when evaluations focus on small and curated datasets. In this large-scale evaluation several computationally intensive algorithms (e.g., NicheNet, ProTINA) are impractical for high-throughput use, with runtimes exceeding 400 seconds per signature or total failure to complete the analyses. For methods intended for integration into large-scale screening pipelines, these constraints represent a critical limitation.

With a strong influence on algorithm performance, the importance of network choice was no exception in our study. Metabase demonstrated the most consistent target recovery results, while others appear better suited for pathway enrichment. Network selection should be aligned with the research goal, with denser and high-coverage networks for precision target recovery and more structured, curated pathway resources for context mapping.

The consistent superiority of drug perturbations over genetic approaches is important to highlight. The assumption that genetic perturbations are more precise can be refuted by our results. One possible explanation is that drugs often have multiple targets, leading to broader and more robust transcriptomic changes that are easier to detect. They also tend to produce coherent, dose-dependent effects on pathways without triggering compensatory or stress responses [108], which can make them easier for algorithms to trace back to upstream regulators. Also, the poor performance of LINCS genetic perturbations highlights



potential limitations of cell line-based genetic screens. The validation approaches used by the original dataset creators varied considerably, affecting result interpretation. Therefore, the quality of the datasets can be a determining factor in the performance of these tools. Among the seven sources of datasets used, based on the validation carried out by the authors, some datasets do not seem to be the most suitable for obtaining reliable results. CREEDS is one such case where the data is extracted from public databases and not generated by the researchers involved in the benchmarking effort [51]. The data is only checked on a technical level, to confirm that the samples are from GEO, and if they are labeled correctly. The attempt at biological evaluation, by looking for specific patterns of signatures (i.e. if two signatures come from perturbing the same gene), showed inconsistent results. Moreover, for some datasets validation was practically non-existent [39, 55], and the filters used during the pre-processing step of this study actually increased the quality of the data (this is the case of ChemPert and PertOrg). Both single-cell datasets, GWPS and Sci-Plex, showed a good validation of the data. GWPS, focused on gene knockdown signatures, used strict criteria to define significant responses [45]. Sci-Plex, with drug perturbations, validated its data through some complementary statistical analysis [41]. The data from LINCS are more controversial. First, L1000 measures around 978 landmark genes and computationally infers the rest (~12,000). Additionally, despite each perturbation is tested in triplicate and z-scores are adjusted to minimize replicate inconsistencies [8], the reliability of these data is still questionable. Recent alternatives, such as DRUG-seq, demonstrate to have some advantages over L1000. Direct comparison between the two datasets [52, 53] showed that DRUG-seq directly measures more than 10,000 genes without relying on inference, at a lower cost. When testing the accuracy of both, DRUG-Seq proved to be more accurate in distinguishing samples between different diseases. Also, the emergence of Perturb-seq, which incorporates CRISPR perturbations with scRNA-seq [45], offers another promising alternative with greater transcriptome coverage at a reduced cost [109]. Finally, the cancer cell line composition is another shortcoming of this dataset limiting applicability to other contexts. This is being addressed by programs such as NeuroLINCS [110], which create signatures using patient-derived induced pluripotent stem cells [109]. Our finding that tissue-derived signatures consistently outperformed cell line data supports the expansion of data derived from other sources.

The critical importance of reference network selection extends beyond simple coverage metrics. Large-scale networks may include interactions irrelevant to specific cellular contexts [9]. For this reason, tools for constraining networks based on tissue-specific expression data from Human Protein Atlas [111] or TissueNet [65] could improve performance [9]. Our results showing MetaBase Linear Path Regulons' excellent pathway enrichment, despite poor coverage, support this tissue-specific approach.

The focus on transcriptomic-based methods, in this thesis, but also other benchmarking studies, represents just one facet of MoA inference. Recent multi-omics integration tools like SignalingProfiler 2.0 [84] and COSMOS [112] demonstrate the value of combining transcriptomics with proteomics, metabolomics, and phosphoproteomics data. The use

of COSMOS to identify known clear cell renal cell carcinoma therapeutic targets was a success. This was achieved by capturing crosstalks within and between different omics layers, illustrating how the integration of multiple data types can generate novel biomarker hypotheses. Additionally, data capturing changes in cell morphology and the chemical structure of specific compounds [113] can become a valuable complement to expression data in DD studies [9].

Some considerations and future work that can be explored are related with certain adjustments or experimental directions that could have been considered, but have been left out of this study for time and resource limitations. One aspect that has not been considered, but could effectively influence the performance of the algorithms is parameter optimization. Algorithm behavior and performance can change dramatically with parameter tuning, and it is suggested as good practice for benchmarking studies [93]. Yet comprehensive optimization across all algorithm-dataset combinations proved computationally unreasonable. The same applies to the methods for filtering DEGs from full signatures. However, the parameters of the algorithms and the preparation of the input data are specified and calculated in the wrapper function and they can be changed if needed. Being this a benchmarking study, the entire workflow is set up precisely to test multiple options, even though not all of them have been included in this work. When evaluating, another metric that could be included was an ensemble score. Instead of relying on a single tool or data set, a score combining the results of several tools could also be generated to produce more robust predictions. Future benchmarking efforts can also weight results by validation confidence, given that the heterogeneity in dataset validation can partially explain some performance variations across query datasets.

The Molecular networks used as a reference in this study, both MetaBase and OmniPath, are of considerable size, with hundreds of thousands of interactions. When using these networks, one of the recommendations is to ensure that the interactions present are within the context of the study, in this case, the type of cell or tissue being studied [9]. This reasoning makes sense if we consider that these large networks include all kinds of interactions, including those specific to other cells or tissues. To this end, instead of using global networks, more specific and context-relevant networks can be used, which may already be prepared by the databases themselves. Alternatively, databases such as TissueNet [65], provide interaction networks that are specific to certain tissues. Another thing that could be considered would be to integrate networks from different sources, for example, MetaBase and OmniPath networks, customizing the subsets of interactions to suit the context and by keeping only those with higher reliability scores. If the focus extends beyond the types of algorithms used in this project, additional software and tools are available for target identification that leverage different computational approaches. Depending on the specific scientific questions being investigated, it could be also worth exploring them. One example is Drug2ways [114], a software implemented in Python, to identify potential drug repositioning and predict the effects of drugs. To do this, it performs causal reasoning through biological networks with three types of vertices:

molecular entities, drugs, and indications. The paths between the drug and the indication are noted, as well as their direction, suggesting positive or negative regulation. The most frequent direction is taken as the predicted effect. Solutions are becoming increasingly innovative, taking advantage of advances in computer efficiency. Such as the use of knowledge graphs and Large language models for drug repurposing [115].

In summary, based on the results obtained and the intersection of previous research, this study recommends the randomWalk algorithm for precise target recovery, as well as wmean for pathway-level context. It is crucial to highlight the importance of high-quality reference data, and for topology-based tools, to consider filtering interactions relevant to the study in question. Regarding query data, the use of experimental data and data from tissues should be preferred. In addition, other algorithm optimization parameters should be considered, rather than the default settings. Finally, given the progresses in the field of system biology and the avalanche of increasingly available data, integration of other levels of omics data and the use of machine learning methods promise to further improve the ability to identify target genes and perturbed pathways.

## REFERENCES

- [1] QuillBot. *QuillBot: Your complete writing solution*. en. URL: <https://quillbot.com/> (cit. on p. i).
- [2] Microsoft. *Microsoft Copilot*. en. URL: <https://copilot.microsoft.com/> (cit. on p. i).
- [3] J. M. Loureno. *The NOVAthesis L<sup>A</sup>T<sub>E</sub>X Template User's Manual*. NOVA University Lisbon. 2021. URL: <https://github.com/joaomlourenco/novathesis/raw/main/template.pdf> (cit. on p. i).
- [4] D. Cook et al. "Lessons learned from the fate of AstraZeneca's drug pipeline: a five-dimensional framework". In: *Nature Reviews Drug Discovery* 13.6 (2014), pp. 419–431. ISSN: 1474-1784. DOI: 10.1038/nrd4309 (cit. on p. 1).
- [5] X. Ji, J. M. Freudenberg, and P. Agarwal. "Integrating Biological Networks for Drug Target Prediction and Prioritization". In: *Computational Methods for Drug Repurposing*. Ed. by Q. Vanhaelen. New York, NY: Springer New York, 2019, pp. 203–218. ISBN: 978-1-4939-8955-3. DOI: 10.1007/978-1-4939-8955-3\_12 (cit. on p. 1).
- [6] M. Wang et al. "TREAP: A New Topological Approach to Drug Target Inference". In: *Biophysical Journal* 119.11 (2020), pp. 2290–2298. ISSN: 0006-3495. DOI: 10.1016/j.bpj.2020.10.021 (cit. on p. 1).
- [7] K. Zhao and H.-C. So. "Using Drug Expression Profiles and Machine Learning Approach for Drug Repurposing". In: *Computational Methods for Drug Repurposing*. Ed. by Q. Vanhaelen. New York, NY: Springer New York, 2019, pp. 219–237. ISBN: 978-1-4939-8955-3. DOI: 10.1007/978-1-4939-8955-3\_13 (cit. on p. 1).
- [8] A. Subramanian et al. "A Next Generation Connectivity Map: L1000 Platform and the First 1,000,000 Profiles". In: *Cell* 171.6 (2017), 1437–1452.e17. DOI: 10.1016/j.cell.2017.10.049 (cit. on pp. 1, 9, 10, 14, 36, 66).

## REFERENCES

---

- [9] M.-A. Trapotsi, L. Hosseini-Gerami, and A. Bender. "Computational analyses of mechanism of action (MoA): data, methods and integration". In: *RSC Chemical Biology* 3.2 (2022), pp. 170–200. DOI: 10.1039/D1CB00069A (cit. on pp. 1, 2, 7, 8, 10, 19, 20, 26, 66, 67).
- [10] G. Bradley and S. J. Barrett. "CausalR: extracting mechanistic sense from genome scale data". In: *Bioinformatics* 33.22 (2017), pp. 3670–3672. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btx425 (cit. on pp. 2, 29, 40, 44).
- [11] J. Lamb et al. "The Connectivity Map: using gene-expression signatures to connect small molecules, genes, and disease". In: *Science* 313.5795 (2006), pp. 1929–35. ISSN: 0036-8075. DOI: 10.1126/science.1132939 (cit. on pp. 2, 9, 14, 26, 27, 41).
- [12] P. Badia-i Mompel et al. "decoupleR: ensemble of computational methods to infer biological activities from omics data". In: *Bioinformatics Advances* 2.1 (2022). ISSN: 2635-0041. DOI: 10.1093/bioadv/vbac016 (cit. on pp. 2, 40, 41).
- [13] Clarivate. *CBDD, Computational biology methods for drug discovery*. Web Page. 20.0.3. URL: <https://clarivate.com/life-sciences-healthcare/consulting-services/research-and-development-consulting/cbdd/> (cit. on p. 2).
- [14] Johansson et al. "Precision medicine in complex diseases-Molecular subgrouping for improved prediction and treatment stratification". In: *J Intern Med* 294.4 (2023), pp. 378–396. DOI: 10.1111/joim.13640 (cit. on p. 6).
- [15] J. Avorn. "The 2.6 Billion Dollar Pill Methodologic and Policy Considerations". In: *New England Journal of Medicine* 372.20 (2015), pp. 1877–1879. DOI: 10.1056/NEJMp1500848 (cit. on p. 6).
- [16] Y. Wang, J. Yella, and A. G. Jegga. "Transcriptomic Data Mining and Repurposing for Computational Drug Discovery". In: *Computational Methods for Drug Repurposing*. Ed. by Q. Vanhaelen. New York, NY: Springer New York, 2019, pp. 73–95. ISBN: 978-1-4939-8955-3. DOI: 10.1007/978-1-4939-8955-3\_5 (cit. on pp. 6, 26).
- [17] ERG. *Drug Development - Final Report*. Report. 2024. URL: <https://aspe.hhs.gov/reports/drug-development> (cit. on p. 6).
- [18] Z. Wu, Y. Wang, and L. Chen. "Network-based drug repositioning". In: *Molecular BioSystems* 9.6 (2013), pp. 1268–1281. ISSN: 1742-206X. DOI: 10.1039/C3MB25382A (cit. on p. 6).
- [19] F. Sohraby, M. Bagheri, and H. Aryapour. "Performing an In Silico Repurposing of Existing Drugs by Combining Virtual Screening and Molecular Dynamics Simulation". In: *Computational Methods for Drug Repurposing*. Ed. by Q. Vanhaelen. New York, NY: Springer New York, 2019, pp. 23–43. ISBN: 978-1-4939-8955-3. DOI: 10.1007/978-1-4939-8955-3\_2 (cit. on pp. 6, 7).
- [20] K. Park. "A review of computational drug repurposing". In: *Transl Clin Pharmacol* 27.2 (2019), pp. 59–63. DOI: 10.12793/tcp.2019.27.2.59 (cit. on pp. 6, 7, 19).

- 
- [21] S. Chakraborti, G. Ramakrishnan, and N. Srinivasan. "Repurposing Drugs Based on Evolutionary Relationships Between Targets of Approved Drugs and Proteins of Interest". In: *Computational Methods for Drug Repurposing*. Ed. by Q. Vanhaelen. New York, NY: Springer New York, 2019, pp. 45–59. ISBN: 978-1-4939-8955-3. DOI: 10.1007/978-1-4939-8955-3\_3 (cit. on p. 7).
- [22] B. B. Kakoti, R. Bezbaruah, and N. Ahmed. "Therapeutic drug repositioning with special emphasis on neurodegenerative diseases: Threats and issues". In: *Front Pharmacol* 13 (2022), p. 1007315. DOI: 10.3389/fphar.2022.1007315 (cit. on p. 7).
- [23] A. E. Kel. "Search for Master Regulators in Walking Cancer Pathways". In: *Biological Networks and Pathway Analysis*. Ed. by T. V. Tatarinova and Y. Nikolsky. New York, NY: Springer New York, 2017, pp. 161–191. ISBN: 978-1-4939-7027-8. DOI: 10.1007/978-1-4939-7027-8\_8 (cit. on p. 7).
- [24] J. K. Ryu and J. G. McLarnon. "Thalidomide inhibition of perturbed vasculature and glial-derived tumor necrosis factor-alpha in an animal model of inflamed Alzheimer's disease brain". In: *Neurobiol Dis* 29.2 (2008), pp. 254–66. DOI: 10.1016/j.nbd.2007.08.019 (cit. on p. 7).
- [25] S. Kato et al. "Challenges and perspective of drug repurposing strategies in early phase clinical trials". In: *Oncoscience* 2.6 (2015), pp. 576–80. DOI: 10.18632/oncoscience.173 (cit. on p. 7).
- [26] A. Antona et al. "Dissecting the Mechanism of Action of Spiperone-A Candidate for Drug Repurposing for Colorectal Cancer". In: *Cancers (Basel)* 14.3 (2022). DOI: 10.3390/cancers14030776 (cit. on p. 7).
- [27] L. Gennari et al. "Raloxifene in breast cancer prevention". In: *Expert Opin Drug Saf* 7.3 (2008), pp. 259–70. ISSN: 1474-0338. DOI: 10.1517/14740338.7.3.259 (cit. on p. 7).
- [28] M. Iwata et al. "Elucidating the modes of action for bioactive compounds in a cell-specific manner by large-scale chemically-induced transcriptomics". In: *Scientific Reports* 7.1 (2017), p. 40164. ISSN: 2045-2322. DOI: 10.1038/srep40164 (cit. on p. 7).
- [29] L. Hosseini-Gerami et al. "Benchmarking causal reasoning algorithms for gene expression-based compound mechanism of action analysis". In: *BMC Bioinformatics* 24.1 (2023), p. 154. ISSN: 1471-2105. DOI: 10.1186/s12859-023-05277-1 (cit. on pp. 7, 30, 31).
- [30] "Mechanism matters". In: *Nature Medicine* 16.4 (2010), pp. 347–347. ISSN: 1546-170X. DOI: 10.1038/nm0410-347 (cit. on p. 7).
- [31] I. Bezprozvanny. "The rise and fall of Dimebon". In: *Drug News Perspect* 23.8 (2010), pp. 518–23. DOI: 10.1358/dnp.2010.23.8.1500435 (cit. on p. 7).

## REFERENCES

---

- [32] M. Pilarczyk et al. "Connecting omics signatures and revealing biological mechanisms with iLINCS". In: *Nature Communications* 13.1 (2022), p. 4678. ISSN: 2041-1723. DOI: 10.1038/s41467-022-32205-3 (cit. on pp. 9, 14).
- [33] A. B. Keenan et al. "Connectivity Mapping: Methods and Applications". In: *Annual Review of Biomedical Data Science* 2. Volume 2, 2019 (2019), pp. 69–92. DOI: 10.1146/annurev-biodatasci-072018-021211 (cit. on pp. 9, 10, 18).
- [34] R. B. Stoughton and S. H. Friend. "How molecular profiling could revolutionize drug discovery". In: *Nature Reviews Drug Discovery* 4.4 (2005), pp. 345–350. ISSN: 1474-1776. DOI: 10.1038/nrd1696 (cit. on p. 9).
- [35] T. R. Hughes et al. "Functional discovery via a compendium of expression profiles". In: *Cell* 102.1 (2000), pp. 109–126. ISSN: 0092-8674. DOI: 10.1016/S0092-8674(00)00015-5 (cit. on p. 9).
- [36] B. Ganter et al. "Development of a large-scale chemogenomics database to improve drug candidate selection and to understand mechanisms of chemical toxicity and action". In: *Journal of biotechnology* 119.3 (2005), pp. 219–244. ISSN: 0168-1656. DOI: 10.1016/j.jbiotec.2005.03.022 (cit. on pp. 9, 14).
- [37] A. Engelberg. "Iconix Pharmaceuticals, Inc. removing barriers to efficient drug discovery through chemogenomics". In: *Pharmacogenomics* 5.6 (2004), pp. 741–744. ISSN: 1462-2416. DOI: 10.1517/14622416.5.6.741 (cit. on pp. 9, 14).
- [38] B. Alexander-Dann et al. "Developments in toxicogenomics: understanding and predicting compound-induced toxicity from gene expression data". In: *Mol Omics* 14.4 (2018), pp. 218–236. ISSN: 2515-4184. DOI: 10.1039/c8mo00042e (cit. on pp. 9, 13).
- [39] M. Zheng et al. "ChemPert: mapping between chemical perturbation and transcriptional response for non-cancer cells". In: *Nucleic Acids Research* 51.D1 (2023), pp. D877–D889. DOI: 10.1093/nar/gkac862 (cit. on pp. 9, 11, 13, 27, 36, 66).
- [40] Z. Liu et al. "CDS-DB, an omnibus for patient-derived gene expression signatures induced by cancer treatment". In: *Nucleic Acids Research* 52.D1 (2024), pp. D1163–D1179. DOI: 10.1093/nar/gkad888 (cit. on pp. 10, 11, 13, 27, 36).
- [41] S. R. Srivatsan et al. "Massively multiplex chemical transcriptomics at single-cell resolution". In: *Science* 367.6473 (2020), pp. 45–51. DOI: 10.1126/science.aax6234 (cit. on pp. 11, 17, 36, 66).
- [42] Z. Wei et al. "PerturBase: a comprehensive database for single-cell perturbation data analysis and visualization". In: *Nucleic Acids Research* 53.D1 (2024), pp. D1099–D1111. ISSN: 1362-4962. DOI: 10.1093/nar/gkae858 (cit. on pp. 11, 16).
- [43] J. Cao et al. "Comprehensive single-cell transcriptional profiling of a multicellular organism". In: *Science* 357.6352 (2017), pp. 661–667. DOI: 10.1126/science.aam8940 (cit. on p. 11).

- [44] B. K. Martin et al. "Optimized single-nucleus transcriptional profiling by combinatorial indexing". In: *Nature Protocols* 18.1 (2023), pp. 188–207. ISSN: 1750-2799. DOI: 10.1038/s41596-022-00752-0 (cit. on p. 11).
- [45] J. M. Replogle et al. "Mapping information-rich genotype-phenotype landscapes with genome-scale Perturb-seq". In: *Cell* 185.14 (2022), 2559–2575.e28. DOI: 10.1016/j.cell.2022.05.013 (cit. on pp. 11, 12, 15, 18, 36, 66).
- [46] B. Ganter et al. "Toxicogenomics in drug discovery and development: mechanistic analysis of compound/class-dependent effects using the DrugMatrix database". In: *Pharmacogenomics* 7.7 (2006), pp. 1025–44. DOI: 10.2217/14622416.7.7.1025 (cit. on pp. 13, 17, 27).
- [47] A. Brazma et al. "ArrayExpressa public repository for microarray gene expression data at the EBI". In: *Nucleic Acids Research* 31.1 (2003), pp. 68–71. ISSN: 0305-1048. DOI: 10.1093/nar/gkg091 (cit. on pp. 13, 17, 18).
- [48] R. Edgar, M. Domrachev, and A. E. Lash. "Gene Expression Omnibus: NCBI gene expression and hybridization array data repository". In: *Nucleic Acids Research* 30.1 (2002), pp. 207–210. ISSN: 0305-1048. DOI: 10.1093/nar/30.1.207 (cit. on pp. 13–15, 17, 18).
- [49] C. Martini et al. "CEBS update: curated toxicology database with enhanced tools for data integration". In: *Nucleic Acids Research* 50.D1 (2022), pp. D1156–D1163. DOI: 10.1093/nar/gkab981 (cit. on p. 13).
- [50] M. Waters et al. "CEBS—Chemical Effects in Biological Systems: a public data repository integrating study design and toxicity data with microarray and proteomics data". In: *Nucleic Acids Research* 36.Database issue (2008), pp. D892–900. DOI: 10.1093/nar/gkm755 (cit. on p. 13).
- [51] Z. Wang et al. "Extraction and analysis of signatures from the Gene Expression Omnibus by the crowd". In: *Nat Commun* 7 (2016), p. 12846. ISSN: 2041-1723. DOI: 10.1038/ncomms12846 (cit. on pp. 14, 18, 36, 66).
- [52] C. Ye et al. "DRUG-seq for miniaturized high-throughput transcriptome profiling in drug discovery". In: *Nature Communications* 9.1 (2018), p. 4307. ISSN: 2041-1723. DOI: 10.1038/s41467-018-06500-x (cit. on pp. 15, 66).
- [53] J. Li et al. "DRUG-seq Provides Unbiased Biological Activity Readouts for Drug Discovery". In: *bioRxiv* (2021), p. 2021.06.07.447456. DOI: 10.1101/2021.06.07.447456 (cit. on pp. 15, 66).
- [54] Y. Igarashi et al. "Open TG-GATEs: a large-scale toxicogenomics database". In: *Nucleic acids research* 43.D1 (2015), pp. D921–D927. ISSN: 1362-4962. DOI: 10.1093/nar/gku955 (cit. on pp. 16, 17).



## REFERENCES

---

- [55] Z. Zhai et al. "PertOrg 1.0: a comprehensive resource of multilevel alterations induced in model organisms by in vivo genetic perturbation". In: *Nucleic Acids Research* 51.D1 (2023), pp. D1094–d1101. DOI: 10.1093/nar/gkac872 (cit. on pp. 16, 18, 36, 66).
- [56] Y. Zhang et al. "PerturbAtlas: a comprehensive atlas of public genetic perturbation bulk RNA-seq datasets". In: *Nucleic Acids Research* 53.D1 (2024), pp. D1112–D1119. ISSN: 1362-4962. DOI: 10.1093/nar/gkae851 (cit. on p. 17).
- [57] D. Hadley et al. "Precision annotation of digital samples in NCBI's gene expression omnibus". In: *Scientific data* 4.1 (2017), pp. 1–11. ISSN: 2052-4463. DOI: 10.1038/sdata.2017.125 (cit. on p. 17).
- [58] S. K. Nair et al. "ToxicoDB: an integrated database to mine and visualize large-scale toxicogenomic datasets". In: *Nucleic Acids Research* 48.W1 (2020), W455–W462. ISSN: 0305-1048. DOI: 10.1093/nar/gkaa390 (cit. on p. 17).
- [59] A. Hill et al. "Benchmarking network algorithms for contextualizing genes of interest". In: *PLOS Computational Biology* 15.12 (2019), e1007403. DOI: 10.1371/journal.pcbi.1007403 (cit. on pp. 19, 26, 31, 64, 65).
- [60] S. Miller-Dott et al. "Expanding the coverage of regulons from high-confidence prior knowledge for accurate estimation of transcription factor activities". In: *Nucleic Acids Research* 51.20 (2023), pp. 10934–10949. ISSN: 0305-1048. DOI: 10.1093/nar/gkad841 (cit. on pp. 19–21).
- [61] B. Schwikowski, P. Uetz, and S. Fields. "A network of protein-protein interactions in yeast". In: *Nature Biotechnology* 18.12 (2000), pp. 1257–1261. ISSN: 1546-1696. DOI: 10.1038/82360 (cit. on p. 19).
- [62] J. Gillis and P. Pavlidis. "'Guilty by association' is the exception rather than the rule in gene networks". In: *PLoS Comput Biol* 8.3 (2012), e1002444. DOI: 10.1371/journal.pcbi.1002444 (cit. on p. 19).
- [63] L. Chindelevitch et al. "Assessing statistical significance in causal graphs". In: *BMC Bioinformatics* 13.1 (2012), p. 35. ISSN: 1471-2105. DOI: 10.1186/1471-2105-13-35 (cit. on pp. 20, 21, 26, 29).
- [64] X. Zhu, M. Gerstein, and M. Snyder. "Getting connected: analysis and principles of biological networks". In: *Genes Dev* 21.9 (2007), pp. 1010–24. DOI: 10.1101/gad.1528707 (cit. on p. 20).
- [65] M. Ziv et al. "The TissueNet v.3 Database: Protein-protein Interactions in Adult and Embryonic Human Tissue contexts". In: *J Mol Biol* 434.11 (2022), p. 167532. ISSN: 0022-2836. DOI: 10.1016/j.jmb.2022.167532 (cit. on pp. 20, 66, 67).
- [66] J. Pollard J. et al. "A computational model to define the molecular causes of type 2 diabetes mellitus". In: *Diabetes Technol Ther* 7.2 (2005), pp. 323–36. DOI: 10.1089/dia.2005.7.323 (cit. on pp. 20, 29).

- [67] D. Yu et al. "Review of biological network data and its applications". In: *Genomics Inform* 11.4 (2013), pp. 200–10. DOI: 10.5808/gi.2013.11.4.200 (cit. on p. 21).
- [68] D. Türei, T. Korcsmáros, and J. Saez-Rodriguez. "OmniPath: guidelines and gateway for literature-curated signaling pathway resources". In: *Nature Methods* 13.12 (2016), pp. 966–967. ISSN: 1548-7105. DOI: 10.1038/nmeth.4077 (cit. on pp. 21, 24, 37).
- [69] A. Valdeolivas et al. *OmnipathR: client for the OmniPath web service*. Bioconductor Package. 2019. DOI: 10.18129/B9.bioc.OmnipathR (cit. on pp. 21, 24).
- [70] Clarivate. *MetaBase, Accelerate drug discovery research*. Web Page. 5.1.0. URL: <https://clarivate.com/life-sciences-healthcare/research-development/discovery-development/early-research-intelligence-solutions/> (cit. on pp. 21, 24, 37).
- [71] M. Zitnik et al. *BioSNAP Datasets: Stanford Biomedical Network Dataset Collection*. Dataset. 2018. URL: <http://snap.stanford.edu/biodata> (cit. on p. 23).
- [72] S. Orchard et al. "The MIntAct project-IntAct as a common curation platform for 11 molecular interaction databases". In: *Nucleic Acids Research* 42.D1 (2013), pp. D358–D363. ISSN: 0305-1048. DOI: 10.1093/nar/gkt1115 (cit. on p. 23).
- [73] M. Kanehisa and S. Goto. "KEGG: kyoto encyclopedia of genes and genomes". In: *Nucleic Acids Research* 28.1 (2000), pp. 27–30. DOI: 10.1093/nar/28.1.27 (cit. on p. 23).
- [74] F. Ceccarelli et al. "Bringing data from curated pathway resources to Cytoscape with OmniPath". In: *Bioinformatics* 36.8 (2020), pp. 2632–2633. ISSN: 1367-4803 (cit. on p. 24).
- [75] A. Luna et al. "PaxtoolsR: pathway analysis in R using Pathway Commons". In: *Bioinformatics* 32.8 (2016), pp. 1262–4. DOI: 10.1093/bioinformatics/btv733 (cit. on p. 24).
- [76] E. G. Cerami et al. "Pathway Commons, a web resource for biological pathway data". In: *Nucleic acids research* 39 (2010), pp. D685–D690. ISSN: 0305-1048. DOI: 10.1093/nar/gkq1039 (cit. on p. 24).
- [77] A. Fabregat et al. "Reactome pathway analysis: a high-performance in-memory approach". In: *BMC Bioinformatics* 18.1 (2017), p. 142. ISSN: 1471-2105. DOI: 10.1186/s12859-017-1559-2 (cit. on p. 25).
- [78] G. Yu and Q. Y. He. "ReactomePA: an R/Bioconductor package for reactome pathway analysis and visualization". In: *Mol Biosyst* 12.2 (2016), pp. 477–9. ISSN: 1742-2051. DOI: 10.1039/c5mb00663e (cit. on p. 25).
- [79] J. Griss et al. "ReactomeGSA - Efficient Multi-Omics Comparative Pathway Analysis". In: *Mol Cell Proteomics* 19.12 (2020), pp. 2115–2125. DOI: 10.1074/mcp.TIR120.002155 (cit. on p. 25).

## REFERENCES

---

- [80] C. von Mering et al. "STRING: known and predicted protein-protein associations, integrated and transferred across organisms". In: *Nucleic Acids Research* 33.Database issue (2005), pp. D433–7. DOI: 10.1093/nar/gki005 (cit. on p. 25).
- [81] D. N. Slenter et al. "WikiPathways: a multifaceted pathway database bridging metabolomics to other omics research". In: *Nucleic Acids Research* 46.D1 (2018), pp. D661–D667. DOI: 10.1093/nar/gkx1064 (cit. on p. 25).
- [82] A. R. Pico et al. "WikiPathways: pathway editing for the people". In: *PLoS biology* 6.7 (2008), e184. ISSN: 1544-9173. DOI: 10.1371/journal.pbio.0060184 (cit. on p. 25).
- [83] J. Bajorath. "Molecular Similarity Concepts for Informatics Applications". In: *Bioinformatics: Volume II: Structure, Function, and Applications*. Ed. by J. M. Keith. New York, NY: Springer New York, 2017, pp. 231–245. ISBN: 978-1-4939-6613-4. DOI: 10.1007/978-1-4939-6613-4\_13 (cit. on p. 26).
- [84] V. Venafrà, F. Sacco, and L. Perfetto. "SignalingProfiler 2.0 a network-based approach to bridge multi-omics data to phenotypic hallmarks". In: *npj Systems Biology and Applications* 10.1 (2024), p. 95. ISSN: 2056-7189. DOI: 10.1038/s41540-024-00417-6 (cit. on pp. 26, 66).
- [85] Q. Wen et al. "Connectivity mapping using a combined gene signature from multiple colorectal cancer datasets identified candidate drugs including existing chemotherapies". In: *BMC Syst Biol* 9 Suppl 5.Suppl 5 (2015), S4. ISSN: 1752-0509. DOI: 10.1186/1752-0509-9-s5-s4 (cit. on p. 27).
- [86] K. Lin et al. "A comprehensive evaluation of connectivity methods for L1000 data". In: *Brief Bioinform* 21.6 (2020), pp. 2194–2205. DOI: 10.1093/bib/bbz129 (cit. on pp. 27, 40, 64).
- [87] A. Subramanian et al. "Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles". In: *Proceedings of the National Academy of Sciences* 102.43 (2005), pp. 15545–15550. ISSN: 0027-8424. DOI: 10.1073/pnas.0506580102 (cit. on p. 28).
- [88] E. Kiciman and A. Sharma. *Causal Reasoning: Fundamentals and Machine Learning Applications*. 2019. URL: <https://causalinference.gitlab.io/Causal-Reasoning-Fundamentals-and-Machine-Learning-Applications/> (cit. on p. 28).
- [89] S. Farahmand et al. "Causal Inference Engine: a platform for directional gene set enrichment analysis and inference of active transcriptional regulators". In: *Nucleic Acids Research* 47.22 (2019), pp. 11563–11573. DOI: 10.1093/nar/gkz1046 (cit. on pp. 29, 40, 45).
- [90] A. Krmer et al. "Causal analysis approaches in Ingenuity Pathway Analysis". In: *Bioinformatics* 30.4 (2013), pp. 523–530. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btt703 (cit. on p. 29).

- 
- [91] L. Chindelevitch et al. "Causal reasoning on biological networks: interpreting transcriptional changes". In: *Bioinformatics* 28.8 (2012), pp. 1114–21. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/bts090 (cit. on pp. 29, 40).
  - [92] C. T. Fakhry et al. "Interpreting transcriptional changes using causal graphs: new methods and their practical utility on public networks". In: *BMC Bioinformatics* 17.1 (2016), p. 318. ISSN: 1471-2105. DOI: 10.1186/s12859-016-1181-8 (cit. on pp. 29, 40).
  - [93] S. Mangul et al. "Systematic benchmarking of omics computational tools". In: *Nature Communications* 10.1 (2019), p. 1393. ISSN: 2041-1723. DOI: 10.1038/s41467-019-09406-4 (cit. on pp. 30, 31, 67).
  - [94] L. M. Weber et al. "Essential guidelines for computational method benchmarking". In: *Genome Biology* 20.1 (2019), p. 125. ISSN: 1474-760X. DOI: 10.1186/s13059-019-1738-8 (cit. on pp. 30, 63).
  - [95] R. D. C. Team. *R: A Language and Environment for Statistical Computing*. Computer Program. 2024. URL: <https://www.R-project.org/> (cit. on p. 32).
  - [96] A. Liu et al. "From expression footprints to causal pathways: contextualizing large signaling networks with CARNIVAL". In: *npj Systems Biology and Applications* 5.1 (2019), p. 40. ISSN: 2056-7189. DOI: 10.1038/s41540-019-0118-z (cit. on pp. 40, 41).
  - [97] H. Noh, J. E. Shoemaker, and R. Gunawan. "Network perturbation analysis of gene transcriptional profiles reveals protein targets and mechanism of action of drugs and influenza A viral infection". In: *Nucleic Acids Research* 46.6 (2018), e34. DOI: 10.1093/nar/gkx1314 (cit. on pp. 40, 43).
  - [98] R. Browaeys, W. Saelens, and Y. Saeys. "NicheNet: modeling intercellular communication by linking ligands to target genes". In: *Nature Methods* 17.2 (2020), pp. 159–162. ISSN: 1548-7105. DOI: 10.1038/s41592-019-0667-5 (cit. on pp. 40, 46, 47).
  - [99] S. Jaeger et al. "Causal Network Models for Predicting Compound Targets and Driving Pathways in Cancer". In: *J Biomol Screen* 19.5 (2014), pp. 791–802. ISSN: 1087-0571. DOI: 10.1177/1087057114522690 (cit. on p. 40).
  - [100] O. Vanunu et al. "Associating genes and protein complexes with disease via network propagation". In: *PLoS Comput Biol* 6.1 (2010), e1000641. DOI: 10.1371/journal.pcbi.1000641 (cit. on p. 40).
  - [101] S. Köhler et al. "Walking the interactome for prioritization of candidate disease genes". In: *Am J Hum Genet* 82.4 (2008), pp. 949–58. DOI: 10.1016/j.ajhg.2008.02.013 (cit. on p. 40).

## REFERENCES

---

- [102] Y. Nikolsky et al. "Genome-wide functional synergy between amplified and mutated genes in human breast cancer". In: *Cancer Res* 68.22 (2008), pp. 9532–40. ISSN: 0008-5472. DOI: 10.1158/0008-5472.Can-08-3082 (cit. on p. 40).
- [103] Z. Dezso et al. "Identifying disease-specific genes based on their topological significance in protein networks." In: *BMC Syst Biol* 3 (2009), p. 36. DOI: 10.1186/1752-0509-3-36 (cit. on p. 40).
- [104] C. L. Hsu et al. "Prioritizing disease candidate genes by a gene interconnectedness-based approach". In: *BMC Genomics* 12 Suppl 3.Suppl 3 (2011), S25. ISSN: 1471-2164. DOI: 10.1186/1471-2164-12-s3-s25 (cit. on p. 40).
- [105] S. D. Zhang and T. W. Gant. "A simple and robust method for connecting small-molecule drugs using gene-expression signatures". In: *BMC Bioinformatics* 9 (2008), p. 258. DOI: 10.1186/1471-2105-9-258 (cit. on p. 41).
- [106] R. L. Davis. "Mechanism of Action and Target Identification: A Matter of Timing in Drug Discovery". In: *iScience* 23.9 (2020), p. 101487. ISSN: 2589-0042. DOI: 10.1016/j.isci.2020.101487 (cit. on p. 63).
- [107] J. Paik et al. *Mechanisms of Action (MoA) Prediction*. Web Page. 2020. URL: <https://kaggle.com/competitions/lish-moa> (cit. on p. 63).
- [108] I. Mellis et al. "Prevalence of and gene regulatory constraints on transcriptional adaptation in single cells." In: *Genome biology* 25 (2024), pp. 347–347. DOI: 10.1186/s13059-024-03351-2 (cit. on p. 65).
- [109] R. Shukla et al. "Signature-based approaches for informed drug repurposing: targeting CNS disorders". In: *Neuropsychopharmacology* 46.1 (2021), pp. 116–130. DOI: 10.1038/s41386-020-0752-6 (cit. on p. 66).
- [110] A. D. Matlock et al. "NeuroLINCS Proteomics: Defining human-derived iPSC proteomes and protein signatures of pluripotency". In: *Sci Data* 10.1 (2023), p. 24. DOI: 10.1038/s41597-022-01687-7 (cit. on p. 66).
- [111] P. J. Thul and C. Lindskog. "The human protein atlas: A spatial map of the human proteome". In: *Protein Sci* 27.1 (2018), pp. 233–244. DOI: 10.1002/pro.3307 (cit. on p. 66).
- [112] A. Dugourd et al. "Causal integration of multi-omics data with prior knowledge to generate mechanistic hypotheses". In: *Mol Syst Biol* 17.1 (2021), e9730. DOI: 10.15252/msb.20209730 (cit. on p. 66).
- [113] M. A. Trapotsi et al. "Comparison of Chemical Structure and Cell Morphology Information for Multitask Bioactivity Predictions". In: *J Chem Inf Model* 61.3 (2021), pp. 1444–1456. ISSN: 1549-9596. DOI: 10.1021/acs.jcim.0c00864 (cit. on p. 67).
- [114] D. Rivas-Barragan et al. "Drug2ways: Reasoning over causal paths in biological networks for drug discovery". In: *PLoS Comput Biol* 16.12 (2020), e1008464. DOI: 10.1371/journal.pcbi.1008464 (cit. on p. 67).

- [115] Z. P. Wang et al. "Drug repurposing for Alzheimer's disease using a graph-of-thoughts based large language model to infer drug-disease relationships in a comprehensive knowledge graph". In: *BioData Mining* 18.1 (2025), p. 51. ISSN: 1756-0381. DOI: 10.1186/s13040-025-00466-5 (cit. on p. 68).

## SUPPLEMENTARY TABLES

Table I.1: Summary of evaluated runs. Tool category: C = Connectivity mapping tools, N = Network tools, E = Enrichment tools; Query dataset characteristics: D = Drug perturbation, G = Gene perturbation, C = Cell lines, T = Tissues/in vivo models, Sc = Single-cell data. \*: Selected impact signatures from 7390 signatures, \*\*: Generated consensus signatures from 82256 signatures.

Query				Reference								
Database	Perturbation	Origin	Signatures	MetaBase	OmniPath	MetaBase (Linear Path)	MetaBase (Regulons)	OmniPath (Regulons)	LINCS CRISPR	LINCS OE	LINCS shRNA	GWPS Perturb-Seq
CDS-DB	D	T	181	N	N	NE	NE	NE	CNE	CNE	CNE	CNE
Sci-Plex	D	C (Sc)	405	N	N	NE	NE	NE	CNE	CNE	CNE	CNE
PertOrg	G	T	951*	N	N	NE	NE	NE	CNE	CNE	CNE	CNE
ChemPert	D	C	1,304**	N	N	N	N	N	CN	CN	CN	CN
GWPS	G	C (Sc)	1,980	N	N	NE	NE	NE	CNE	CNE	CNE	-
CREEDS	DG	T	2,642	N	N	N	N	N	CN	CN	CN	CN
LINCS Compounds	D	C	3,540	N	N	NE	NE	NE	CNE	CNE	CNE	CNE
LINCS OE	G	C	3,780	N	N	NE	NE	NE	CNE	-	CNE	CNE
LINCS shRNA	G	C	4,854	N	N	NE	NE	NE	CNE	CNE	-	CNE
LINCS CRISPR	G	C	5,156	N	N	NE	NE	NE	-	CNE	CNE	CNE

Table I.2: Algorithm execution success and run rate across all benchmarking runs. For each of the 24 algorithms evaluated, the number of successful runs, failures, total runs attempted, and runs skipped are reported. Success rate represents the proportion of successful executions among attempted runs (Number of success / Number of run). The run rate is the proportion of time where the algorithm was launched (Total number of runs / (Total number of runs + Number of skipped runs)). Algorithms are ordered by decreasing run rate.

Algorithm	Number success	Number fail	Number run	Number skipped	Success rate	Run rate
causalReasoning	84	2	86	0	0.98	1.00
randomWalk	84	2	86	0	0.98	1.00
ulm	51	1	52	0	0.98	1.00
overconnectivity	83	0	83	3	1.00	0.97
wsum	43	7	50	2	0.86	0.96
wmean	43	8	51	1	0.84	0.98
SigNet	60	18	78	8	0.77	0.91
KS	23	0	23	13	1.00	0.64
GSEAwight0	22	1	23	13	0.96	0.64
GSEAwight1	22	1	23	13	0.96	0.64
XSum	22	1	23	13	0.96	0.64
Zhang	22	1	23	13	0.96	0.64
udt	30	1	31	21	0.97	0.60
interconnectivity	47	0	47	39	1.00	0.55
networkPropagation	47	0	47	39	1.00	0.55
mdt	29	2	31	21	0.94	0.60
XCos	20	3	23	13	0.87	0.64
viper	20	4	24	28	0.83	0.46
fgsea	13	3	16	36	0.81	0.31
NicheNet	2	0	2	84	1.00	0.02
mlm	8	6	14	38	0.57	0.27
CARNIVAL	0	10	10	76	0.00	0.12
CIE	0	9	9	77	0.00	0.10
ProTINA	0	0	0	66	-	0.00





**NOVA**

UNIVERSIDADE NOVA  
DE LISBOA

2025

BENCHMARKING COMPUTATIONAL ALGORITHMS FOR ENHANCED DRUG DISCOVERY

**Maria Inês Nunes Vilar Gomes**

MASTER IN **Computational Biology and Bioinformatics**

**SPECIALIZATION** Multi-Omics for Life and Health Sciences