# Discovering Causal Relations and Equations from Data

## A review of methods, challenges and opportunities

Gustau Camps-Valls, Andreas Gerhardus, Urmi Ninad, Gherardo Varando, Georg Martius, Ricardo Vinuesa, Emili Balaguer-Ballester, Emiliano Diaz, Laure Zanna, Jakob Runge

May 24, 2023

*First version, May 24, 2023*

# Discovering Causal Relations and Equations from Data

Gustau Camps-Valls[1,†], Andreas Gerhardus[2,*], Urmi Ninad[3,*], Gherardo Varando[1,*], Georg Martius[4,5], Emili Balaguer-Ballester[6,7], Ricardo Vinuesa[8], Emiliano Diaz[1], Laure Zanna[9], Jakob Runge[2,3]

[1]Universitat de València, València, Spain

[2]German Aerospace Center, Jena, Germany

[3]Technische Universität Berlin, Berlin, Germany

[4]University of Tübingen, Tübingen, Germany

[5]Max Planck Institute for Intelligent Systems, Tübingen, Germany

[6]Bournemouth University, Bournemouth, UK

[7]Medical Faculty Mannheim and Heidelberg University, Mannheim, Germany.

[8]KTH Royal Institute of Technology, Stockholm, Sweden

[9]New York University, New York, USA

[*] These authors contributed equally.

[†] Corresponding author at Image Processing Laboratory (IPL), E4 building - 4th floor, Parc Científic Universitat de València. C/ Cat. Agustín Escardino Benlloch, 9. 46980 Paterna (València). Spain. E-mail address: gustau.camps@uv.es

# Contents

# Preface

*"As in Mathematics, so in Natural Philosophy, the Investigation of difficult Things by the Method of Analysis ought ever to precede the Method of Composition. This Analysis consists in making Experiments and Observations, and in drawing general Conclusions from them by Induction, and admitting of no Objections against the Conclusions, but such as are taken from Experiments, or other certain Truths ... By this way of Analysis, we may proceed from Compounds to Ingredients, and from Motions to the Forces producing them; and in general, from Effects to their Causes, and particular Causes to more general ones, till the Argument end in the most general. This is the Method of Analysis"* (Newton, 1718).

Physics is a field of science that has traditionally used the scientific method to answer questions about why natural phenomena occur and to make testable models that explain the phenomena. Discovering equations, laws and principles that are invariant, robust and causal explanations of the world has been fundamental in physical sciences throughout the centuries. Discoveries emerge from observing the world and, when possible, performing interventional studies in the system under study. With the advent of big data and the use of data-driven methods, causal and equation discovery fields have grown and made progress in computer science, physics, statistics, philosophy, and many applied fields. All these domains are intertwined and can be used to discover causal relations, physical laws, and equations from observational data. This paper reviews the concepts, methods, and relevant works on causal and equation discovery in the broad field of Physics and outlines the most important challenges and promising future lines of research. We also provide a taxonomy for observational causal and equation discovery, point out connections, and showcase a complete set of case studies in Earth and climate sciences, fluid dynamics and mechanics, and the neurosciences. This review demonstrates that discovering fundamental laws and causal relations by observing natural phenomena is being revolutionised with the efficient exploitation of observational data, modern machine learning algorithms and the interaction with domain knowledge. Exciting times are ahead with many challenges and opportunities to improve our understanding of complex systems.

# 1. Introduction

This paper reviews the recent advances in *causal discovery* and *equation discovery* from data. Both problems are conundrums for scientists and philosophers of science. After all, Science is about studying, discovering, and understanding the structure and behaviour of the physical and natural world through observation and experimentation. Understanding the system's structure involves performing interventions on the systems to evaluate their responses. However, interventional experiments are often not feasible for economic or ethical reasons, so relying on observations, simulations, and domain knowledge must be exploited. In recent decades, discovering causal relations and underlying governing laws from data have emerged as exciting fields of research that promise advancing science.

## 1.1 Understanding in the physical sciences

A pertinent question arises here; what is understanding? Understanding is the ability to comprehend and make sense of processes to gain a deeper knowledge of a system. Understanding involves analysing information, making (eventually causal) connections, and coming to conclusions. Whether the conclusions should be falsifiable has been the subject of active discussion in the Philosophy of Science [Popper, 2005, Munz, 2014]. We aim to *understand* complex systems by following the *scientific method*; make an observation, ask a scientific question, form a hypothesis, theory, model, or explanation of the phenomena, and make predictions, which are ultimately tested and whose results are used to make new hypotheses or predictions (see Fig. 1.1).

Understanding involves reasoning and thinking critically about a subject, a system's behaviour, or a problem. Understanding and explaining how a system works is more complicated than making predictions about the system's behaviour. The Oracle of Delphi gave accurate predictions of the future and the optimal course of action, but the lack of understanding frequently led to disaster. True understanding is about making (truly accurate) predictions and, more importantly, gaining knowledge of the causal chain. Science generally aims to answer causal questions, infer causal relations, and attain mathematical models (mainly laws and equations) that work well in most situations, explain the system and underlying processes, and are invariant across space and time. Without it, we cannot predict the consequences of our actions (*interventions*) or analyse when,

Figure 1.1: Standard loop in understanding complex systems following the standard scientific method. Understanding involves experimentation by refining a descriptive mechanistic model. The initial model hypothesis is tested in practice and, through experiments, yields observations that are confronted with the model's predictions. The unexplained processes are then used to improve the model's misspecification and predictions.

where and why things went wrong (*counterfactuals*) [Pearl, 2009c].



Figure 1.2: The Oracle of Delphi at inference time.

Yet, what do we want to understand? And how do we generally do it? In physical sciences, one typically analyses phenomena and instantiations of the physical world, uses observations, and refines and tests models. For learning about the system, one aims to (1) characterise its complexity in terms of trajectories, persistence, stability and collapse, bifurcations and viability boundaries [Sterman, 1994, Kwapień and Drożdż, 2012, Salcedo-Sanz et al., 2022, May, 1972], (2) obtain explanatory and causal models of their behaviour [Pearl, 2009a, Peters et al., 2017a, Runge et al., 2019a], and (3) discover and formalise general laws, governing equations, and parameterizations [Richardson, 1996, Brunton et al., 2016b, Chen et al., 2022]. These three components allow us to advance science and technology. However, in many systems, governing equations and causal relations are (partially) unknown, and recourse to first principles is untenable. Resorting to algorithms that can discover laws, governing equations, and causal relations from data may thus constitute a paradigm shift that promises to accelerate science.

## 1.2 Scientific discovery

Before revising the fields, let us take a step back and review the formalism in the logic of science and the key elements, definitions, and challenges we face in addressing such problems. What are the cornerstones for understanding how science and the scientific method work? The philosophy of science studies science's assumptions, foundations, and implications. It examines the implications of scientific theories and methods for understanding the world. It concerns the fundamental questions at the heart of science, such as the nature of knowledge, the limits of scientific inquiry, and the relationship between science and values.

Scientific discovery is the process or product of successful scientific inquiry. Objects of discovery can be things, events, processes, causes, properties, theories, hypotheses, and their characteristics. Philosophical discussions of scientific discovery vary widely in scope and definition, from the narrowest sense of a "eureka moment" to the broadest sense of a "successful scientific endeavour". The utilisation of data sets to create and test new hypotheses in philosophical discourse

has led to a multifaceted and intricate discussion regarding the precise definition and potential misuse of the term "discovery".



Figure 1.3: Archimedes (287-212 BC) in his bath. A woodcut by Flörtner (1490-1546).

Human nature aims to discover. Always. Since the Bronze Age[1]. Generations have created and discovered new principles, techniques, and operations through millennia. Right after the Neolithic Revolution, the world stopped except for some remarkable technological advances, like the invention of the water wheel (476-221 BC) and the windmill (ca 644 BC). Romans were amazed by stories of what Archimedes (287-212 BC) had been able to do. But, bold as it may sound, one may claim that modern science was invented between 1572 when Tycho Brahe saw a *nova* or new star, and 1704 when Newton published his Opticks [Wootton, 2015]. What happened in that period prepared humanity and scientists for a New Era: a research program endorsed with a scientific method that allowed scrutinising new theories and validating or refuting hypotheses and models of the world, and all that in the light of evidence and observations. After fitting many ovoids to observational data, Kepler discovered the laws of planetary motion (1609) and needed four years to discover Mars' orbit was an ellipse. The scientific method was slow but sure. Galilei discovered the law of falling bodies (1638) by dropping two cannonballs of different masses from the tower of Pisa and measuring the effect of mass on the fall rate to the ground. And in 1662, Boyle discovered the law of ideal gases. Only ten years later, in 1672, Newton discovered that white light is a mixture of distinct coloured rays, and in 1687 he formulated the classical mathematical description of the fundamental force of universal gravitation and the three physical laws of motion. The triumph of Newtonianism marks the end of the beginning of scientific discovery.

The history of science is a long and complex narrative punctuated by moments of major scientific revolutions [Kuhn, 1962]. Kuhn identified a general pattern: A discovery is not a simple act but an extended, complex process that culminates in paradigm changes. These events mark a fundamental shift in how science is understood and practised and have greatly impacted scientific progress. The first scientific revolution occurred in the 16th and 17th centuries when the Copernican revolution challenged the traditional Ptolemaic view of the universe [Copernicus et al., 1965]. The next major scientific revolution was the Enlightenment during the 18th century, which saw the emergence of the scientific method and the development of modern physics and chemistry, from formulating the laws of motion to discovering electricity. This period also saw the emergence of scientific societies, which helped to propagate and popularise scientific ideas. The 19th century saw the emergence of the theory of evolution, which revolutionised the field of biology [Darwin, 1859]. This revolution was followed by the rise of modern genetics, which further expanded our understanding of the evolution of life [Watson and Crick, 1953]. The 20th century saw the emergence of the quantum revolution, which revolutionised our understanding of the physical world as it could not be fully explained by classical physics [Heisenberg, 1925]. Revolutions happen only gradually, as it takes time for the scientific community to recognise *"both that something is and what it is"* [Kuhn, 1962]. Eventually, a new paradigm becomes established, and the strange phenomena become the expected phenomena.

The idea that there is such a thing as 'the Scientific Revolution' and that it took place in the 17th century is thus a fairly recent one [Butterfield, 1965]; some have argued that it can be seen as the construction of intellectuals looking back from the 20th century [Wootton, 2015]. Like the term 'Industrial Revolution', the idea of a scientific revolution brings problems of multiplication (how many scientific revolutions?) and periodisation (how often?). Some philosophers of science

---

[1] https://en.wikipedia.org/wiki/Timeline_of_scientific_discoveries

have argued for continuity, others have sought multiple revolutions: the Darwinian revolution, the Quantum revolution, the DNA revolution, and so on, while others claim that the real Scientific Revolution came in the 19th century when science and technology married. Recently, we are witnessing the so-called 4th Industrial Revolution, which conceptualises rapid change to technology, industries, and societal patterns due to increasing interconnectivity, smart automation and the amalgamation of artificial intelligence and automated machines. Yet, is it only a technological revolution, or can machines discover and explain new science? Are we facing the emergence of machine discovery of science?



Figure 1.4: Left: The title page of Johannes Stradanus's New Discoveries ('*Nova reperta*', c.1591) summarises the knowledge that marks off the modern world from the ancient: the discovery of America, the invention of the compass, the printing press, the gunpowder, the clock, silk weaving, distillation and the saddle with stirrups. Right: Can AI start a new scientific revolution? Is AI itself the scientific revolution?

### 1.2.1 Elements of scientific discovery

During the 18th and 19th centuries, the different elements of discovery gradually became separated and discussed in more detail. Discussions concerned the nature of observations and experiments, the act of having insight and the processes of articulating, developing, and testing the novel insight. For Whewell, for example, discovery comprised three elements: the happy thought, the articulation and development of that thought, and the testing or verifying it. In contrast to many 20th-century approaches, Whewell's philosophical conception of discovery articulates happy thoughts and integrates the process of verification as an integral part of discovery. Thus, to verify a hypothesis, the scientist needs to show that it accounts for the known facts, that it foretells new, previously unobserved phenomena, and that it can explain and predict phenomena which are explained and predicted by a hypothesis that was obtained through an independent happy thought-cum-colligation [Ducasse, 1951]. Until the late 20th century, most philosophers operated with a narrower notion of discovery than Whewell's. Controversies in the 20th century moved around whether or not the discovery process should include the articulation and development of a novel thought.

### 1.2.2 Elements of knowledge

Throughout history, many scientific accounts have described methods of knowledge generation and scientific reasoning without explicitly labelling them as such. These methods include using the senses to gather knowledge, observation and experimentation, analysis and synthesis, induction and deduction, hypotheses, probability, and certainty (Table 1.1). By exploring these methods, scientists have generated new knowledge and developed theories about the nature of matter and natural forces. These methods are integral to scientific inquiry and continue to be used today.

In the early modern period, authors such as Bacon and Newton advanced ideas about generating and verifying empirical knowledge and the difficulties that may arise in scientific inquiry. These theories were closely linked to theories about matter and force. However, by the 18th and 19th

centuries, authors on scientific method and logic began citing early modern approaches to model proper scientific practice and reasoning. At the same time, the connection between the two was gradually severed. It was common in 20th-century philosophy of science to draw a sharp contrast between those early theories of the scientific method and modern approaches. Yet, recent research in the history of the philosophy of science has shown that the development of the scientific method was a gradual and ongoing process rather than a sudden shift from one approach to another.

Yet, how can we build knowledge? *Knowledge* is a well-studied, yet elusive, concept in the philosophy of science [Popper, 2005]. Knowledge can be acquired through observation, experimentation, and reasoning and is often expressed through language, symbols, and concepts. In quantitative physical sciences, we encapsulate our domain knowledge in mathematical constructs and equations that describe the system under consideration. Knowledge is built on *facts*, that is, *observations* about the physical world that are accepted as true and are often the starting point for deeper analysis and understanding. Empiricism offers resorts to perform *experiments* that are used to test the validity of facts and to gain further insight into the underlying *causes and effects* of phenomena. From there, scientists construct *laws*, which are generalisations based on many observations and experiments, the rules describing the behaviour of a system that can be used to predict future outcomes. In the scientific method, *hypotheses* are testable explanations for facts and laws based on *evidence* but have not been proven true. See Table 1.1.

We must use judgement to evaluate the validity of facts, experiments, laws, hypotheses, theories, and evidence. *Judgement* involves forming an opinion or deciding based on evidence and knowledge through reasoning and considering our conclusions' potential implications. Judgement is important in all aspects of knowledge, from generating new ideas to applying existing knowledge. It is important to remember that knowledge is not static but an ongoing discovery, analysis, and interpretation process. New theories and hypotheses can be formed as new evidence is discovered, and existing knowledge can be refined or overturned. Formalising new knowledge needs *arguments*, which can be inductive, deductive, analogical, and statistical. *Inductive* arguments are based on observations and generalisations, while *deductive* arguments are based on logical reasoning. *Analogical* arguments are based on similarities between two or more cases, while *statistical* arguments are based on data and evidence and formally linked to the probability theory.

**Objectivity and subjective conviction.**

Objectivity is the idea that scientific inquiry should be based on factual evidence and logical reasoning, independent of personal bias or opinion. Objectivity, as a (causal) statement, should be inter-subjectively tested according to Popper [Popper, 2005] or be valid, justifiable or verifiable in the sense of Dickerson [Dickerson, 2003]. Objectivity can be challenged because some scientific theories can never be fully justifiable or verifiable but testable (i.e. reproducibility). On the contrary, conviction relies on believing the inquiry's findings are true and trustworthy, stating a scientific argument without proofs [Popper, 2005]. Conviction implies arguments based on mere association, which are not necessarily causal, verifiable, testable or reproducible and cannot justify a scientific statement [Pearl, 2009c].

**Universality, falsifiability, and consistency.**

The concepts of universality, falsifiability, and consistency are essential tools for understanding scientific method [Popper, 2005] (Table 1.1). *Universality*, the principle of *universalism*, is the idea that scientific truths can be applied to all situations. It is a fundamental tenet of scientific inquiry, implying that the same principles can be used to explain phenomena in different contexts. This allows scientists to conclude the natural world, not limited to a single instance or context. It also makes experiments more reliable and reproducible, as results can be compared across multiple contexts. The latter relates to *falsifiability* introduced by Popper [Popper, 2005]: a scientific hypothesis must be tested and disproved. Falsifiability is important for scientific research, as it

allows for the testing and refining of hypotheses. The concept of falsifiability is used to evaluate the validity of scientific claims, as it implies that a hypothesis can be disproved if it does not fit the evidence. Yet, scientific theories must remain *consistent* over time; any change or modification to a scientific theory must be able to be explained by existing evidence. Newton first proposed this concept in the 17th century [Newton, 1833]. All three concepts remain an important part of the scientific method.

**Empiricism.**

Empiricism is a philosophical position that asserts that knowledge must be acquired through sensory experience [Sellars et al., 1956]. This means knowledge is not innate and can only be gained through observation and experimentation. Empirical evidence is necessary for scientific progress as the only reliable way of assessing the validity of theories and hypotheses. Empiricism is closely associated with the scientific method, as it gathers data and evidence to understand the world comprehensively. This data then form hypotheses and theories explaining why certain phenomena occur worldwide. Empiricism has been challenged most notably by Feyerabend [Feyerabend, 1981], where realism is promoted to enable the proliferation of new and incompatible theories. This way, scientific progress comes through "theoretical pluralism" allowing a plurality of incompatible theories, each of which will contribute by competition to maintaining and enhancing the testability, and thus the empirical content, of the others.

### 1.2.3   Elements of models, governing equations and laws

**Identifiability and equifinality.**

A relevant problem in model development is identifiability, the property by which learning the true values of a model's underlying parameters is theoretically possible from an infinite number of observations [Ljung and Glad, 1994, Peters et al., 2011]. The impossibility arises when two more model parameterisations are observationally equivalent, i.e., indistinguishability. This situation is sometimes called equifinality; one can achieve the same result or state description by many potential solutions and model parameterisations. This is a very active field of research in statistics, machine learning and functional analysis, where regularisation helps. Combining domain knowledge with data-driven approaches, for example, may alleviate cases of model misspecification, as the model solution would be confined to a solution subspace with little expressive power and even implausible solutions.

**Compressibility, sparsity and compositionally.**

A desirable property of models and governing equations is compositionality, by which theory/models are typically a composition of a small set of elementary functions [Udrescu and Tegmark, 2020]. The problem of attaining compositional models is challenging, as misspecification arises when defining model components or facing multidimensional compositions. Compositionality also advocates for sparse models, where combinations of simpler explanations are preferred, directly implementing Occam's razor in the model search. Besides, intimately related, we find the desirable property of compressibility by which data/models are compressible if and only if they exhibit a pattern: this way, high-level models are much simpler than their low-level counterparts [Dennett, 1991, Shannon, 1948, Kolmogorov, 1963]. As for attaining compositional models or achieving the model's sparsity, estimating multivariate information-theoretic estimates is challenging here.

**Generalisation, robustness, invariance and extrapolation.**

Models should operate well in unseen but similar situations. The generalisation property (or generalizability) gives confidence in the model's performance, expressive power, and predictions. Models operating in out-of-the-sample regimes tend to extrapolate, thus compromising trustworthiness. Extrapolation is about inferring the unknown from the known, estimating beyond the original

observation range, and predicting future data by relying on historical data [Scott Armstrong and Collopy, 1993]. Since the governing equations and laws of Physics are invariant through space and time, inferred models should be. However, attaining robust models that generalise well and implement properties of invariance is typically challenged by the lack of reliability of the data source, the accuracy of the extrapolation process itself, the complexity of the data, and the potential for errors in the extrapolation.

## 1.3 Knowledge discovery from data

### 1.3.1 Discoverability and heuristic strategies

The questions of what and how phenomena and mechanisms can be discovered have been the subject of intense research and philosophical discussion. In the philosophy of science, *discoverability* is the concept that scientific knowledge must be discoverable and verifiable [Ducasse, 1951, Popper, 2005, Langley, 2019, Langley et al., 1987b, Klahr and Simon, 1999]. This means that hypotheses or theories must be supported by evidence and based on empirical observations and data. Furthermore, scientific knowledge must be available for anyone to see and verify, ensuring that it is not biased or limited to a specific group of people.

Recent advancements in the philosophy of science have seen a revival of interest in *heuristic strategies* to discover knowledge; these strategies are seen as problem-solving activities, whereby a discovery is a solution to a problem. Heuristics-based discovery methodologies are neither completely subjective and intuitive nor algorithmic or formalisable. This view has shifted the scientific researcher from being viewed as a 'puzzle solver' to a 'problem solver' and 'decision maker' in complex, variable, and changing environments [Wimsatt and Wimsatt, 2007]. In this paper, we will review mathematical models that address equation discovery and causal discovery by generally formalising the problems as concrete statistical inference tasks; regression, conditional dependence or density estimation.

### 1.3.2 Modern approaches to data-driven discovery

Observational discovery relies on modelling. Yet, what types of models? Table 1.2, cf. [Peters et al., 2017b], gives a simple categorisation of models from mechanistic/physical models based on first principles and (rigid) equations and laws but with desirable properties of interpretability, invariance and robustness to distribution shifts, to purely statistical (machine learning) models that excel in prediction and are learned from data. In the middle, we have structural causal models, which can answer counterfactual questions but do not necessarily capture physical knowledge. All three models are used in quantitative data-driven science and map to different levels of discovery: learning statistical associations in data streams, identifying causal relations between variables, and discovering equations from data. Note the resemblance to the Ladder of Causation proposed by Pearl and Mackenzie [2018] with three rungs: association, intervention, and counterfactual. Discovering causal and physical laws from observations is a paradigm shift in AI and can impact the physical sciences and other disciplines. The fields of discovery of scientific knowledge and causal models of scientific phenomena are intertwined and tightly connected: our scientific endeavour is constantly challenged with causal questions, robust model building, intervention analysis, and hypothesis testing. The fields also share important theoretical challenges, where generalizability, compressibility, robustness, invariance, and extrapolation come into play.

#### Level 1 – Learning statistical associations.

The most rudimentary approach to building association links from multivariate time series data involves computing pairwise Pearson's correlations or mutual information, which capture relationships between variables at lag zero. Networks derived from these measures have found applications

Table 1.1: Key concepts, definitions, and challenges for understanding complex systems from data.

| Concept | Definition | Challenge in causal/equation discovery |
|---|---|---|
| *Logic of scientific discovery* | | |
| Consistency | theory not leading to logical contradictions; refutes/includes/generalises previous theories; (Gödel, 1931, Tarski and Tarski, 1994, Popper, 2005) | impossibility to perform controlled experiments repeated observations, or survey research; cf. hypothesis testing |
| Deduction | type of inference where the conclusion follows logically from its premises (Johnson-Laird and Byrne, 1991, Popper, 2005) | deduction breaks under confounders |
| Empiricism | principle that only `experience' can decide about the truth or falsity of a factual statement (Popper, 2005); central role of empirical evidence to form ideas, rather than innate ideas or traditions; cf. induction | impossibility to perform experiments or access observations |
| Explanation | deduction of a statement that describes the event using as premises one or more universal laws, together with certain singular statements (initial conditions) (Popper, 2005, Pearl, 2009c) | no formal principle of causality or universal causation available |
| Falsifiability | capacity for some proposition, statement, theory, or hypothesis to be proven wrong (Popper, 2005) | discovered models and laws can be misspecified, untestable or incomplete |
| Induction | method of reasoning where a general principle is derived from a set of observations (Popper, 2005, Munz, 2014, Reichenbach, 1991, Keynes, 2013), cf. generalisation (sample to population), probabilistic inference | impossibility, data scarcity or conditions; unlike in deduction, the truth of the conclusion is only probable |
| Objectivity | (causal) statement that can be inter-subjectively tested (Popper, 2005) or that is valid, justifiable or verifiable (Dickerson, 2003) | elusive term, scientific theories can never be fully justifiable or verifiable, but testable, cf. reproducibility |
| Conviction | stating a scientific argument without proofs (Popper, 2005) | mere association is not causal, verifiable, testable or reproducible and cannot justify a scientific statement |
| Universality | universal facts exist and can be progressively discovered, as opposed to relativism (Popper, 2005) | model misspecification, reproducibility, theory testability; Universal laws transcend any finite number of their observable instances |
| *Equation discovery: model properties and data challenges* | | |
| Compositionality | theory models are typically a composition of a small set of elementary functions (Udrescu and Tegmark, 2020) | misspecification, multidimensional compositions |
| Compressibility | data/models are compressible iff exhibit a pattern, high-level models are much simpler than their low-level counterparts (Dennett, 1991, Shannon, 1948, Kolmogorov, 1963) | multivariate information-theoretic estimates are difficult, cf. sparsity |
| Extrapolation | inferring the unknown from the known, estimation beyond the original observation range, predicting future data by relying on historical data (Scott Armstrong and Collopy, 1993) | reliability of the data source, the accuracy of the extrapolation process itself, the complexity of the data, and the potential for errors in the extrapolation |
| Generalization | model's ability to adapt to new, unseen data drawn from the same/similar distribution as the one used to create the model (Vapnik, 1999) | hard to evaluate, pointless in regime/domain shift, cf. extrapolation |
| Identifiability | learning the true values of a model's underlying parameters is theoretically possible from an infinite number of observations (Ljung and Glad, 1994, Peters et al., 2011) | two or more model parameterisations are observationally equivalent, cf. indistinguishability, equifinality |
| Invariance | probability of a certain event occurring remains the same, regardless of any changes to the underlying system or process (Olver, 1999, Kallenberg, 2005) | strong assumption in many complex systems |
| Robustness | property of models to perform well in noisy environments (presence of outliers, mixed noise distributions, missing data), extrapolation regimes, covariate shift | hard-to-identify regimes, design model's regularisers, cf. inductive bias, priors, Bayesian inference |
| Separability | model can be written as a sum or product of two parts with no variables in common (Udrescu and Tegmark, 2020) | misspecification, oversimplification |
| Smoothness | model is continuous and perhaps even analytic in its domain (Schölkopf et al., 2021) | oversimplification, inappropriate to deal with regime shifts and anomalies |
| Symmetry | model exhibits translational, rotational, or scaling symmetry concerning some of its variables (Udrescu and Tegmark, 2020) | inappropriate for unstructured data, macroscale-vs-microscale systems |
| Uncertainty | epistemic situations involving imperfect or unknown information; it applies to predictions of future events, to physical measurements that are already made, or to the unknown (Lindley, 2013, Keynes, 2013) | partially observable or stochastic environments, yet difficult to disentangle the sources and to quantify it properly |
| Units | models & corresponding variables have known physical units (Udrescu and Tegmark, 2020) | dealing with unknown factors and confounders |
| *Causal inference: discovery, cause-effect estimation & counterfactuals* | | |
| Causal discovery | Reconstructing the full causal graph or pairwise causation from data and assumptions (Pearl and Mackenzie, 2018, Peters et al., 2017a) | challenging, limited data, poor assumptions, incomplete knowledge, untestable |
| Causal effects | Estimating causal effects (of hypothetical interventions) in terms of expectations or full interventional distributions (Winship and Morgan, 1999, Pearl, 1995) | incomplete knowledge, untestable |
| Confounder | variable causally associated with both exposure and outcome and not on the causal pathway between $X$ and $Y$ (Pearl, 2009c, Morgan and Winship, 2015, Spirtes et al., 1995) | unobserved processes, hidden confounding, bidirected edges, periodicity, slow-time scales, discrete states |
| Counterfactuals | Attribution of observed events or mediation analysis (Pearl, 1995, Morgan and Winship, 2015) | incomplete knowledge, untestable |
| Covariate shift | The system's probability density function changes through time or space (Sugiyama and Kawanabe, 2012) | hard to detect and correct models for, cf. invariance |
| Dependencies | models trained to maximise statistical measures of linear or nonlinear association, or in-distribution prediction (Kotz and Drouet, 2001) | non-causal, right-for-the-wrong-reasons |
| Faithfulness | only the variables that are $d$-separated in a DAG will be independent (i.e., all others will be dependent) (Pearl, 2009c) | requires an understanding of how different variables work together to cause an outcome |
| Missingness | unobserved quantities or random variables, either missing at random or through particular mechanisms (McKnight et al., 2007) | potentially leading to selection bias |

Table 1.2: Simple taxonomy of models and their level. Figure partly reproduced from (Peters et al., 2017a).

| Model/Level | i.i.d. | distribution shifts | Counterfactuals | Physical insight | Data-driven |
|---|---|---|---|---|---|
| 3 - Mechanistic | ✓ | ✓ | ✓ | ✓ | ✓ |
| 2 - Structural causal | ✓ | ✓ | ✗ | ? | ? |
| 1 - Statistical/ML | ✓ | ✗ | ✗ | ✗ | ? |

in numerous scientific and engineering domains, including climate network analysis [Donges et al., 2009b], financial market network analysis [Fiedor, 2014], and brain network analysis [Johnson et al., 2007, Takagi, 2020]. However, mutual information-based associations at lag zero cannot be interpreted directionally since the term "information" implies a lack of directionality. Other regions or variables may influence the nodes under investigation, or the association could be due to a common driving process. Lagged association measures are commonly used to account for directional links and quantify the time lag of associations.

Lagged correlation analysis has a long history in climate research [Walker, 1923] and neuroscience [Medaglia et al., 2015, Richiardi et al., 2013], with the delay at the maximum of the cross-correlation function being used to interpret the delay of the underlying physical mechanism coupling two processes. Other lagged measures of association, such as mutual information [Granger, 1969], have been proposed to determine lags in nonlinear processes. In addition to analysing time lags, the magnitude of the cross-correlation is often used as a measure of the impact of one process on another or as a measure of the strength of an association. This aligns with the statistical interpretation of the square of correlation as the proportion of variance in one process that can be linearly represented by another [Von Storch and Zwiers, 2002, Chatfield, 2013].

However, relying solely on association measures, even those that account for lags and nonlinearity, cannot uncover directionality, detect the delay of the underlying mechanism, or provide a physically or causally interpretable estimate [Runge et al., 2014]. The widespread use and abuse of association measures in engineering and science throughout the 20th century have impeded the exploration and development of meaningful causation measures and hindered the discovery of new and alternative explanatory laws from data.

**Level 2 – Learning causal relations from observations.**

A fundamental objective in the scientific enterprise is understanding the causes behind the phenomena we observe [Pearl, 2009a, Peters et al., 2017a]. This is particularly challenging in disciplines dealing with complex dynamical systems, where experimental interventions are expensive, unethical, or practically impossible. In some fields (e.g., climate sciences, economics, cardiology, and neurosciences), the current alternative is to rely on computationally expensive simulation experiments. Still, those do not adequately represent all relevant physical processes involved. At the same time, a rapidly increasing amount of time series data is generated by observations and also models. How can we use this wealth of data to gain new insights into our fundamental understanding of these systems?

In recent years, rapid progress has been made in computer science, physics, statistics, philosophy, and applied fields to infer and quantify potential causal dependencies from data without intervening in the systems. Although the truism that correlation does not imply causation holds, the key idea shared by several approaches follows Reichenbach's common cause principle [Reichenbach, 1991]: if variables are dependent, then they are either causal to each other (in either direction) or driven by a common driver. To estimate causal relationships among variables, different methods take different, partially strong, assumptions. Granger [Granger, 1969] addressed this question quantitatively using prediction. At the same time, in the last decades, several complementary con-

cepts emerged, from nonlinear dynamics [Sugihara et al., 2012] based on attractor reconstruction to computer science exploiting statistical independence relations in the data [Pearl, 1995, 2009a]. More recently, statistics and machine learning research utilised the framework of structural causal models (SCMs) [Peters et al., 2017a] for this purpose. Causal inference from data is becoming a mature scientific approach [Pearl, 2009a].

Causal inference combines domain knowledge, ML models, and data to learn the underlying system's causal structure and quantify causal effects [Pearl, 2009b, Pearl et al., 2016, Pearl and Mackenzie, 2018, Spirtes et al., 2000, Peters et al., 2017a]. Causal inference can leverage observational or (interventional) model output data [Pearl, 2009b] to learn, understand and evaluate the plausibility of specific causal relations among the considered variables. Today, many methods and tools are available to address challenges in complex systems [Wagner, 1999, Sugihara et al., 2012] and many other fields. Causality is pivotal not only for a better academic understanding of processes in science but also for more robust forecasts, attributing the causes of events, and improving the physics embedded in physics models. Many fields of science and engineering are using causal inference/discovery methods, from Earth and climate sciences [Runge et al., 2019a, Pérez-Suay and Camps-Valls, 2019, Ebert-Uphoff and Deng, 2017, Raia, 2008, Reitsma, 2010, Niemeijer and de Groot, 2008, Shepherd, 2019, Goodwell et al., 2020], neurosciences [Reid et al., 2019, Siddiqi et al., 2022, Stokes and Purdon, 2017], social sciences [Marini and Singer, 1988, Russo, 2010], health and epidemiology [Hernán, 2018, Glass et al., 2013, Hernan and Robins, 2020], or economics [Hicks et al., 1980, LeRoy, 2004].

**Level 3 – Equation discovery in physical systems.**
The scientific enterprise distinctly differs from other intellectual endeavours by relying on formal theories, laws, and models to explain and predict observations and using such observations to construct, revise, and evaluate its formal statements [Langley et al., 1987b, Langley, 2019, Klahr and Simon, 1999]. Many of these activities have been studied by philosophers of science for over a century. The Logic of Science [Popper, 2005] (or justification) aims to characterise how observational data, simulations, and experiments can collectively support or refute laws, models, or theories.

A common claim was that scientific discovery requires some 'creative spark', which cannot be analysed rationally or logically [Simon et al., 1989]. Popper, Hempel, and many other philosophers of science maintained that the discovery process was inherently irrational and beyond any formal understanding. The key insight came from Simon [Popper, 2005, Hempel, 2001, Simon et al., 1989], who proposed that scientific discovery, rather than *"depending on some unknown mystical ability, is a variety of problem-solving that involves searching through a space of problem states generated by applying mental operators and guided by heuristics to make the search tractable."* Such observation established the first heuristic programming methods of hypothesis (model) search to automate the creative process and law discovery. Discovering numeric laws from data has been approached by many authors in the past using grammars, logic rules, propositional bases, entailment, and genetic algorithms, to name a few [Falkenhainer and Michalski, 1986, Kokar, 1986, Żytkow et al., 1990, Schaffer, 1990, Nordhausen and Langley, 1990, Langley, 2019, Moulet, 1992, Gordon et al., 1994, Murata and Tanaka, 1994, Džeroski and Todorovski, 1995, Washio and Motoda, 1997, Bradley et al., 2001, Koza et al., 2001]. Later, [Schmidt and Lipson, 2009b] proposed an automated algorithm to discover Hamiltonians, Lagrangians, and other geometric and momentum conservation laws without prior knowledge of physics, kinematics, or geometry. A new field was born; learning explicit mathematical laws from observations, which was often referred to as *equation discovery* or *data-driven system identification* [Simidjievski et al., 2020].

Inspired by earlier work on the DENDRAL system [Feigenbaum et al., 1971], which inferred structural models of organic molecules from their mass spectra, the community developed different systems that created models of other scientific phenomena (e.g., [Langley et al., 1987a]). The

field was named computational scientific discovery, and the challenge of automating it has been approached by many researchers since then [Evans and Rzhetsky, 2010, Fortunato et al., 2018, Bongard and Lipson, 2007, Schmidt et al., 2011, Waltz and Buchanan, 2009, King et al., 2009]. Efforts in this paradigm differ from mainstream work in machine learning by producing scientific formalisms [Langley et al., 2002a,b, 1987b], ranging from componental models in particle physics [Kocabas, 1991] to reaction pathways in chemistry and to regulatory models in genetics [King et al., 2004]. Reviews have been edited by Shrager and Langley [1990], Dzeroski and Todorovski [2007], Simidjievski et al. [2020].

Recently, the field has been approached by scientists in AI, functional analysis and mathematical operators, nonlinear control, and system identification. Modern approaches that we will review in this paper consider: Automated reverse engineering of nonlinear dynamical systems [Bongard and Lipson, 2007], sparse-promoting solutions that identify parsimonious models of nonlinear dynamics; e.g. relevance vector machine, a sparse Bayesian regression method [Tipping, 2001] or the SINDy method [Brunton et al., 2016b, Brunton and Kutz, 2022], which has been combined with deep neural networks [Champion et al., 2019], reduced-order models [Mezić, 2005] and Koopman operators based on kernel theory and autoencoders [Kaiser et al., 2020, 2021b, Kostic et al., 2022, Klus et al., 2020, Lusch et al., 2018], differentiable networks that can learn the true underlying equation and extrapolate to unseen domains [Sahoo et al., 2018], a constrained symbolic regression methodology, named AI Feynman, that enforces desirable constraints in equation learning (compositionality, units, separability, symmetry, smoothness) [Udrescu and Tegmark, 2020], genetic programming to distil laws of physics [Schmidt and Lipson, 2009a], or a transformer-based architecture using massive pretraining can predict formulas from data [Biggio et al., 2021].

A relevant challenge in this context is the discovery of state variables from experimental/observational data. Almost in all equation and causal discovery algorithms, the variables are given or assumed, which is impossible when trying to understand new or highly complex systems. Some approaches exist in the literature based on the combination of learning compact and expressive feature representations, manifold learning for the determination of the intrinsic dimensionality of the system representation coordinates, and SINDy as a regulariser that enforces dynamic equations in the state space [Brunton and Kutz, 2022, Chen et al., 2022, Pukrittayakamee et al., 2009, Schütt et al., 2017]. These approaches are somewhat related to the discovery of causal relations under noise and latent functions, which is also an active field of research [Stegle et al., 2010, Monti et al., 2020, Diaz et al., 2023].

### 1.3.3 AI for scientific discovery

The debate about AI-based theories of scientific discovery has been ongoing for decades, beginning with whether computers can devise new concepts or merely process the concepts already included in a given computer language. However, the discussion has been revived with the development of new computational tools for data analysis. It is now largely uncontroversial that machine learning tools can aid discovery, though there is still debate about whether they generate new knowledge or merely speed up data processing. Moreover, there is the question of whether data-intensive science fundamentally differs from traditional research and the ethical implications of "superhuman AI". Philosophers have also focused on the opacity of machine learning, asking whether we can say that humans and machines are "co-developers" of knowledge. Ultimately, the debate about AI-based theories of scientific discovery is still ongoing, with researchers considering both the potential benefits and ethical implications of such tools.

> *"With the advent of data-driven methods that learn patterns and relations from data, the tedious human endeavour of scientific discovery (laws, equations and causes of phenomena) is being revolutionised... and accelerated in many fields."*

Figure 1.5: Roadmap of the article.

The fields of equation discovery and causal inference that we will review in this paper promise to shed light on the previous fundamental questions: what can an algorithm learn and discover, what can AI explain, and what new science may emerge through a collaborative AI-machine dialogue about science? An integrative approach seems necessary, where domain experts, data and machine work together in a data-driven framework that formulates and answer causal questions and discover new laws [Pearl, 2009a]. At a more general level, it becomes pertinent to ask if machines can start a new scientific revolution and even if AI itself is the ultimate scientific revolution [Russell, 2021, Boden, 2014, Gillies, 1998].

## 1.4  Outline

This paper aims to review the most important concepts, methods, and previous works on causal inference and discovery in the physical sciences. We use statistical learning techniques to discover causal relations, physical laws, and governing equations from data. Sec. 2 and Sec. 3 present general frameworks and taxonomies for causal discovery and learning physical laws from data, respectively. Both sections categorise the field, reviewing concepts and methods, their specific characteristics, challenges, and opportunities in the physical sciences. Sec. 4 provides examples of causal discovery and equation discovery in a wide range of fields of the physical sciences: dynamical systems, neuroscience, classical and quantum systems, fluid mechanics, geosciences and climate sciences. We pay attention to how causality concepts and methods can improve our knowledge of a given physical system from observations. Sec. 5 outlines the most promising future lines of research in this area of study at the intersection of machine learning and nonlinear physical processes.

# 2. Causal discovery in the physical sciences

Causal discovery, see for example [Spirtes et al., 2000, Peters et al., 2017a] for extended expositions of the topic, has become increasingly popular in the last years as a tool to discover the underlying causal structure of physical systems [Runge et al., 2019a]. There is an abundance and ever-growing number of methods designed to work under different assumptions and tackle other use cases. This section reviews several methods for causal discovery, focusing on methods for time series and their potential use in physical sciences. To this end, in Sec. 2.1, we first provide a taxonomy for many available causal discovery methods. Sec. 2.2 discusses causal discovery's challenges in real-world applications. We conclude in Sec. 2.3 by discussing opportunities for applications of causal discovery in the physical sciences. We would also like to point to other reviews focusing on causal discovery of time series [Runge et al., 2019a, Assaad et al., 2022b, Moraffah et al., 2021]. In addition, Runge et al. [2023] provides a shorter accessible summary of methods for causal discovery and causal effect estimation with practical case studies to illustrate typical challenges, such as contemporaneous causation, hidden confounding and non-stationarity.

## 2.1 A taxonomy of causal discovery methods

In this section, we structure the zoo of existing causal discovery methods to guide method users in finding a method suitable for their application and guide method developers in identifying open challenges. To this end, we summarise the central formal aspects of the graphical-model-based causal inference framework in Sec. 2.1.1. Then, in Sec. 2.1.2, we discuss several characteristics (hereafter referred to as "axes") by which methods can be conceptually distinguished. Lastly, in Sec. 2.1.3, we present an extensive (but not exhaustive) list of causal discovery methods and characterise these methods according to previously introduced axes.

### 2.1.1 Preliminaries

At the heart of the graphical-model-based causal inference framework are *structural causal models (SCMs)*, e.g. [Bollen, 1989, Pearl, 2009c, Peters et al., 2017a]. An SCM serves as a causal model for the data-generating process and specifies how the system reacts to *interventions*, that is, to

Table 2.1: Taxonomy of methods for causal discovery. The entries in parentheses $(\cdot)$ indicate that there are versions of the algorithm in which the assumption is relaxed. We use the abbreviations 'TSG' for time series graph, 'summary' for summary graph and 'ext. sum. graph' for an extended summary graph.

| Method | Target of Inference | | Approach | Process Assumptions | | | | Data Assumption |
|---|---|---|---|---|---|---|---|---|
| | Bi-/ Multi-variate | Graph type | Indep./ Asymm./ Score | non-/linear | Stoch. / Det. | Con-temp. | Cycles | Hidden var. |
| GC [Granger, 1969] | Bi. | summary | indep. | linear | stoch. | ✗ | lagged-only | ✗ |
| Mult-GC [Geweke, 1982] | Multi. | summary | indep. | linear | stoch. | ✗ | lagged-only | ✗ |
| Mult-nonlin-GC [Diego Bueso, 2020] | Multi. | summary | indep. | nonlin. | stoch. | ✗ | lagged-only | ✗ |
| TE [Schreiber, 2000] | Bi. | summary | indep. | nonlin. | stoch. | ✗ | lagged-only | ✗ |
| Multi-TE [Barnett et al., 2009] | Multi. | summary | indep. | nonlin. | stoch. | ✗ | lagged-only | ✗ |
| CCM [Sugihara et al., 2012] | Bi. | summary | indep.(?) | nonlin. | det. | ✓ | ? | part-ially |
| Ext.-CCM [Diaz et al., 2022] | Bi. | summary | indep.(?) | nonlin. | det. | ✓ | ? | part-ially |
| tsPC [Runge, 2020] | Multi. | TSG | indep. | both | stoch. | ✓ | (✓) | ✗ |
| PCMCI [Runge et al., 2015b] | Multi. | TSG | indep. | both | stoch. | ✗ | lagged-only | ✗ |
| PCMCI$^+$ [Runge, 2020] | Multi. | TSG | indep. | both | stoch. | ✓ | (✓) | ✗ |
| PCGCE [Assaad et al., 2022a] | Multi. | ext. sum. graph | indep. | both | stoch. | ✓ | (✓) | ✗ |
| FCIGCE [Assaad et al., 2022a] | Multi. | ext. sum. graph | indep. | both | stoch. | ✓ | (✓) | ✓ |
| tsFCI | Multi. | TSG | indep. | both | stoch. | (✓) | (✓) | ✓ |
| SVAR-FCI [Assaad et al., 2022a] | Multi. | TSG | indep. | both | stoch. | ✓ | (✓) | ✓ |
| SVAR-GFCI [Malinsky and Spirtes, 2018] | Multi. | TSG | score & in-dep. | both | stoch. | ✓ | (✓) | ✓ |
| LPCMCI [Gerhardus and Runge, 2020] | Multi. | TSG | indep. | both | stoch. | ✓ | (✓) | ✓ |
| (F)GES [Meek, 1997, Chickering, 2002b,a, Ramsey et al., 2017] | Multi. | summary | score | linear | stoch. | ✓ | lagged-only | ✗ |
| DYNOTEARS [Pamfil et al., 2020] | Multi | TSG | score | linear | stoch. | ✓ | lagged-only | ✗ |
| IDYNO [Gao et al., 2022] | Multi | TSG | score | linear and non-linear | stoch. | ✓ | lagged-only | ✗ |
| NTS-NOTEARS [Sun et al., 2023] | Multi | TSG | score | linear and non-linear | stoch. | ✓ | lagged-only | ✗ |
| TiMiNo [Peters et al., 2013] | Multi. | TSG | indep. | both | stoch. | ✓ | lagged-only | ✗ |
| RHINO [Gong et al., 2023] | Multi. | TSG | indep. | both | stoch. | ✓ | lagged-only | ✗ |
| VARLiNGAM [Shimizu et al., 2006] | Multi. | TSG | asymm. | linear | stoch. | ✓ | lagged-only | ✗ |

idealised experimental manipulations that deliberately hold fixed a subset of the system's variables while not perturbing the system in any other way.

An SCM for a system described by the set of variables $\mathbf{V} = \{V^1, \ldots, V^n\}$ consists of $n$ so-called *structural assignments*

$$V^i := f^i(pa^i, \varepsilon^i) \quad \text{with } 1 \leq i \leq n, \tag{2.1}$$

together with a product distribution $p_\varepsilon(\varepsilon^1, \ldots, \varepsilon^n) = p_\varepsilon^1(\varepsilon_1) \cdot \ldots \cdot p_\varepsilon^n(\varepsilon_n)$ of the random variables

$\varepsilon^i$. Formally, the $f^i$ are measurable functions that depend non-trivially on all of their input arguments, and the $pa^i \subseteq \mathbf{V} \setminus \{V^i\}$ are subsets of the system variables $V^1, \ldots, V^n$. The functions $f^i$ are interpreted as the causal mechanisms by which the values of the respective variable $V^i$ are determined from the value of $\varepsilon^i$ and the values of the variables in $pa^i$. Consequently, the variables in $pa^i$ are referred to as the *causal parents* of $V^i$. The random variables $\varepsilon^i$ are interpreted as noise that summarises all factors that are not modelled explicitly, and the factorisation of $p_\varepsilon(\varepsilon^1, \ldots, \varepsilon^n)$ amounts to the assumption that the $\varepsilon^1, \ldots, \varepsilon^n$ are jointly independent. This joint independence is motivated by the view that any dependence between the noise variables must be due to a causal relationship between them and that such a dependence should then rather be modelled explicitly by enlarging the set $\mathbf{V}$ of system variables.

The *causal graph* of an SCM with system variables $V^1, \ldots, V^n$ is the directed graph whose vertices are the variables $V^i$ and with a directed edge $V^i \rightarrow V^j$ if and only if $V^i$ is a causal parent of $V^j$, that is, if and only if $V^i \in pa^i$. Consequently, the causal graph of an SCM shows the *qualitative* cause-and-effect relationships as specified by the sets $pa^i$. If the causal graph is acyclic, that is, if the causal graph is a directed acyclic graph (DAG), then the SCM is said to be *acyclic*.

An SCM obtains causal meaning by asserting how the modelled system reacts to interventions. Formally, an intervention on the variable $V^k \in \mathbf{V}$ is a mapping, conventionally denoted as $do(V^k := v^k)$, that maps the original SCM and a number $v_k$ to a new SCM in which the original structural assignment for $V^k$ is replaced by the new structural assignment $V^k := v^k$ and the noise variables $\varepsilon^k$ is removed. This new SCM is referred to as an *intervened SCM*, and $do(V^k := v^k)$ is interpreted as an idealised experimental manipulation by which the value of $V^k$ is held fixed at $v_k$ while leaving the system unaltered else. This specification of how the system reacts to interventions is why using the symbol ":=" instead of "=" in (2.1) is conventional. On the level of causal graphs, $do(V^k := v^k)$ amounts to removing all edges that point into $V^k$ because, in the intervened SCM, the variable $V^k$ has no causal parents. Interventions on subsets of variables are defined similarly.

In an acyclic SCM, the combination of noise distribution $p_\varepsilon$ and functions $f^i$ uniquely determines a distribution of the system variables $V^1, \ldots, V^n$. This distribution is often referred to as the *entailed distribution of the SCM*. The entailed distribution $p(\cdot)$ of the original SCM (that is, of the SCM that models that system without interventions) is often referred to as the *observational distribution*. The entailed distributions of the intervened SCMs are often referred to as *interventional distributions* and are conventionally often denoted as $p(\cdot \,|\, do(V^k := v^k))$; and similarly for interventions on subsets of system variables.

When using this notation, it is important to keep in mind that $p(\cdot \,|\, do(V^k := v^k))$ is, in general, not equal to $p(\cdot \,|\, V^k = v_k)$. Indeed, $p(\cdot \,|\, do(V^k := v^k))$ is the distribution of the *intervened* SCM, whereas $p(\cdot \,|\, V^k = v_k)$ corresponds to *observing* $V^k = v^k$; put differently: Correlation is not equal to causation.

The article [Bongers et al., 2021] discusses in much detail the more complicated case of cyclic SCMs. As shown there, cyclic SCMs need not entail a unique distribution for the system variables. However, [Bongers et al., 2021] defines a restricted class of cyclic SCMs, termed *simple SCMs*, that entail a unique distribution and that are closed under interventions. Moreover, acyclic SCMs are a special case of simple SCMs.

In the time series case, which we are predominantly interested in this paper, (2.1) can be generalised by putting a time index $t$ on $V^i$, $f^i$, $pa^i$ and $\varepsilon^i$. The commonly used term *causal stationary* then refers to time-invariance of the qualitative cause-and-effect relationships, that is, to the situation that $pa^i_{t+\Delta t} = \{V^j_{s+\Delta t} \mid V^j_s \in pa^i_t\}$ for all $t$ and $\Delta t$.

## 2.1.2 Axes for categorising causal discovery methods

This section introduces and explains several axes for categorising and distinguishing causal discovery methods. While it would be possible to consider more axes yet, the authors believe that

the choice of axes presented here is a reasonable compromise between a sufficiently fine-grained categorisation on the one hand and clarity of exposition on the other hand. Table 2.1 lists many causal discovery methods and categorises them according to the aforementioned axes. In Figure 2.1, we graphically illustrate some axes.



Figure 2.1: Illustration of some axes along which causal discovery methods are categorised in this review; see references in the main text to the respective subparts of the figure.

**Bivariate vs multivariate causal discovery.**

This axis concerns the number of variables that are being considered. Bivariate causal discovery aims to discover the causal relationship between exactly two variables $X$ and $Y$ (in the non-temporal case) or between exactly two component time series $X^i$ and $X^j$ (in the time series case). Multivariate causal discovery aims to discover the causal relationships between any number of variables or component time series, respectively. Bivariate causal discovery often (but not necessarily) assumes *causal sufficiency* (see axis on causal sufficiency below). In the time series case, bivariate causal discovery often (but not necessarily) targets to infer the *summary graph* rather than the *time series graph* or *extended summary graph* (see axis on time series graph discovery below). If time lags are at least partially resolved in the bivariate time series case, that is, if the target of inference is the time series graph or the extended summary graph, then one effectively deals with a multivariate causal discovery problem.

**Time series graph discovery vs summary graph discovery vs extended summary graph discovery.**

This axis is specific to the temporal setting and concerns the target of inference. Some methods are designed to learn the *time series graph* [Runge et al., 2012], also known as *full-time graph* [Peters et al., 2013] and *time series chain graph* [Dahlhaus and Eichler, 2003], that is, the collection of all causal links $X_{t-\tau}^i \to X_t^j$ including the respective lags $\tau$ of these links. Part $a)i)$ of Figure 2.1 shows an example of a time series graph with four component time series. As indicated by the grey edges, the pattern of edges in this graph is implicitly assumed to repeat both to the left (past) and right (future). Due to this repetitive structure of the edges, a time series graph is uniquely specified by the collection of edges that point into a vertex at an arbitrary reference time step $t$. Other methods disregard the information about the time lags and instead learn the *summary graph* [Peters et al., 2017a]. In the summary graph, there is exactly one vertex per component time series $X^i$ and an edge $X^i \to X^j$ if and only if there is an edge in the time series graph $X_{t-\tau}^i \to X_t^j$ at any lag $\tau$. Part $a)iii)$ of Figure 2.1 shows the summary graph associated with the time series graph in part $a)i)$ of the same figure. Another option is to learn *extended summary graphs* [Assaad et al., 2022a]. These graphs go midway between learning time series and summary graphs by distinguishing between contemporaneous and lagged links but disregarding the information about the specific time lags of lagged links. Specifically, the extended summary graph contains exactly two vertices per component time series $X^i$, namely the vertex $X_t^i$ for the present time steps and the vertex $X^{i,-}$ for all past time steps. There is an edge $X_t^i \to X_t^j$ if and only if, this same edge is also in the time series graph, there is an edge $X^{i,-} \to X_t^j$ if only if there is at least one $\tau \geq 1$ such that $X_{t-\tau}^i \to X_t^j$ is in

the time series graph, and there is no edge between $X^{i,-}$ and $X^{j,-}$. By this construction, extended summary graphs do distinguish between $X^i_{t-\tau} \to X^j_t$ with $\tau > 0$ and $X^i_t \to X^j_t$ but do not distinguish between, for example, $X^i_{t-2} \to X^j_t$ and $X^i_{t-1} \to X^j_t$. Part $a)ii)$ of Figure 2.1 shows the extended summary graph associated with the time series graph in part $a)i)$ of the same figure. Resolving the lag structure does yield more information but also implies a more complex target of inference. Learning more complex graphs (e.g., a time series graph vs a summary graph) is conceptually and statistically more challenging.

Methods for time series causal discovery typically require the user to specify a maximal lag $\tau_{max}$ up to which the method is supposed to be sensitive. If the target of inference is the time series graph, then this choice is apparent as the learned graph has exactly $\tau_{max} + 1$ steps. For example, if the user specifies the maximal lag to be $\tau_{max} = 3$, then the learned graph has exactly $4 = 3 + 1$ time steps, and the method cannot find causal links of lag $\tau \geq 4$ even if such links are present. If time series graph discovery is combined with the assumption of causal sufficiency, then this combination implicitly implies no causal links with a lag larger than $\tau_{max}$. If causal sufficiency is not assumed, then this additional implicit assumption is not made. Also, methods for summary and extended summary graph discovery typically require the specification of a maximal lag up to which the method is sensitive. As opposed to that, when learning summary graphs or extended summary graphs the choice of $\tau_{max}$ is not apparent from the learned graph. However, this choice is not apparent from the learned graph in the summary graph and extended summary graph cases.

For the case of causal models for time series, the causal arrow between any two variables can be at any time lag ranging from zero, i.e. a contemporaneous link, to infinity. A causal discovery algorithm can thus aim to either infer the full graph with edges being classified by time lags or to infer the summary graph where the time lag of the causal effect is irrelevant.

### Methods based on independence, asymmetry, scores and context.

This axis distinguishes causal discovery methods by the type of information/signal that they use to learn the causal graphs from data. In this review, as is common in the literature (for example, see [Glymour et al., 2019]), we distinguish the independence-based, asymmetry-based, and score-based approaches. Further, we consider the context-based approach to distinguish those methods that employ the invariance of causal mechanisms across different environments. An exact delineation between these four approaches is not always possible as there are *hybrid* methods that combine more than one approach.

> *"The wide variety of causal discovery methods can be structured into independence-based, asymmetry-based, score-based, and context-based approaches."*

First, *independence-based causal discovery*, sometimes called *constraint-based causal discovery*, utilises marginal and conditional independencies between variables to learn the causal graph or a set of causal graphs consistent with those independencies. Recall that an SCM is defined by a collection of structural assignments for each variable, where each assignment is a function of the variable's parents and a noise term. The collection of noise variables is assumed to be jointly independent. Independence-based causal discovery relies on the fact that, for data generated by an SCM, the structure of the SCM's causal graph imprints some independencies onto the data [Verma and Pearl, 1990a, Geiger et al., 1990, Pearl, 2009c]. This property is known as *causal Markov condition* [Spirtes et al., 2000]. Alternatively, if one does not assume that an SCM generates the data, then the causal Markov condition is not automatically implied but needs to be assumed separately, leading to the so-called *causal Markov assumption*. The *d*-separation criterion [Pearl, 1988] allows to graphically determine all independencies that are necessarily implied in a given causal graph [Verma and Pearl, 1990b, Geiger et al., 1990, Pearl, 2009c]. The basic idea then is to run statistical tests of marginal and conditional independencies on the data and, second, use the results of these tests to constrain the causal graph's structure.

For the second of these two steps to hold, one further needs to make the *causal faithfulness assumption* [Spirtes et al., 2000]. This assumption says there are no independencies beyond those necessarily implied by the causal Markov condition in the observed data. For example, causal faithfulness excludes that a variable $X$ causally influences another variable $Y$ along multiple pathways which in total cancel out exactly—because otherwise there would be statistical independence, namely the marginal independence of $X$ and $Y$, that would not be implied by the causal Markov condition.

For linear models and infinite samples, causal faithfulness is guaranteed except for Lebesgue measure zero sets in parameter space [Spirtes et al., 2000]. However, [Robins et al., 2003] shows that causal faithfulness ensures point-wise consistency but not uniform consistency of, for example, the famous PC algorithm [Spirtes et al., 2000] (see below for the details on the PC algorithm). Uniform consistency of the PC algorithm is proven under the stronger *strong-faithfulness assumption* [Zhang and Spirtes, 2002, Kalisch and Bühlman, 2007]. However, [Uhler et al., 2013] shows that the set of not strongly-faithful distributions has a non-zero Lebesgue measure for various graph structures. This analysis indicates that causal faithfulness should be regarded as a strong assumption. There are also weaker types of faithfulness assumptions, for example, the *adjacency faithfulness* and *orientation faithfulness* assumptions [Ramsey et al., 2006, Zhang and Spirtes, 2008].

Independence-based causal discovery is non-parametric in that no assumption on the SCM's functional relationships and/or noise distributions needs to be made. However, choosing a particular method for (conditional) independence testing may implicitly impose a parametric assumption. For example, testing for (conditional) independence by (partial) correlation implicitly makes the assumption that the data-generating process is linear Gaussian. Conversely, if a parametric assumption can be made, this assumption might favour specific methods for (conditional) independence testing. For example, suppose one can assume linear Gaussian data. In that case, it is reasonable to use a (partial) correlation instead of more general (conditional) independence tests like, for example, a test based on (conditional) mutual information as given in [Runge, 2018].

Many independence-based causal discovery are proven to be *sound and complete*, meaning they provably learn the respective target of inference if they are given ground-truth knowledge of (conditional) independencies. However, the statistical tests of (conditional) independence are expected to make errors in finite samples even if all assumptions are met. For the famous PC algorithm [Spirtes et al., 2000] (see below for more details), there are also probabilistic finite-sample guarantees [Zhang and Spirtes, 2002, Kalisch and Bühlman, 2007].

Typically, there are multiple graphs that by means of the causal Markov condition, imply the exact same set of (conditional) independencies. For example, the three graphs $X \to Y \to Z$ and $X \leftarrow Y \leftarrow Z$ and $X \leftarrow Y \to Z$ by means of the causal Markov condition all imply exactly the same independence, namely that $X$ and $Z$ are conditionally independent given $Y$ (and no further independencies). Such graphs are said to be *Markov equivalent* to each other and constitute a *Markov equivalence class*. Consequently, independence-based causal discovery algorithms cannot distinguish between Markov equivalent graphs. Instead, these algorithms target learning a graphical representation of the entire Markov equivalence class to which the true causal graph belongs. Since the graphs $X \to Y$ and $X \leftarrow Y$ are Markov equivalent, independence-based causal discovery can not infer the causal direction in the fundamental non-temporal bivariate case (see axis I). Independence-based causal discovery becomes meaningful for at least three variables; however, in the time series case, even two variables suffice to make non-trivial causal inferences (because then, if the variables are resolved in time, there are effectively more than just two variables).

Second, *asymmetry-based causal discovery* makes and relies on parametric assumptions on the form of the functional relationships and/or noise distributions of the data-generating SCM [Peters et al., 2017a].

This approach is motivated by the elementary bivariate case, that is, by finding the causal

relationship between two variables $X$ and $Y$. As explained above, with independence-based causal discovery, it is not possible to distinguish the Markov equivalent graphs $X \to Y$ and $X \leftarrow Y$. This impossibility is not a shortcoming of independence-based causal discovery but rather is fundamental unless stronger assumptions are made [Peters, 2012, Peters et al., 2017a]. The proof of the impossibility of distinguishing between $X \to Y$ and $X \leftarrow Y$ works by showing that if the true data-generating SCM goes in the direction $X \to Y$, then one can always construct an alternative SCM in the direction $X \leftarrow Y$ that gives rise to the same data distribution as the true SCM.

The basic idea for removing this fundamental ambiguity is as follows: For certain choices of *restricted* SCMs, defined by certain restricted parametric assumptions, it is impossible to have a restricted SCM in the first direction $X \to Y$ and simultaneously a restricted SCM in the second direction $X \leftarrow Y$. Hence, given the assumption that the true SCM lies in the restricted class of models, it becomes possible to distinguish the causal and anti-causal direction. A restricted class of SCMs with this property is said to be *identifiable*. This approach to causal discovery relies on the expectation that the SCM in the causal direction generically has lower complexity than any alternative SCM in the anti-causal direction. As explained in Section 4.1.2 of [Peters et al., 2017a], this expectation can be motivated by the principle of *independence of cause and mechanism* [Daniušis et al., 2010, Peters, 2012]. There are also asymmetry-based causal discovery methods for learning the causal graph between two or more variables for multivariate causal discovery [Peters et al., 2017a].

Third, *score-based causal discovery* chooses one or multiple best-scoring graphs with respect to a predefined scoring function. This scoring function is typically built on the likelihood of the observed data given a particular graph and an assumed parametric statistical model [Peters et al., 2017a]. This approach requires searching over the space of causal graphs. Even if causal sufficiency (see axis on causal sufficiency) and acyclicity (see axis on cycles below) are assumed, in which case the causal graph is a *directed causal graph (DAG)*, the search space of graphs already grows super-exponentially, e.g. [Chickering, 2002b]. An exact search is thus infeasible even for a moderate number of variables. Instead, greedy search techniques are often used, e.g. in the famous GES algorithm [Chickering, 2002b,a] (see below for details on this algorithm). If the assumed statistical model does not yield identifiability beyond the Markov equivalence class, then the scoring function must be chosen such that Markov equivalent graphs have the same score. Hence, one can search over the space of Markov equivalence classes rather than over the space of graphs.

Fourth, *context-based causal discovery* requires access to data of the same system in different contexts. The term *different context* is understood rather broadly: Its meaning ranges from, for example, observing the same physical system at other locations to, for example, observing a system both before and after an intervention. The basic assumption and idea of this approach to causal discovery are that the causal mechanisms, that is, the functional mappings from causes to effects and hence also the conditional distributions of the effects given their causes, remain unchanged across all contexts (unless the effect variable is the target of an intervention in one of the contexts). In contrast, marginal distributions and hence also the conditional distributions of causes given their effects can change [Peters et al., 2017a]. A prime example of a context-based causal discovery method is Invariant Causal Prediction [Peters et al., 2016] (see below for more details). The *joint causal inference (JCI)* framework [Mooij et al., 2020] proposes to model all contexts with one graph by including one or multiple so-called *context variables* whose values determine the context and subsequently pooling the data from the different contexts into one joint dataset. This enlargement of the system by context variables effectively reduces the case of multiple contexts to the standard case of a single context. Consequently, it is possible to apply standard causal discovery algorithms to the pooled data (if the respective assumptions are met). One might thus argue that context-based causal discovery should not be considered a distinct approach to causal discovery but should instead be characterised by the requirement of data from multiple contexts.

### Linear or nonlinear dependencies.

This axis concerns the form of the functional relationships in the data-generating structural causal model. Broadly, see part $b$) of Figure 2.1, one can distinguish between linear and nonlinear functional relationships. In independence-based and score-based causal discovery, an assumption of linearity can enter implicitly by using partial correlation for testing conditional independence (independence-based approach) or by choice of the statistical model (score-based approach). In asymmetry-based causal discovery, an assumption of linearity is, if made, typically explicit by choice of the functional model. Various asymmetry-based causal discovery methods do not assume linearity but still use restricted functional model classes that do not allow for entirely generic dependencies, for example, the functional model class of nonlinear additive noise models [Hoyer et al., 2008]. However, for simplicity, we here only distinguish the methods by whether or not they assume linearity.

### Deterministic vs stochastic systems.

This axis concerns an assumption on the type of data-generating process. Some methods assume the data are generated by a deterministic process, for example, a deterministic dynamical system. In contrast, other methods make explicit use of the assumption that the data-generating process is inherently stochastic, see part $c$) of Figure 2.1. In the case of stochastic data-generating processes, the stochasticity is interpreted as dynamical noise that arises due to factors outside of the model. Dynamical noise needs to be distinguished from measurement noise: The former is an inherent property of the data-generating process, and the latter arises from uncertainty in the data-collection process. The causal inference and discovery frameworks have also been extended to dynamical systems, both deterministic and stochastic, without stable equilibrium distribution [Mooij et al., 2013b, Bongers and Mooij, 2018, Rubenstein et al., 2018].

### Contemporaneous links.

This axis is specific to the temporal setting and concerns a connectivity assumption on the causal time series graph. Some methods make the assumption that all causal links in the time series graph are *lagged*, meaning that all causal links are of the form $X_{t-\tau}^i \to X_t^j$ with $\tau > 0$, whereas *contemporaneous* links, that is, links of the form $X_t^i \to X_t^j$ are assumed to be absent. Other methods do not make this assumption. For example, in the time series graph in part $a)i$) of Figure 2.1, there are the contemporaneous edges $Z_t \to W_t$ and $Y_t \to X_t$. Consequently, methods that assume the absence of contemporaneous links would, by assumption, disallow this particular time series graph. Contemporaneous causal links correspond to causal influences that act on a time scale shorter than the measurement interval; for example, a causal influence on a time scale of six hours in daily measured data.

### Causal cycles.

This axis concerns a connectivity assumption on the causal graph. Many methods assume the absence of cyclic causal relationships. This assumption means that a variable $X_t^j$ cannot be a causal ancestor of another variable $X_{t-\tau}^i$ if that second variable $X_{t-\tau}^i$ is a causal ancestor of the first variable $X_t^j$. For example, the lower graph in part $d$) of Figure 2.1 has the causal cycle $X_t \to Y_t \to X_t$. Because causation cannot go backwards in time, the assumption of acyclicity only restricts the contemporaneous section of the causal time series graph. The assumption is thus only relevant for $\tau = 0$. In particular, even with the assumption of acyclicity, it is possible to model temporal feedbacks. For example, although the upper graph in part $d$) of Figure 2.1 is acyclic, it displays a causal influence of time series $X$ on $Y$ (by the edge $X_t \to Y_t$) and a causal influence of time series $Y$ on $X$ (by the edge $Y_{t-1} \to X_t$). For example, if the component time series $X^i$ causally influences another component time series $X^j$ at, say, lag $\tau = 1$ (i.e., in the time series graph, there is the link $X_{t-1}^i \to X_t^j$ or an indirect directed path from $X_{t-1}^i$ to $X_t^j$), then one can still allow a causal influence

of $X^j$ on $X^i$ (e.g., $X^j_{t-2} \to X^i_t$ or $X^j_t \to X^i_t$ or indirect directed paths from $X^j_{t-2}$ or $X^j_t$ to $X^i_t$). It is not allowed, however, that $X^i$ and $X^j$ causally influence each other both at lag $\tau = 0$. For example, if both $X^i_t \to X^j_t$ and $X^j_t \to X^k_t \to X^i_t$ are present, then the time series graph is cyclic. There are also causal discovery methods that allow cyclic causal relationships, e.g. [Forré and Mooij, 2018, Strobl, 2019, M. Mooij and Claassen, 2020, Bongers et al., 2021]. These methods typically infer less informative graphs than those inferred by methods that do not allow causal cycles. The early work [Richardson, 1996] considers the special case of causal discovery in *linear* cyclic systems.

**Causal sufficiency.**
This axis concerns an assumption that can be viewed as an assumption on the data-generating process or the data-collection process. The assumption of *causal sufficiency* [Spirtes et al., 2000] says that there are no *latent confounders*, also called *unobserved confounders* or *hidden common causes*. A latent confounder is an unobserved variable that (potentially indirectly through other unobserved variables) causally influences two observed variables $X^i_{t-\tau}$ and $X^j_t$. For example, in the lower graph in part $e$) of Figure 2.1 the unobserved variable $Z$ acts as a latent confounder of the variables $X$ and $Y$. Consequently, this graph violates causal sufficiency whereas the upper graph in part $e$) of Figure 2.1 satisfies causal sufficiency. Methods that do not assume causal sufficiency typically infer graphs that are less informative than the graphs inferred by methods that do assume causal sufficiency. For example, consider the elementary non-temporal bivariate case with two variables $X$ and $Y$. If these variables are dependent and causal sufficiency is assumed, then either $X$ causes $Y$ ($X \to Y$), or $Y$ causes $X$ ($X \leftarrow Y$). If causal sufficiency is not assumed, then there is the third possibility that neither $X$ causes $Y$ nor vice versa but that there rather is an unobserved variable $L$ which causes both $X$ and $Y$ ($X \leftarrow L \to Y$ or, for example, $X \leftarrow L' \leftarrow L \to L'' \to Y$ with other unobserved variables $L'$ and $L''$).

### 2.1.3  Description and categorisation of causal discovery methods

Here, we list and briefly explain several existing causal discovery methods for time series data. We summarise this list and the placement of each method with respect to the axes presented above in Table 2.1. In Figure 2.2, we illustrate and compare an example time series graph and the respective graphical objects obtained by some of the discussed causal discovery algorithms when applied to data generated from an SCM with that time series graph.

**Granger Causality.**
Granger Causality (GC) [Granger, 1969, 1980] is originally a statistical test to decide whether a time series $X_t$ is a *cause* of another time series $Y_t$, in the sense that past values of $X_t$ have significant predictive power in forecasting $Y_t$. GC is thus, in principle, a simple test of temporal (or lagged) relationship and predictability. Nevertheless, under causal sufficiency and no contemporaneous effects assumptions, it can be formally shown that GC testing detects actual causal links (see e.g. Peters et al. [2017a] for a formal derivation of these results in the SEM setting).

In a multivariate setting, which is seldom the case, testing if $X_t$ causes $Y_t$ requires controlling for all possible confounders. Therefore the conditional GC [Granger, 1980, Geweke, 1982, Chen et al., 2004, Barrett et al., 2010], includes in the restricted and full models the past of all other *relevant* time series in the system that are not $X_t$ and $Y_t$. Classically, GC considers linear models for which standard $t$-tests or $F$-tests can be employed, but non-linear extensions have been considered both in econometrics [Bell et al., 1996, Hiemstra and Jones, 1994, Abhyankar, 1998, Warne, 2000, Diks and Panchenko, 2006] and in physical and biological applications [Ancona et al., 2004, Marinazzo et al., 2008b, Diego Bueso, 2020].

Figure 2.2: Figure illustrating a time series graph (TSG) and the respective graphs discovered by applying various causal discovery methods to data generated from an SCM with that time series graph. (a) Time series graph. (b) The discovered undirected graph by considering (lagged) correlations, where spurious correlations are highlighted as dashed grey lines. (c) The directed graph discovered by multivariate Granger causality does not consider contemporaneous links and retains a spurious link from $Y$ to $W$. (d) Graph discovered by CCM. (e) The graph discovered by applying the plain PC and (F)GES algorithms fail to show lagged links and, in addition, fail to orient a link that the time series adapted algorithms can orient. (f) The time series version of PC (tsPC), PCMCI+, and the time series version of GES (tsGES) discover both lagged and contemporaneous links and orient edges up to the Markov equivalence class. (g) Plain FCI has the same drawbacks as plain PC or (F)GES. (h) FCI-based time series causal discovery algorithms account for latent confounders and thus discover causal arrows up to latent confounding. In particular, the algorithm cannot exclude that the association between $Y_{t-1}$ and $Z_t$ is due to latent confounding rather than a causal relationship. (i) The PCGCE algorithm discovers the extended summary graph up to its Markov equivalence class. (j) Var-LiNGAM discovers all causal relationships correctly if the assumptions of linear relationships and additive non-Gaussian noise are satisfied.

---

**Linear and nonlinear Granger causality in dynamic systems.**

A unified view of (nonlinear) GC with kernel methods for the physical sciences was introduced in [Diego Bueso, 2020]. Two examples are given here:

1. Bivariate system with coupled non-linear and autoregressive relations given by $x_{t+1} = 3.4x_t(1 - x_t^2)\exp(-x_t^2) + \varepsilon_t^x$ and $y_{t+1} = 3.4y_t(1 - y_t^2)\exp(-y_t^2) + 0.5x_ty_t2 + \varepsilon_t^y$, where $\varepsilon$ is white Gaussian noise with zero mean and variance 0.4. The causal direction is $X \to Y$. The histogram of the estimated causal index $\delta$ on the left figure reveals GC's insensitivity to the causal direction, the high false positive rate of KGC [Marinazzo et al., 2008a], and a higher detection power by the XKGC [Diego Bueso, 2020].

2. Two logistic maps system defined as $x_{t+1} = 1 - 1.8x_t^2$ and $y_{t+1} = (1 - \alpha)(1 - 1.8y_t^2) + \alpha(1 - 1.8x_t^2)$, where $\alpha \in [0,1]$ controls the coupling strength. The causal relationship implemented is $X \to Y$, and the challenge is to assess the detection power of methods without introducing any external variable, just using $X$ and $Y$. We analyse segments of length $n = 2000$ and fixed $p = 2$. The left figure shows the prediction skills for varying $\alpha$. Note that the system is completely synchronised at $\alpha = 0.37$. T XKGC [Diego Bueso, 2020] improves detection power over GC/KGC for any $\alpha$.

Transfer entropy [Schreiber, 2000] between $Y_t$ and $X_t$ measures the amount of unique information contained in the past of $X_t$, denoted $X^- = \{X_{t-1}, X_{t-2}, \ldots\}$, about the state of $Y_t$ and is defined as $\mathscr{T}_{X \to Y|Z} = H(Y_t|Y^-, Z^-) - H(Y_t|Y^-, X^-, Z^-)$. Transfer entropy can be considered as the generalisation of GC by extending the implicit conditional independence test to arbitrary orders of dependence. Indeed, Barnett et al. [2009] proved that (linear) GC and transfer entropy causality are equivalent under the assumptions of VAR model class and Gaussian error distributions.

**CCM.**

Convergent cross-mapping (CCM) [Sugihara et al., 2012] is based on the simple observation that if data from a deterministic dynamical system is generated by a system of ordinary differential equations (ODEs), then the explicit form of these equations directly defines the causes of each variable in the system: $x$ is a cause of $y$ if the dynamics (any of the derivatives of $y$) is expressed in terms of the state of $x$. For example, in the Lorenz attractor system:

$$x' = \sigma y - \sigma x$$
$$y' = -xz + \rho x - y$$
$$z' = xy - \beta z \qquad\qquad (2.2)$$

$x_t$ is caused by $y_t$, $y_t$ is caused by $x_t$ and $z_t$ and $z_t$ by $x_t$ and $y_t$ as summarised by the summary directed graph of Fig. 2.3.



Figure 2.3: Summary graph for Lorenz attractor system.

Learning the generating ODE from time series data would allow us to recover the causal relations and summary-directed graph. Nevertheless, relying on Takens' theorem [Takens, 1981], Sugihara et al. [2012] concluded that it is not necessary to recover the exact ODE to recover its causal properties: if one has a cause and effect variable within an ODE system, then a qualitative description of the dynamics of the cause based on the dynamics of the effect can be recovered. Surprisingly, while one would need a large number of lags to estimate an ODE where no parametric assumptions have been made, Takens' theorem states that a good enough estimate, i.e. one that retains the causal properties of the ODE, can be made with at most $2d + 1$ lags where $d$ is the number of variables in the ODE.

---

**The CCM pseudo-algorithm for checking if two variables $X$ and $Y$ are causally related:**

1. Choose embedding dimension $E$: number of lags to use with $1 \leq E \leq 2d + 1$.
2. Estimate cross-map skill $\rho(l)$ for a sequence of several observations $l_1, \ldots, l_N$ with $l_N \leq L$, $L$ is the maximum number of available observations in the time series. For each $l_i$:
   (a) Construct shadow manifold $M_x$: in practice represented by matrix $Y \in R^{l \times E}$ with time series $y_t, y_{t-1}, y_{t-2}, y_{t-E+1}$
   (b) Assume the shadow manifold satisfies Takens' theorem condition and retains the metric properties of manifold $M$. Thus estimate euclidean distance $d_i$ of $E+1$ nearest points on manifold $M_x$ to point $(x_t, x_{t-1}, \ldots, x_{t-E+1})$. Denote the time indices corresponding to these points as $t_1, t_2, \ldots, t_{E+1}$.

(c) Construct estimate of $y_t$ using simplex projection in shadow manifold: weighted average of $E + 1$ nearest points (on $M_x$) with weights determined according to the exponentially weighted distance on $M_x$ of each point (calculated in the previous step):

$$\hat{y}_t = \sum_{i=1}^{E+1} w_i y_{t_i} \quad \text{where} \quad w_i = \frac{\exp(-\frac{d_i}{d_1})}{\sum_i \exp(-\frac{d_i}{d_1})} \tag{2.3}$$

(d) Construct cross-map skill $\rho(l) = Corr(y_t, \hat{y}_t | M_x)$

3. Check if cross-map skill $\rho(l)$ converges as $l$ tends to $L$. As the number of observations used increases, the manifold estimation should be denser, so cross-map skill should improve and converge, provided our assumption that $M_x$ retains the metric properties of $M$ is true.

The algorithm should also be applied symmetrically to establish the convergence of the cross-map skill $\hat{x}_t | M_y$. If both cross-map skills converge, we can establish that both variables belong in the same ODE system, and the causal relations are bi-directional. Note that in step 2c, we only use the shadow manifold $M_x$ to determine which points and with which weights should be used to estimate $y_t$. If the convergence of the cross-map skill happens in only one direction, the proper conclusion is that a uni-directional causal relationship exists between the two variables. If the cross-map skill of $\hat{y}_t | M_x$ converges, the proper conclusion is that $y$ causes $x$. This is somewhat counterintuitive, at least from the point of view of more classical causal discovery methods, because to establish that $x$ is an effect of the cause $y$, we must be able to predict the cause $y$ using the effect $x$, where for all other methods discussed in this work it is the other way around.

### PC-based methods.

In the following, we start with an exposition of the PC algorithm [Spirtes et al., 2000]. We then explain both a naive (tsPC) and more sophisticated (PCMCI) time series adaption.

1. *PC.* The original PC algorithm (named after Peter and Clark's authorship [Spirtes et al., 2000]) was constructed for i.i.d random variables and thus, in particular, for non-time series data. Below, we will also describe an extension to the time series case. The PC algorithm assumes that the underlying causal graph is a (*directed acyclic graph DAG*). A DAG has only directed edges ($\rightarrow$ and $\leftarrow$) and no cycles. As for independence-based algorithms in general, the PC algorithm assumes the causal Markov condition and causal faithfulness to infer $d$-separations on the causal graph from conditional independencies alone. Moreover, the algorithm assumes causal sufficiency. Consequently, the algorithm cannot distinguish between two graphs with the same set of $d$-separations. The graphical representation of an equivalence class of DAGs with the same $d$-separations is known as a (*completed partially-directed acyclic graph CPDAG*), which is the object of discovery of PC. As compared to DAGS, CPDAGs can contain undirected edges ($\circ\!-\!\circ$). These undirected edges signify that both orientations ($\rightarrow$ or $\leftarrow$) are compatible with the set of conditional independencies.

**The PC algorithm starts from a fully connected undirected graph and consists of three phases:**

(a) The *skeleton phase* uses statistical (conditional) independence tests to infer the adjacencies of the underlying causal graph. If two variables $X$ and $Y$ are found to be independent conditional on a (possibly empty) set of variables $\mathbf{Z}$, then the edge between $X$ and $Y$ is removed.

(b) The *collider orientation phase* then orients all *collider motifs*, that is, motifs of the form $X \rightarrow Y \leftarrow Z$ where $X$ and $Z$ are non-adjacent. These orientations can be inferred because collider motifs impose a particular pattern of (conditional) (in-)dependencies.

(c) The *orientation phase* finally uses graphical rules [Meek, 1995] to infer the orientation of as many remaining unoriented edges as possible using the acyclicity assumption and the fact that all colliders have been found in the previous step.

As compared to a brute-force execution of the skeleton phase as done in the so-called SGS

algorithm [Spirtes et al., 2000], the PC algorithm improves this phase as follows: Instead of searching for conditioning sets $\mathbf{Z}$ that renders a pair of variables $X$ and $Y$ independent (and thereby non-adjacent) within the set of all variables in the system, it is sufficient to restrict the search to those variables that are adjacent to $X$ and $Y$. The sufficiency of this specified search is a direct consequence of assuming the causal Markov condition.

The PC algorithm is fully non-parametric by construction and lends itself to various application cases. It performs optimally when the underlying causal graph is sparse and the number of variables is much less than the sample size. However, for sufficiently sparse graphs, PC was shown to be computationally feasible for very high-dimensional graphs as well [Kalisch and Bühlman, 2007]. In the worst-case scenario of very dense graphs, the number of conditional independence tests to be performed grows exponentially with the number of variables.

A variant of PC called *PC-stable* [Colombo et al., 2014] was shown to be *order-independent*, that is, invariant under permutations of the variables (this invariance is a desired property because it should not matter which variable is considered to the be "first" and which is considered to be the "last" variable). Another variant called for *Conservative PC* [Ramsey et al., 2006] is sound even if instead of the (full) causal faithfulness assumption, only the weaker *adjacency faithfulness* assumption is made. It is also possible to infer whether or not *orientation faithfulness* is violated [Ramsey et al., 2006].

Although the PC algorithm was originally developed for the acyclic case, the work [M. Mooij and Claassen, 2020] shows that PC is also consistent in the presence of cycles if the learned graph is interpreted in a slightly different way using the so-called $\sigma$-separation [Bongers et al., 2021]. The proof of this consistency uses the $\sigma$-separation [Bongers et al., 2021], which is a generalisation of $d$-separation to the cyclic case. Subsequently, it imposes a modified version of the Markov and faithfulness assumptions, namely the $\sigma$-Markov and the $\sigma$-faithfulness conditions. Under these conditions, the PC algorithm is proven to correctly learn the $\sigma$-separation equivalence class of the true (potentially cyclic) causal graph.

2. *tsPC.* The naive extension of the PC algorithm to the time series case is called tsPC, an example implementation is given in [Runge, 2020]. The general idea is to fix an integer $\tau_{\max}$ that is supposed to be equal to or larger than the maximum time-lag of any edge in the causal time series graph and to learn the finite segment of the time series graph on a time window $[t - \tau_{\max}, t]$. Here, $t$ is an arbitrary reference time step, and samples are created by sliding the time window over all recorded time steps. This approach implicitly assumes that the causal relationships do not change throughout the recorded time steps. Suppose the assumption of no contemporaneous causal influences is made. In that case, it is sufficient only to run the skeleton phase of PC because then the orientation of all edges is determined by time order (an effect cannot precede its cause). As discussed in [Runge, 2020], tsPC suffers from a sub-optimal finite-sample performance due to autocorrelation. This issue is remedied by the *PCMCI* algorithm, explained below. See Figure 2.2 for an illustration of the graphs learned by PC and tsPC for the time series case example.

3. *PCMCI.* The PCMCI algorithm [Runge et al., 2019b] is a time series causal discovery algorithm that addresses some of the shortcomings of the naive time series adaption of PC, in particular the issue of low detection power. PCMCI assumes the time series graph to have no contemporaneous causal influences. This assumption implies the absence of contemporaneous cycles, but feedback cycles involving time lags are possible. Additionally, PCMCI assumes the time series data are generated by a causally stationary process (that is, the causal relationships are assumed to not change over time).

As discussed in [Runge et al., 2019b, Runge, 2020], two main challenges with time series data hamper the performance of independence-based discovery algorithms in time series. These challenges are related to autocorrelation, a common feature in time series. First, using non-

i.i.d. samples (created in a sliding window fashion as explained above) typically leads to ill-calibrated conditional independence tests, that is, uncontrolled type I errors, because the degrees of freedom are reduced and cannot be easily measured. This ill-calibratedness leads to inflated false positives, that is, the discovery of dependence when, in fact, independence is true. Secondly, high autocorrelation implies that there is little new information in the next time step compared to the previous step. Depending on how the conditioning sets in conditional independence tests are selected, this results in low effect sizes leading to low detection power of true links. The effect size of a (conditional) independence test is defined as the absolute value of the population value of the test statistic; for example, in a (partial) correlation test, the effect size is the absolute value of the population value of the (partial) correlation. While there is a trade-off in addressing both of these challenges, the PCMCI algorithm and its generalisation PCMCI$^+$ (see below) algorithms remedy these challenges to an extent by using a particular choice of conditioning sets in the independent tests that decides about the presence versus the absence of an edge between a given pair of variables. We spell out the details below.

> **The PCMCI algorithm unfolds in two phases:**
>
> (a) The first phase, referred to as PC$_1$, is a *condition-selection phase* that aims to infer a superset $\hat{\mathscr{P}}(X_t^j)$ of the parents of each variable $X_t^j$ at time step $t$. The PC$_1$ algorithm is a variant of the PC and works as follows: Each sub-step of the skeleton phase is indexed by the integer $p$, starting at $p = 0$ and successively increasing $p$ in increments of one. Within each sub-step, the algorithm tests for independence of $X_{t-\tau}^i$ and $X_t^j$ given a conditioning set that consists of those $p$ potential parents of $X_t^j$ (less $X_{t-\tau}^i$) that have the highest association with $X_t^j$ according to the previous (conditional) independence tests. This particular choice of conditioning sets increases the effect sizes of the (conditional) independence tests, as can be understood information-theoretically [Runge, 2015]. A higher effect size leads to a higher statistical power (equivalently, to a lower probability of a type II error), that is, makes it more likely to detect dependence if dependence is true. However, since the effects of autocorrelation have not been dealt with yet, this phase of PCMCI is affected by the same issues as the naive time series extension of PC in terms of false positives.
>
> (b) The second phase conducts for each pair of variables $X_{t-\tau}^i$ and $X_t^j$ the so-called *momentary conditional independence (MCI) test* that tests the null hypothesis
>
> $$X_{t-\tau}^i \perp\!\!\!\perp X_t^j \mid \hat{\mathscr{P}}(X_t^j) \backslash \{X_{t-\tau}^i\}, \hat{\mathscr{P}}(X_{t-\tau}^i) .$$
>
> If this hypothesis is not rejected, the edge between $X_{t-\tau}^i$ and $X_t^j$ is removed. While the condition on $\hat{\mathscr{P}}(X_t^j)$ only would suffice to condition out confounded and indirect connections, the additional conditioning on $\hat{\mathscr{P}}(X_{t-\tau}^i)$ removes auto-dependencies from $X_{t-\tau}^i$ such that the conditional independence tests are well-calibrated and false positives are controlled at the desired level [Runge et al., 2019b]. Note that no orientation phase is required because, by the assumption of no contemporaneous causal influences, all edges are time-lagged and oriented by time order.

The numerical studies in [Runge et al., 2019b] show that in combination with the MCI tests in its second phase, PCMCI improves detection power and false positive control compared to the naive time series adaption.

4. *PCMCI$^+$*. The PCMCI$^+$ algorithm [Runge, 2020] extends PCMCI to the case where contemporaneous edges are allowed by suitably modifying the MCI phase (but still disallowing contemporaneous causal cycles and latent confounders). It consists of three phases: First, a PC$_1$ lagged phase. Second, an MCI contemporaneous phase. Third, an orientation phase. The first phase applies the PC$_1$ algorithm (see above) to the lagged edges. Its goal is to infer a superset of the *lagged* parents of each variable $X_t^j$. However, due to the potential presence of contemporaneous links, the first phase converges to a superset $\hat{\mathscr{P}}(X_t^j)$ of the lagged parents of $X_t^j$ plus the parents of the contemporaneous ancestors of $X_t^j$. The MCI contemporaneous phase is

initialised with all the links found in the previous phase and all possible contemporaneous links. It then conducts the conditional independence tests $X^i_{t-\tau} \perp\!\!\!\perp X^j_t \mid \mathbf{S}, \hat{\mathscr{P}}(X^j_t) \backslash \{X^i_{t-\tau}\}, \hat{\mathscr{P}}(X^i_{t-\tau})$, where, in addition to the lagged conditions of the MCI tests in PCMCI, $\mathbf{S}$ are sets of contemporaneous adjacencies of $X^j_t$ (and $X^i_{t-\tau}$ for $\tau = 0$). This removes all remaining spurious links due to contemporaneous confounding or indirect paths. Finally, the orientation phase orients all lagged links according to time order and then applies PC orientation rules [Meek, 1995] to orient as many contemporaneous edges as possible. See Figure 2.2 for an illustration.

### FCI-based methods.

A major success of causal discovery is the development of methods for learning causal relationships without assuming causal sufficiency. One famous example is the FCI algorithm [Spirtes et al., 1995, 2000, Zhang, 2008b]. We review the standard FCI algorithm and several of its time series adapations.

1. *FCI.* The FCI algorithm generalises the PC algorithm (see above) to the causally insufficient case. In addition to latent confounders, the algorithm can also deal with *selection variables*. These variables influence whether a given sample point belongs to the observed population. For example, a certain satellite observation might be more likely to be made if the cloud cover is not too dense. Consequently, the statistical dependence relations will be biased (giving it the name *selection bias*) as only one segment of the entire population of possible satellite observations is being considered. Below, we assume the absence of selection variables and explain the specialisation of FCI in this case.

> **Quick Introduction to Maximal Ancestral Graphs.**
>
> To deal with latent confounders (and selection variables), FCI works with a larger class of graphical models than PC does: Instead of DAGs, FCI works with *maximal ancestral graphs* (*MAGs*) [Richardson and Spirtes, 2002]. Maximal ancestral graphs can be interpreted as projections of the underlying DAG (which consists of observed variables, latent confounders, and selection variables) to a graph over the observed variables only. When assuming the absence of selection variables (as we do here), it is sufficient to work with a subset of MAGs that are called *directed maximal ancestral graphs* (*DMAGs*) [M. Mooij and Claassen, 2020]. Parts of the literature for notational simplicity do not explicitly distinguish between MAGs and DMAGs, that is, speak of MAGs although referring to DMAGs.
>
> Directed maximal ancestral graphs can have two types of edges: directed edges ($\rightarrow$) and bidirected edges ($\leftrightarrow$). A directed edge $X \rightarrow Y$ says that variable $X$ causally influences variable $Y$. This causal influence can be direct or indirect through one or multiple unobserved variables. A bidirected edge $X \leftrightarrow Y$ says that $X$ and $Y$ are subject to latent confounding and that, at the same time, neither $X$ causally influences $Y$ nor the other way around. Being subject to latent confounding means that there is an unobserved variable $L$ that (potentially indirectly through other unobserved variables) causally influences both $X$ and $Y$. In addition, an edge between $X$ and $Y$ (i.e., $X \rightarrow Y$ or $X \leftarrow Y$ or $X \leftrightarrow Y$) means that $X$ and $Y$ are not (conditionally) independent given any set of observed variables. A subtle part of the interpretation of DMAGs is that directed edges $X \rightarrow Y$ can "hide" latent confounding. That is to say, while $X \rightarrow Y$ does say that $X$ causally influences $Y$, it is possible that $X$ and $Y$ are also subject to latent confounding. Even if there is such additional latent confounding, then the DMAG does not, in addition, also contain the edge $X \leftrightarrow Y$ because there is at most one edge between any pair of variables. However, for a certain type of directed edges, called *visible* and which can be determined graphically from the DMAG, one can assert with certainty that there cannot be additional latent confounding [Zhang, 2008a].

The FCI algorithm works in a way that is similar to the PC algorithm: First, a sequence of (conditional) independence tests is performed to find the skeleton (that is, the adjacencies) of the graph. Second, several orientation rules are applied to determine the direction of as many links as possible. For these details, we refer to the original works [Spirtes et al., 1995, 2000, Zhang, 2008b] or to more technical reviews of FCI, for example, see Section S2 in the supplementary

material of [Gerhardus and Runge, 2020].

As in the case of the PC algorithm, FCI does not learn a unique DMAG but rather a Markov equivalence class of DMAGs. These equivalence classes are graphically represented by *directed partial ancestral graphs* (*DPAGs*) [Zhang, 2008a,b, M. Mooij and Claassen, 2020]. In addition to directed ($\rightarrow$) and bidirected edges ($\leftrightarrow$), DPAGs can also contain edges of the types $X \circ \rightarrow Y$ and $X \circ - \circ Y$. An edge $X \circ \rightarrow Y$ says that $Y$ does not have a causal influence on $X$ while $X$ might or might not have a causal influence on $Y$, whereas an edge $X \circ - \circ Y$ does not make any claim about whether or not $X$ or $Y$ have a causal influence on each other. See Figure 2.2 for an illustration of the graph that the FCI algorithm discovers when applied to time series data. The work [M. Mooij and Claassen, 2020] has shown that FCI, originally developed with the assumption of acyclicity, can also be consistently applied to data that is generated by a cyclic SCM with certain regularity conditions.

2. *tsFCI.* The tsFCI algorithm [Entner and Hoyer, 2010] adapts FCI to causally stationary time series. Here, the term *causal stationarity* means that the causal relationship between the variables $X_{t-\tau}^i$ and $X_t^j$ is the same as the causal relationships between the variables $X_{s-\tau}^i$ and $X_s^j$ for any time steps $t$ and $s$. In other words, the causal relationships are assumed invariant in time. As compared to the FCI algorithm, tsFCI applies the following two conceptual modifications: First, lagged links ($\tau \geq 1$) are by default oriented as $X_{t-\tau}^i \circ \rightarrow X_t^j$. These default orientations are valid because an effect cannot precede its cause. Note that it would not be valid to orient all lagged links as $X_{t-\tau}^i \rightarrow X_t^j$ because $X_{t-\tau}^i \leftrightarrow X_t^j$ (i.e., latent confounding) is a possibility. Second, so-called *homologous* edges are by default oriented in the same way. That is if the edge between $X_{t-\tau}^i$ and $X_t^j$ has been found to have a certain orientation (for example, $X_{t-\tau}^i \rightarrow X_t^j$ or $X_{t-\tau}^i \leftarrow X_t^j$) and if in addition there is an edge between $X_{s-\tau}^i$ and $X_s^j$ for $s \neq t$, then this latter edge is immediately oriented in the same way as the former edge (for example, oriented as $X_{s-\tau}^i \rightarrow X_s^j$ or $X_{s-\tau}^i \leftarrow X_s^j$). This copying of edge orientations is valid because of causal stationarity. In addition to these modifications, tsFCI uses the knowledge of time order and causal stationarity to apply further modifications that are useful from a computational and/or statistical point of view. There are two versions of tsFCI, both of which have been introduced in the original work [Entner and Hoyer, 2010]: One version does not allow for contemporaneous causal influences in the data-generating process and another version in which such influences are allowed. Figure 2.2 shows an illustration of the output of tsFCI (its version that allows contemporaneous causal influences) and other FCI-based time series causal discovery algorithms (see below) in the case of an example time series graph.

In the graph learned by tsFCI (both variants), there can be an edge between the pair of variables $X_{s-\tau}^i$ and $X_s^j$ (for example, $X_{s-\tau}^i \leftrightarrow X_s^j$ or, only if $\tau = 0$, $X_s^i \leftarrow X_s^j$) and at the same time no edge between the pair of variables $X_{t-\tau}^i$ and $X_t^j$ for $t > s$. At first sight, this non-repetition of edges seems to contradict the assumption of a causally stationary time series. However, there is, in fact, no contradiction: The "additional" edge between $X_{s-\tau}^i$ and $X_s^j$ can result from *temporal confounding*, that is, result from confounding by variables that are before the observed time window $[t - \tau_{\max}, t]$.

3. *SVAR-FCI.* The SVAR-FCI algorithm [Malinsky and Spirtes, 2018] is another adaption of FCI to causally stationary time series. As compared to tsFCI (variant with contemporaneous causation), SVAR-FCI removes the "additional" edges that have been discussed in the last paragraph in the section on tsFCI. More specifically, whenever SVAR-FCI detects marginal or conditional independence of $X_{t-\tau}^i$ and $X_t^j$ and hence removes the edge between these two variables, then the algorithm automatically and immediately also removes the edges between the variables $X_{s-\tau}^i$ and $X_s^j$ for all time steps $s$. Consequently, in the graph learned by SVAR-FCI, the variables $X_{t-\tau}^i$ and $X_t^j$ are adjacent if and only if for all $s$ the variables $X_{s-\tau}^i$ and $X_s^j$ are adjacent.

The removal of the additional edges is justified by the fact that these edges are known to result

from confounding variables that are before the observed time window. Hence, these edges would also disappear in tsFCI (both variants) for an appropriately increased length of the observed time window. Moreover, due to this modification, SVAR-FCI requires a smaller number of (conditional) independence tests than tsFCI. However, as has been realised and explained in [Gerhardus, 2021], SVAR-FCI, in theory, discovers fewer edge orientations than tsFCI (variant with contemporaneous causation). To the authors' knowledge, an empirical comparison of tsFCI (variant with contemporaneous causation) and SVAR-FCI on finite samples that would shed light on this trade-off has not yet been performed.

4. *SVAR-GFCI.* The SVAR-GFCI algorithm [Malinsky and Spirtes, 2018] is another adaption of FCI to causally stationary time series. This algorithm is a hybrid method that combines independence-based and score-based causal discovery. In the first step, the algorithm employs a time series variant of the score-based GES algorithm (see below for an explanation of the GES algorithm). The adjacencies found in this first step are then passed as starting conditions to the independence-based SVAR-FCI algorithm.

5. *LPCMCI.* The LPCMCI algorithm [Gerhardus and Runge, 2020], also known as Latent-PCMCI, is another adaption of FCI to causally stationary time series. As compared to SVAR-FCI, Latent-PCMCI applies the ideas of PCMCI and PCMCI$^+$ (see above) generalised to the causally insufficient case. Simulation studies in [Gerhardus and Runge, 2020] show that Latent-PCMCI strongly outperforms SVAR-FCI on finite samples, especially for strongly autocorrelated and continuously valued variables.

6. *Learning extended summary graphs.* Extended summary graphs [Assaad et al., 2022a] are time-compressed representations of time series graphs that distinguish between contemporaneous and lagged links (so they are more informative than summary graphs) but that do not distinguish between two different non-zero lags (so they are less informative than time series graph), see Sec. 2.1.2 above for more details. To learn extended summary graphs, [Assaad et al., 2022a] propose the **PCGCE** and **FCIGCE** algorithms. The basic idea of both these algorithms is to group the past time steps of each component time series $X^i$ into a vector variable $X^{i,-}$ and then, using conditional independence tests that can handle vectors of variables, respectively apply the PC and FCI algorithms directly to the extended summary graph. Specifically, [Assaad et al., 2022a] use a conditional independence test based on conditional mutual information. Moreover, for practical reasons, the vector variable $X^{i,-}$ does not contain *all* past time steps of $X^i$ but all those time steps within the time window $[t - \gamma, t - 1]$ where $\infty > \gamma \geq 1$ is a user-specified positive integer that corresponds to the maximal lag up to which the method is supposed to be sensitive (this parameter $\gamma$ corresponds to the parameter $\tau_{max}$ in PCMCI, for example).

**Greedy Equivalence Search (GES).**

GES [Meek, 1997, Chickering, 2002b,a] is probably the most famous score-based causal discovery method for i.i.d data assuming that the true causal graph is a DAG. It performs greedy steps directly on the CPDAG, thus searching in the Markov equivalent class space. Chickering [Chickering, 2002b] proved that an efficient two-phase greedy search, combined with the BIC score is sufficient to find the true CPDAG in the large sample limit (the Meek conjecture [Meek, 1997]) assuming causal sufficiency. Ramsey *et al.* [Ramsey et al., 2017] developed Fast GES (FGES) an optimised and parallelised version of the GES algorithm, which they were able to scale up to a million variables. Additionally, GES has been improved by bounding polynomially the score evaluations [Chickering and Meek, 2015]; to obtain statistical efficiency [Chickering, 2020]; to obtain finite-sample correction of confidence intervals [Gradu et al., 2022] and to deal with latent variables [Claassen and Bucur, 2022]. Similarly to other methods developed for i.i.d. data, (F)GES can be applied to uniformly-sampled causally stationary time series by simply considering the transformed lagged variables and imposing that the effects cannot precede the causes in time.

**Continuous optimisation methods.**

A recent advance in structure learning has been the development of so-called continuous optimisation methods, particularly score-based methods which avoid the explosion of the discrete space of DAGs by employing continuous optimisation. The first such method is the NOTEARS algorithm proposed by Zheng *et al.* [Zheng et al., 2018] which proposed $h(\mathbf{W}) = \text{tr}(\exp(\mathbf{W} \circ \mathbf{W}) - d$ as differentiable characterisation of acyclicity for a weighted adjacency matrix $\mathbf{W}$, that is $h(\mathbf{W}) = 0$ if and only if the associated graph is a DAG. Such differentiable characterisations allow the plug-in use of different continuous optimisation methods and even complex function parameterisations [Zheng et al., 2020, Ng et al., 2022, Lachapelle et al., 2020, Ng et al., 2020, Yu et al., 2021, Bello et al., 2022]. Various implementations of the continuous optimisation framework for structure learning from time series are available: (1) DYNOTEARS [Pamfil et al., 2020] considers structural linear VAR (SVAR) models, allowing for contemporaneous links and enforcing acyclicity among the instantaneous edges. (2) IDYNO [Gao et al., 2022] is designed to perform structural discovery from both observational and interventional data. Moreover, both linear and non-linear relationships are considered. (3) NTS-NOTEARS [Sun et al., 2023] models cause-effect relationships through one-dimensional convolution neural networks (CNN) and allows prior knowledge to be encoded and exploited directly by the optimisation procedure.

**(VAR)LiNGAM.**

Linear non-Gaussian acyclic model (LiNGAM) [Shimizu et al., 2006] is a classical method for causal discovery which assume acyclicity, causal sufficiency, linear relationships and non-Gaussian additive independent noises. Under those assumptions, the model is shown to be identifiable thanks to classical results from independent component analysis (ICA) [Hyvärinen et al., 2004]. Specifically, a linear SEM can be represented by,

$$\boldsymbol{X} = \mathbf{B}\boldsymbol{X} + \boldsymbol{\varepsilon}, \tag{2.4}$$

where $\boldsymbol{X}$ is the vector of system variables, $\boldsymbol{\varepsilon}$ the noise vector and $\mathbf{B}$ is the matrix of coefficients. Thus solving Equation 2.4 for $\boldsymbol{X}$ we obtain,

$$\boldsymbol{X} = (\mathbf{I} - \mathbf{B})^{-1}\boldsymbol{\varepsilon}.$$

The above equation is the an ICA problem, and its theory states that when $\boldsymbol{\varepsilon}$ are non-Gaussian noises variables, the mixing matrix $\mathbf{A} = (\mathbf{I} - \mathbf{B})^{-1}$ is identifiable in the large sample limit [Hyvärinen et al., 2004]. LiNGAM-ICA works by first exploiting the ICA results to obtain the mixing matrix $\mathbf{A}$, and secondly by permuting and normalising $\mathbf{A}$ to obtain the appropriate matrix $\mathbf{B}$ and the corresponding causal DAG. The LiNGAM can also be solved directly without using ICA [Shimizu et al., 2011]. The direct-LiNGAM iteratively selects the variable which is the most independent from the residual (in a linear regression onto the remaining variables) and replaces the original data with the residuals matrix. In the end, the procedure gives a causal order; and ultimately, simple linear regressions (e.g. estimated with least squares) are used to obtain the final estimation of the lower triangular $\mathbf{B}$ matrix.

VAR-LiNGAM [Hyvärinen et al., 2008, Zhang, 2008a, Hyvärinen et al., 2010] is the extension of the classical LiNGAM model to time series data. It considers a linear structural vector autoregressive (VAR) model, with possibly acyclic instantaneous relationships, and assumes non-Gaussian disturbances. In particular, the process $\boldsymbol{X}_t$ is assumed to evolve following

$$\boldsymbol{X}_t = \sum_{\tau=0}^{k} \mathbf{B}_\tau \boldsymbol{X}_{t-\tau} + \boldsymbol{\varepsilon}_t.$$

Where $\mathbf{B}_0$ is the matrix of instantaneous relationships and its sparsity pattern corresponds to a DAG while $\mathbf{B}_\tau$ for $\tau > 0$ are the matrices of lagged relationships; moreover the non-Gaussian noise vector

$\varepsilon_t$ has independent components and is assumed independent over time. Hyvärinen et al. [2010] propose two methods to estimate the coefficients of the model above: (1) a two-stage method which combines least-square estimation of the autoregressive model and classical LiNGAM estimation; (2) a method based on multichannel blind deconvolution. In Figure 2.2, we illustrate that if the assumptions are satisfied, VAR-LiNGAM is able to discover and orient all causal links.

**Structural equation models for time series, TiMiNo.**
Peters et al. [2013] extended the classical structural equation modelling framework to the time series setting by introducing time series models with independent noise (TiMINo). In particular, a multi-variate time series $X_t = (X_t^i)_{i \in V}$ satisfies a TiMINo if there exists a $p > 0$ and, for every $i \in V$, there are subsets $\mathbf{PA}_0^i \subseteq X^{V \setminus \{i\}}$ and $\mathbf{PA}_k^i \subseteq X^V$ for $k = 1, \ldots, p$ such that,

$$X_t^i = f_i \left( (\mathbf{PA}_p^i)_{t-p}, \ldots, (\mathbf{PA}_1^i)_{t-1}, (\mathbf{PA}_0^i)_t, N_t^i \right), \tag{2.5}$$

where $N_t^i$ are assumed to be jointly independent over $i$ and time and for each $i$, $N_t^i$ are identically distributed in time. Thus, the assumed model is extremely general, instantaneous relationships are allowed and the noise contribution is not restricted by allowing arbitrarily mixing through the $f_i$ functions. Nevertheless, to prove the identifiability of the full-time graph and the causal summary time graph, additional assumptions have to be made. Specifically, in Peters et al. [2013] it is assumed that either (i) Equations 2.5 follows an identifiable functional model class (e.g. nonlinear functions with additive Gaussian noise or linear functions with additive non-Gaussian noise [Peters et al., 2011]) or (ii) each function $f_i$ exhibits a time structure, that is the union of the causal parents of $X_t^i$ contains at least one $X_{t-k}^i$ and moreover the joint distribution is faithful with respect to the full-time graph with the summary time graph being acyclic. Under either assumption (i) or (ii), the complete causal graph is proved to be identifiable in the large sample limit. Practically, to estimate a TiMINo, methods for causal discovery for additive noise models with i.i.d. data [Mooij et al., 2009] can be adapted to the time series setting. As in the direct-LiNGAM, it iteratively selects the causal order between the variables by fitting regression models and evaluating the independence between the variables and the residuals. To fit the regression models for $f_i$, various methods can be used, such as vector autoregressive models (linear), generalised additive models or Gaussian processes. Moreover, to test independence from the residuals, the HSIC [Gretton et al., 2007] can be applied to all possible shifted time series up to the maximum lag order.

**Invariant Causal Prediction.**
Invariant causal prediction deals with the setting of independent and identically distributed samples of the random vectors $X = (X_1, X_2, \ldots, X_p)^\top \in \mathbb{R}^p$, $E = (E_1, E_2, \ldots, E_q)^\top \in \mathbb{R}^q$, and $Y \in \mathbb{R}$. For a variable or vector of interest $Y$, $E$ is a set of environment variables that may be causes of $X$ but are not direct causes or effects (direct or indirect) of $Y$. ICP assumes that $Y$ is generated *causally* from a subset $S^* \subseteq \{1, \ldots, p\}$ of the $p$ variables considered, so that there is causal sufficiency, and $Y$ is generated from a Structural Causal Model obeying:

$$Y = g(X_{S^*}, \varepsilon), \quad \varepsilon \sim F, \quad \varepsilon \perp\!\!\!\perp X_{S^*}, \tag{2.6}$$

where $g$ and $F$ are arbitrary functions and distributions, respectively. A set $S \subseteq \{1, \ldots, p\}$ is a generic subset of the full set of candidate causes. We refer to $S$ and $X_S$ interchangeably for brevity. The task of ICP is to infer the set $S^*$ of direct causes of the variable of interest $Y$.

The Invariance Causal Prediction (ICP) framework [Peters et al., 2016] is based on the observation that $Y$ is independent of $E$ given $X_{S^*}$, denoted $Y \perp\!\!\!\perp E | X_{S^*}$. Assuming we have a set of candidate causes $X$ that includes $S^*$, the causal subset, and an environment variable $E$ that we know does not directly cause $Y$ or is an effect of $Y$. We can search for $S^*$ by applying a conditional independence test $Y \perp\!\!\!\perp E | X_S$ on $Y$, $E$ and subsets $S \subseteq \{1, \ldots, p\}$ of $X$. ICP then selects as causal variables the

intersection of all those subsets $S$ where the corresponding conditional independence test is not rejected:

$$\hat{S}^* = \bigcap_{S:p_S > \alpha} S. \tag{2.7}$$

Here $p_S$ is the $p$-value associated with the conditional independence test of $Y \perp\!\!\!\perp E | X_S$, with the null hypothesis corresponding to conditional independence. We do not reject conditional independence at significance level $\alpha$ if $p_S > \alpha$.

One way to interpret the problem setting is that causal associations are more robust than other associations. So ICP finds the causes of a variable of interest by investigating which associations are invariant across environments. Another interpretation is that the environment variables define data generated under different interventions to the system (SCM).

In physical sciences, regional and temporal variables are good candidates since these often describe changes in environments that alter the conditions under which physical processes occur.

Peters et al. [2016] introduce an algorithm to implement ICP that assumes linear relationships between causes $X$ and effects $Y$ and a categorical, univariate variable $E$. In [Heinze-Deml et al., 2018a], more general algorithms are presented that allow for nonlinear relationships between cause and effects and for continuous and multivariate environment variables $E$. Pfister et al. [2019] provide an ICP variant for time series data. The proposed time-series ICP relies on the causal invariance assumption across time points, thus removing the requirement of environment knowledge. This setting is also partially robust to hidden confounders, similar to the original ICP [Peters et al., 2016] framework, and in general, the ICP is expected to be conservative with respect to violations of its assumptions [Pfister et al., 2019].

**Causal frameworks for continuous-time systems.**

As we already reviewed, discrete-time causal systems fit directly as an extension of i.i.d. and DAG framework. When considering discrete-time systems, we can express the value of a variable at a time $t$ as a function of other variables (and itself) observed at past instants, thus the complete causal graph can be seen as a DAG extended (possibly infinitely) in time. Classical methods for causal discovery in the i.i.d. setting can then be adapted to discrete-time systems quite straightforwardly.

Conversely, considering continuous-time systems, raise the issue that a time-extended DAG is not feasible, since the included variables would be uncountable [Hansen and Sokol, 2014]. Nevertheless, modelling continuous dynamical system helps in dealing with non-uniform sampling and extrapolating among different sampling frequencies, two major drawbacks of causal discovery methods for discrete-time systems. The following are the major causal frameworks available for continuous-time systems:

1. *Causal interpretation of ODE and SDE.* Various efforts have been made to describe causal systems with ODEs and SDEs. First, causal discourses around ODE were used to obtain different justifications for the cyclic SEM [Hyttinen et al., 2012, Mooij et al., 2013a, Bongers et al., 2018, Peters et al., 2022]. Rubenstein et al. [2018] described Dynamical Structural Causal Models (DSCM) as extensions of SEM where each equation or assignment is a relationship between a set of causal parent trajectories and an effect trajectory. Under some stability conditions, such DSCM can be obtained from ODE systems. SDEs have also been studied from a causal perspective [Hansen and Sokol, 2014, Peters et al., 2022]. The advantage of SDEs in modelling physical systems is that they allow incorporating an inherent source of stochasticity, a common assumption in numerous real-world systems [Abbati et al., 2019]. Graphical parameterisations of SDE equilibrium distributions leading to models allowing for cycles have been investigated and different structure learning algorithms have been proposed [Varando and Hansen, 2020].

2. *Dynamic Causal Models.* Dynamic Causal Models (DCM) [Friston et al., 2003, Marreiros et al., 2010, Stephan et al., 2010, Abbati et al., 2019, Friston et al., 2013], in short, is a Bayesian framework for fitting and comparing causal models for coupled dynamical systems. DCM was introduced and applied mostly in Neuroscience, and particularly in the problem of estimating the connectivity between brain regions from neuroimaging data, as discussed in detail in Section 4.2.2. Recently, and it has been even employed to model the COVID-19 pandemic dynamics [Friston et al., 2020, 2022].

3. *Local independence graphs.* Local independence is a notion of conditional independence for stochastic processes (both discrete and continuous time ones) which can be (in)dependent on each other pasts. In detail, for real-valued stochastic processes $X_t = (X_t^1, \ldots, X_t^p)$ and $A, B, C \subseteq \{1, \ldots, p\}$), we say that $X^B$ is locally independent of $X^A$ given $X^C$ at time $t$ if the past of $X^C$ until time $t$ provides the same information, to predict $E[X_t^\beta | \mathscr{F}_t^{A \cup C}]$ [1], as the past of $X^{A \cup C}$ until time $t$, for each $\beta \in B$.

Didelez [Didelez, 2008, 2006] studied graphical representations of local independence with directed graphs together with $\delta$-separation and proved the equivalence of the pairwise and global Markov properties for multivariate counting processes. Directed graphs and $\delta$-separation has been then extended to mixed graphs and $\mu$-separation [Mogensen and Hansen, 2020] to model partially unobserved systems. A constrained-based algorithm has been proposed [Mogensen et al., 2018], which is proven to be sound and complete under faithfulness assumption. Local independence graphs can be applied to multivariate processes which are solutions of SDEs (such as the multivariate Ornstein-Uhlenbeck process) or event and counting processes such as Hawkes processes [Mogensen, 2022].

## 2.2 Challenges

In this section, we discuss several challenges for causal discovery that are frequent in real-world applications. Following Runge et al. [2019a], we distinguish between challenges related to the data-generating process itself (see Sec. 2.2.1), challenges associated with the available data (see Sec. 2.2.2), and challenges of statistical or computational nature (see Sec. 2.2.3). Users should carefully consider the challenges they face in their application and choose a suitable causal discovery method. More generally, how to reason and deal with typical challenges in causal inference is further discussed in a time series context in Runge et al. [2023].

### 2.2.1 Process challenges

Non-linearities pose challenges for causal discovery in both independence-based and score-based methods. Non-linear conditional independence tests, such as the GPDC test [Rasmussen and Williams, 2006] and tests based on conditional mutual information [Runge, 2018], are computationally more expensive and tend to have lower sta-

> *"Discovering causal relations from observational data is impossible without assumptions about the mechanisms and faces important challenges related to data and statistical characteristics. The field will need to incorporate domain knowledge and post-selection inference."*

tistical power than linear tests. Non-linear functional relationships in score-based methods require more complex score functions, which can decrease finite-sample performance. However, non-linear functional relationships can enhance identifiability in some asymmetry-based causal discovery methods [Peters et al., 2017a]. Overall, nonlinearities increase model complexity and require careful consideration when selecting appropriate causal discovery methods.

---

[1] where $\mathscr{F}_t^{A \cup C}$ is a right-continuous and complete filtration which represents the history of the processes $X^{A \cup C}$ see e.g. [Mogensen and Hansen, 2020] for a detailed description

Most time series causal discovery methods assume data generated from a causally stationary process with a unique equilibrium distribution. However, many real-world processes are non-stationary. Recently there have been works on devising causal discovery methods that first detect the variables afflicted with non-stationary driving mechanisms and subsequently infer the entire causal graph, including possibly proxy variables corresponding to the driving force of non-stationarity [Mooij et al., 2020, Huang et al., 2020]. Furthermore, techniques exist to detect regime or context changes in non-stationary data [Huang et al., 2015, Saggioro et al., 2020, Huang et al., 2020], but it remains an an active area of research. Domain experts may be able to identify the source of non-stationarity in time series data and preprocess the data to remove it.

Time series data is common in physical sciences, and it has a distinctive feature of auto-correlation. Many causal discovery algorithms are not designed for time series data and show decreased performance when applied without modification [Runge et al., 2019b, Runge, 2020]. However, some methods are specifically designed for time series data, reducing the detrimental effect of auto-correlation. See the PCMCI algorithm above for a discussion on the challenges of auto-correlated data. In many domains, space adds to time as well. In principle, one could feed different spatial locations of the same variable as distinct variables into causal discovery methods. However, this naive approach ignores spatial correlations and quickly results in a high-dimensional problem. Another workaround, employed for example in [Tibau et al., 2022, Diego Bueso, 2020], is to perform dimension-reduction as a preprocessing step and then perform causal discovery on the dimensionally-reduced space. The development of causal discovery inherently designed for spatio-temporal data is an active area of research; for example, see [Christiansen et al., 2022]. Dealing with variables whose dynamics operate on different time scales, such as e.g. fast atmospheric and slow oceanic processes, pose important challenges. Granger causality in the frequency domain, e.g. [Bressler and Seth, 2011, Chicharro, 2011, Faes et al., 2012], and a combination of wavelet analysis with transfer entropy [Lungarella et al., 2007] are examples of approaches to deal with the time scales of causal influences.

Finally, it is worth mentioning that many causal discovery methods assume the data-generating SCM to be acyclic (see axis VII in Sec. 2.1.2). However, in real-world applications, one can often not exclude the existence of feedback that acts on time scales below the measurement interval. Such feedbacks make the time series graph **cyclic**. In the non-temporal setting, one might often not be able to exclude the existence of causal cycles. As discussed and referenced above, there are causal discovery methods that can handle cyclic causal graphs.

## 2.2.2  Data challenges

Various data challenges arise when tackling the problem of causal discovery in practice. This is mainly due to the discrepancy between the assumed hypothesis needed from each method or framework and the real-world data. One of the most common assumptions of most causal discovery methods is causal sufficiency, which is the hypothesis that all relevant variables are observed. Unobserved variables are especially problematic when they are possible confounders between system variables since omitting confounders from the causal discovery could lead to learning spurious or wrong relationships. There are some available methods which do not assume causal sufficiency, such as LPCMCI [Gerhardus and Runge, 2020], SVAR-FCI [Malinsky and Spirtes, 2018] and GPS [Claassen and Bucur, 2022] (see Sec. 2.1.3).

Missing data and selection bias are other common issues in real-world applications, and there have been some efforts in developing causal discovery methods which are resilient to these challenges [Strobl et al., 2018, Gain and Shpitser, 2018, Tu et al., 2019, Versteeg et al., 2022]. Not always, especially in the physical realm, the data follow predictable and well-behaved distributions. One such example is zero-inflated data, which is common, for instance, in gene expression data, where single-cell expressions lack detectable values of transcripts that appear abundant on bulk

(thousands of cells) gene expression experiments. Recent advances have developed graphical models and causal discovery methods in such scenarios [McDavid et al., 2019, Yu et al., 2020].

### 2.2.3 Statistical and computational challenges

The high dimensionality of data in physical systems, such as spatiotemporal data, and small sample sizes are central statistical challenges for causal discovery. On the other hand, large sample sizes raise issues of unaffordable computational time, which can scale up to cubically for kernel methods typically used for independence testing [Schölkopf and Smola, 2008, Gretton et al., 2007]. High-dimensional data leads to large conditioning sets in particular algorithms, effectively reducing the sample size available to test the hypothesis.

As noted in Sec. 2.2.1, non-linearity is a common characteristic of processes in the physical sciences. For the case of independence-based or hybrid casual discovery techniques, this calls for devising non-parametric tests of independence, for instance, tests based on measures of conditional mutual information [Runge, 2018], or on Gaussian process regression or other kernel-based measures on independence [Gretton et al., 2007], or using quantile regression [Petersen and Hansen, 2021] and copula-based methods [Bouezmarni et al., 2012] (also applied to Granger causality), etc. The no-free-lunch theorem of Shah and Peters [2020] states that no single conditional independence test can have power against all alternatives. Here, the challenge lies in devising and applying (a combination of) conditional independence tests that are the most suited for a particular physical system.

The concept of post-selection inference [Berk et al., 2013] involves performing statistical inference on a model that was selected based on data-driven methods rather than being pre-selected. While there are some advances in solving the post-selection inference problem for regression and causal effect estimation, few solutions have been proposed for the inference after causal discovery setting [Berk et al., 2013, Belloni et al., 2014, Rinaldo et al., 2019]. One possible solution is sample-splitting, but this is often statistically inefficient. A recent development is the randomised version of the greedy equivalence search (GES) algorithm, which allows for finite-sample correction of classical confidence intervals [Gradu et al., 2022].

## 2.3 Opportunities for the physical sciences

The field of causal discovery from observational data is still in its infancy but growing in methodologies, theoretical guarantees of performance, and empirical evidence. Causal inference, in general, is a vast field that offers alternatives and scientific opportunities that we review in what follows. Only revisiting the whole body of empirical science based on association would take a village, but the advances would pay off.

### 2.3.1 Causal hypothesis testing and targeted interventions

Scientists need a principled way to test different hypotheses against each other. A causal hypothesis is a supposition or theory about how things interact, specifically on whether one thing causes another. Causal studies aim to confirm or reject any given causal hypothesis. The problem is that hypotheses in the physical sciences are often presented as narratives giving a chain of causal factors that lead to the studied phenomenon.

> *"Observational causal discovery offers revolutionary opportunities to test hypotheses, evaluate the impact of interventions, attribute extreme events with counterfactuals, and characterise complex systems by deriving causal pathways and robust forecasting models."*

Without a causal vocabulary and analytical tools, it is often impossible to precisely state the hypothesis, which leads to several competing hypotheses or, even worse, a false hypothesis, which is

accepted as true due to its compelling narrative quality. Testing hypotheses have been conducted in myriad ways [Pearl, 2009c, Peters et al., 2017a, Robins and Hernan, 2020].

Causal graphs, as graphical representations of assumed or learned causal relations, provide a more principled way to talk about causal hypotheses. Learned graphs imply causal links and pathways and provide evidence for deciding between rivalling causal hypotheses, in Kretschmer et al. [2016], for instance, regarding competing hypotheses of Arctic climate teleconnections.

The conclusiveness and interpretability of discovered causal graphs from purely observational data sets on the often untestable validity of the methods' assumptions and the statistical complexity of the task. But observational causal discovery can help inform more targeted subsequent interventions, which are often too expensive to employ on a large scale [Pearl, 2009c, Robins and Hernan, 2020]. Incorporating interventions, if performed meaningfully, could thus make the causal discovery process much more efficient and robust (discovered DAGs not being confined to the Markov equivalence class). Interestingly, interventions could be differentiable, i.e., 'learnable' from data [Brouillard et al., 2020].

### 2.3.2 Cause-effect estimation

Causal discovery results in qualitative causal graphs, or often Markov equivalence classes of graphs. But often, the target question is a quantitative estimate of a causal effect of one variable $X$ on another variable $Y$, as pioneered by Pearl [2009c]. This topic is discussed in a time series context in Runge et al. [2023]. The quantity of interest then is the (interventional) distribution of $Y$ given an intervention in $X$, $p(Y = y \mid do(X = x))$. The fundamental problem is that typically $p(Y = y \mid do(X = x)) \neq p(Y = y \mid X = x)$. Confounders, for example, can introduce a non-causal association between the treatments and the outcome. Randomised experiments would be the gold standard by eliminating the unwanted non-causal associations [Fisher, 1935, Imbens and Rubin, 2015, Pearl, 2009c]. The goal of causal effect estimation is to do so without access to interventions by expressing $p(Y = y \mid do(X = x))$ as a function of the observational distribution $p(\mathbf{x})$:

$$p(Y = y \mid do(X = x)) = \text{function of } p(\mathbf{x}). \tag{2.8}$$

If such a re-expression is possible, one calls the causal effect identifiable and obtains a causal estimand, which involves only the observational distribution. The most well-known method for causal effect estimation from data without parametric assumptions is *covariate adjustment* [Pearl, 2009c], which refers to de-confounding the causal relationship by adjusting for a set of variables $\mathbf{Z}$. In the general case, the adjustment formula is

$$p(y \mid do(X = x)) = \int p(y \mid x, \mathbf{z}) p(\mathbf{z}) d\mathbf{z}. \tag{2.9}$$

Recent work has focused on finding statistically optimal adjustment sets [Runge, 2021], i.e., for which the estimators have minimal variance. Using the *do*-calculus [Pearl, 1995, Huang and Valtorta, 2006, Shpitser and Pearl, 2006, 2008], it is possible to determine whether a causal effect is, in principle, identifiable from observational data or not. To this end, causal effect estimation requires fully specified causal graphs from causal discovery (with its inherent reliance on further assumptions) or domain expertise that can qualitatively specify a causal graph. For example, it is known that temperature influences ecosystem respiration, but one may want to quantify how much when given a graph of other observed and unobserved confounding variables. The graph then encodes assumptions about the absence and presence of causal relations.

Different variants of causal effects can be defined based on the interventional distribution $p(y \mid do(X = x))$, and an estimate then involves further parametric assumptions. For example, in a linear model, the total causal effect on the expected value of $Y$ when setting $X$ by intervention to $x'$ as opposed to $x$ is given by

$$\Delta_{X \to Y}(x', x) = \Delta x \cdot \beta_{X \to Y}, \tag{2.10}$$

where $\Delta x = x' - x$ and $\beta_{X \to Y}$ can be estimated as the regression parameter of $X$ in the linear regression of $Y$ on $X \cup \mathbf{Z}$.

### 2.3.3 Causal pathway analysis and mediation

Next to quantifying the overall causal effect of $X$ on $Y$, a relevant follow-up question is often about the causal pathways: the mechanisms by which this effect propagates. In complex systems, it is often interesting to analyse how perturbations spread throughout the systems and through which subprocesses perturbations are mediated [Runge et al., 2015b, Runge, 2015]. Within the structural causal model framework, mediation formally leads to counterfactual quantities, see, for example, VanderWeele [2015] and briefly below in Sec. 2.3.8. But for linear models, the mediated causal effect (MCE) of $X$ on $Y$ that passes through a mediator $M$ (here $X, Y, M \in \mathbf{V}$) can be computed by summing up the contributions along all paths passing through it:

$$\mathrm{MCE}(X, Y | M) = \sum_{\pi_k^M} \prod_{\lambda_{i \to j} \in \pi_k^M} \beta_{i \to j}, \qquad (2.11)$$

where the summand iterates over causal paths $\pi_k^M$ from $X$ to $Y$ through $M$ and the product is over all links $\lambda_{i \to j}$ on each path. The link coefficient $\beta_{i \to j}$ can be estimated as the regression coefficient of $V_t^j$ in the linear regression of $V_t^j$ on the parents of $V_t^j$. Mediation analysis can also answer the complementary question: how strong is the direct effect of $X$ on $Y$.

### 2.3.4 Identifying causes and pathways leading to anomalies

Anomaly detection [Chandola et al., 2009] is an important problem in engineering and the physical sciences. In Earth sciences, extreme events form a subclass of anomalies and can be structured across different dimensions, such as compound extremes [Zscheischler et al., 2020]. While detecting anomalies is an important problem, it does not answer the often relevant question of what causes a particular anomaly or, more generally, what causes the anomalous process. In engineering, business science, and healthcare, a related problem is *root cause analysis* [Andersen and Fagerhaug, 2006]. Causal discovery can address the problem of identifying causal drivers (parents) or indirect mediating pathways and facilitate quantitative analyses to analyse the contribution of different physical drivers in causing an extreme.

### 2.3.5 Causal complex network analysis

Complex systems are often viewed as networks of interacting subprocesses, for example, the human brain [Bullmore and Sporns, 2009], or the Earth system [Donges et al., 2009a, Ludescher et al., 2021]. Tools of network theory [Boccaletti et al., 2006] have been used to analyse quantities such as the information flow as it propagates through the system or the stability of subprocesses [Gozolchiani et al., 2011]. A common network measure is the node degree, which quantifies the number of processes linked to a node. A more involved measure is betweenness centrality, which quantifies the number of shortest paths through a particular node. A crucial question is then to define what these paths mean. In works where the networks are based on pairwise correlation or mutual information [Donges et al., 2009a, Ludescher et al., 2021], one may associate paths with a transfer of information.

However, there is a difference between information being transferred versus perturbations propagating through the network. Here a question can be to identify how critical individual subprocesses are in spreading and mediating perturbations in such dynamic complex systems. The propagation of perturbations, aka interventions, relates to a causal question requiring a causal definition of network links able to distinguish direct from indirect interactions.

In addition, the toolbox of classical network measures is not rich enough for quantifying gateways and mediators of perturbations. Essentially, these measures—with many originating from

the social sciences [Freeman, 1977]—are based on a different definition of links, for example, two persons knowing each other, as opposed to dynamical interactions in a complex system. Hence, here measures based on causal pathways on which perturbations propagate in a complex system's interaction network can be utilised, such as those studied in Runge et al. [2015b], Runge [2015]: Identifying the nodes of the causal graph with the components of the complex system, the average causal effect can be defined as a causal version of the out-degree or closeness centrality, which quantifies by how much an individual component causes any of the remaining components. This serves as a quantitative measure of how much a component is a gateway of perturbations. On the other hand, the average causal susceptibility measures how much a component is changed on average by a perturbation in any of the remaining components as a causal version of the in-degree or in-closeness. Finally, the average mediating causal effect measures how much of the pairwise causal effects between any pair are mediated through a particular variable, which can be seen as a causal version of betweenness centrality. In Runge [2015], these measures are generalised in an information-theoretic framework.

### 2.3.6  Causally robust forecasting models

Forecasting a time series from multivariate predictors constitutes another problem where causal knowledge helps. Even considering the case of forecasting inside the same distribution, that is, assuming a stationary distribution, it can be proven information-theoretically that causal predictors maximise the mutual information with the target variable and, by the Markov property, any further predictors do not add further information. More formally [Brown et al., 2012], the negative log-likelihood can be decomposed as follows

$$\lim_{n\to\infty} -l = \lim_{n\to\infty} -\frac{1}{n}\sum_{t=1}^{n}\log \widehat{p}(X_{t+1}\mid \mathscr{P};\theta) \qquad (2.12)$$

$$= \mathbb{E}\left[\log\frac{p(X_{t+1}\mid \mathscr{P})}{\widehat{p}(X_{t+1}\mid \mathscr{P};\theta)}\right] \quad + \quad \underbrace{I\left(X_{t+1};\mathbf{X}_{t+1}^{-}\setminus\mathscr{P}\mid \mathscr{P}\right)}_{=0} \quad + \quad H\left(X_{t+1}\mid \mathbf{X}_{t+1}^{-}\right), \qquad (2.13)$$

where $n$ is the sample size, $\widehat{p}(X_{t+1}\mid \mathscr{P};\theta)$ the prediction model for $X_{t+1}\in\mathbf{X}_{t+1}$ of the true underlying $p(X_{t+1}\mid \mathscr{P})$ given its causal parents $\mathscr{P}$ and model parameters $\theta$. As shown, the log-likelihood decomposes into the model approximation error given $\mathscr{P}$ (first term), the conditional mutual information between the target and unselected variables $\mathbf{X}_{t+1}^{-}\setminus X_{t+1}$ given $\mathscr{P}$ (second term, zero by the Markov condition), and the irreducible entropy or uncertainty (last term). This is especially relevant for finding optimal sets of predictors in the case where greedy selection strategies do not work because the predictors cause the target variable synergistically, for example, $X_{t+1}=Z_t^1\cdot Z_t^2+\eta_{t+1}$. As shown in Runge et al. [2015a], an optimal subset selection can be better performed on the smaller subset of causal predictors. In Kretschmer et al. [2017], Di Capua et al. [2019], causal pre-selection was used in a climate context. Beyond stationary distributions, Huang et al. [2019] address the task of causally-informed forecasting under nonstationary environments through state-space models.

### 2.3.7  Physical simulation model evaluation

Causal graphs and causal effects can be utilised to intercompare the output of physical models and evaluate and validate them against observations at the level of causal dependencies [Eyring et al., 2020, 2019, Nowack et al., 2020, Pérez-Suay and Camps-Valls, 2019]. One approach in this direction is to compare the causal graph obtained from the observational data to those obtained from simulated data. This procedure has been proposed in the climate sciences to compare climate model simulations and observational data through their corresponding causal graphs derived from

PCMCI [Nowack et al., 2020], cf. Section 4. The methodology could be adapted and applied to other physical science problems where one typically has complex datasets and simulations to confirm hypotheses.

### 2.3.8   Counterfactual causal attribution of extreme events

Counterfactual questions are not about the distributions of a target variable due to possible (future) interventions, but about the distribution of a target variable for an alternative past intervention, given that a particular outcome was observed. Formally, just like interventional causal queries are represented by interventional SCMs, counterfactual queries are represented in counterfactual SCMs [Pearl, 2009c, Peters et al., 2017a]. Given an SCM over variables $\mathbf{V}$ and observations $\mathbf{v}$, in the *counterfactual SCM*, the noise distribution is updated such that the $\mathbf{V} = \mathbf{v}$ holds. Then the noise terms may not be independent anymore. Counterfactual queries are then *do*-statements in the counterfactual SCM. One example of a counterfactual distribution query is $p(y'_{x'}|y_x)$, which specifies the probability of observing $Y = y' \neq y$ under the hypothetical past intervention $do(X = x')$ when, in fact, $Y = y$ was observed under the intervention $do(X = x)$. Such queries can be computed in different ways [Pearl, 2009c, Correa et al., 2021] and generally require more assumptions about the underlying structural causal model than causal effect questions or causal discovery. Next to counterfactual distributional queries, Halpern [2016] discusses causation and counterfactuals regarding single events.

An example of a counterfactual question in climate is the causal attribution of extreme events [Hannart et al., 2016]. The above query $p(y'_{x'}|y_x)$ is one specific type of a counterfactual question and is sometimes called "probability of necessity" (PN), which is typically the quantity of interest in lawsuits. Extreme event attribution requires to study of anthropogenic forcings compared to their absence, that is, solely natural forcings or internal variability of the climate system. If the probability of necessity is high enough, then a human-caused extreme event is established.

### 2.3.9   Signal tracking for the discovery of proximal causes

Many phenomena, such as extreme events in complex systems, such as El Niño events in the climate system and extreme volatility in the financial system, are caused by an initial anomaly that triggers a travelling cascade of events [Press, 1967]. This phenomenon is often called the "butterfly effect"[2], characterised by an anomaly in one part of a system having extreme consequences in another space and time. Such cascaded events are challenging to detect, predict, understand and characterise [Zscheischler et al., 2018, Menzly et al., 2004], and has led to the development of the field of the science surrounding the concept of predictability in complex systems [Grunberg and Modigliani, 1954, Bialek et al., 2001, Boffetta et al., 2002]. A potential strategy for uncovering the cause of notable events is causal discovery, for instance, by conducting a simulation which begins from the start of the event in question and tracing the initiating signal back to its source. However, this presents difficulty in determining the initiating trigger. It is challenging to have a dialogue about recognising the drivers of intense impacts because the amount of correlated drivers is usually much greater than the number of causally pertinent drivers, which may only have a substantial effect when combined (synergy) [Runge et al., 2019a]. These kinds of relations are hard to portray with a pairwise network.

---

[2]The term is attributed to Lorenz when he noted that a weather model failed to reproduce the results of runs with the unrounded initial conditions. However, the idea was earlier recognised by Poincaré and further formalised by Wiener. The analogy became popular and originated the quantitative science of characterising *instability* in complex systems undergoing nonlinear dynamics and deterministic chaos.

Table 2.2: Methods and open-source software for causal discovery.

| Method | Software |
|---|---|
| Granger causality (GC) [Granger, 1969], kernel GC [Marinazzo et al., 2008a], explicit KGC [Diego Bueso, 2020] | causal-learn, statsmodels, KGC, XKGC |
| CCM [Sugihara et al., 2012, Ye et al., 2015] | rEDM |
| PC [Spirtes and Glymour, 1991, Spirtes et al., 2000], FCI [Entner and Hoyer, 2010] | Tetrad, causal-learn, pcalg, Tigramite, MXM, bnlearn, dbnlearn, PyWhy |
| PCMCI [Runge et al., 2019b], PCMCI$^+$ [Runge, 2020], LPCMCI [Gerhardus and Runge, 2020] | Tigramite |
| DYNOTEARS [Pamfil et al., 2020] | Causalnex |
| TiMINo [Peters et al., 2013] | R script |
| VARLiNGAM [Hyvärinen et al., 2010] | causal-learn, lingam, original R code |
| ICP [Peters et al., 2016, Heinze-Deml et al., 2018b, Pfister et al., 2019] | seqICP |

### 2.3.10  Causal benchmarks, software and platforms

Method development and comparison require benchmark datasets with known causal ground truth for validation. Ideally, such ground truth comes from expert knowledge of real data or actual experiments that can also be used to falsify causal relationships predicted from observational causal inference methods. Unfortunately, in many fields, such as Earth system sciences, such datasets exist only for expert-labelled causal relations among a few variables (e.g., some bivariate examples [Mooij et al., 2016]). A tractable approach is to generate synthetic data with simple model systems that mimic properties and challenges of data from the system under study but where the underlying ground truth is known. These can then be used to study the performance of causal discovery (and causal inference methods more generally) for different challenges in realistic finite sample situations. From a practitioner's perspective, it is essential to determine which method is best suited for a particular task with particular challenges and for a specific set of assumptions. Synthetic data, adapted to the problem at hand, can be used to choose the suitable method, including method parameters. A list of key methods for causal discovery and the available software and platforms is given in Table 2.2. An example is the SAVAR model [Tibau et al., 2022] that mimics spatio-temporal features of climate data. The website causeme.net [Runge et al., 2019a, 2020] aims to provide an open platform with synthetic models mimicking real data challenges on which causal discovery methods can be compared. Next to method comparison, the platform also calls for submissions of actual and modelled data sets where the causal structure is known with high confidence and was used on the Causality 4 Climate NeurIPS competition [Runge et al., 2020]. That competition sparked the investigation of a particular property of synthetic data and models called *var-sortability*, which led to new insights in causal discovery methods [Reisach et al., 2021].

## 2.4  Perspectives

This section reviewed causal discovery in the physical sciences, describing the main methods, challenges, and opportunities for future research. We laid out the fundamental elements of the causal discovery framework –SCMs, graphs, and associated distributions–, and gave an overview of the methodological concepts of learning qualitative causal graphs [Runge et al., 2019a, 2023]. We deliberated on commonplace difficulties encountered in the field, such as determining and

preprocessing causal variables, addressing non-stationarity, contemporaneous causation and hidden confounding, and selecting parametric models for nonlinear dependencies and non-Gaussian distributions. Section 4 will illustrate causal discovery methods in neurosciences and Earth sciences case studies.

The body of causal inference has traditionally been embedded in several communities, mainly statistics, social sciences, econometrics and health sciences. Irreconcilable positions and long-standing discussions exist [Dowd, 2011]. Pearl argues that it is essential to distinguish between causal and statistical information, as they refer to two separate concepts[3], and suggests that clear distinctions should be established in the notation used, and each should be subject to different means of calculation [Pearl, 2011]. Arguably, nonparametric SCMs (as a natural generalisation of those used by econometricians and social scientists in the 1950-1960s) have developed the field of causality in new mathematical underpinnings: explicate and enumerate causal assumptions, test implications, decide measurements and experiments, recognise and generate equivalent models, recognise instrumental variables, generalise structural equation models and solve the mediation and external validity problems. These tools, methods and solutions help to determine the accuracy and validity of causal claims in the analysis. The machine learning community is approaching the field of causal discovery in innovative ways by leveraging data, assumptions and models collectively. In recent decades, mathematical foundations have been established to address questions of causality in various scientific fields, mainly emerging for statistics and machine learning [Pearl, 2009c, Spirtes et al., 2000, Peters et al., 2017a]. The causal machine learning (CausalML) field has recently introduced [Kaddour et al., 2022] as an umbrella for machine learning methods based on SCMs. It aims to advance the field in several directions: causal supervised learning, causal generative modelling, causal explanations, causal fairness, and causal reinforcement learning. Applications of the new methods are vast and promise advances in computer vision, natural language processing, and graph representation learning. Therefore, the field of causal discovery is growing in methods, approaches and impactful applications. A unified agenda for Causal Inference is built and deployed in the wild.

However, despite the significant advances in the last decades, many unresolved philosophical and methodological issues remain for causal discovery from observational data. All such challenges also create avenues of research. On the one hand, we identified and discussed algorithmic and data challenges and summarised possible ways to address them in §2.2. Indeed, we must develop more effective methods for incorporating (uncertain) expert knowledge, determining the spatio-temporal complexity of the underlying dynamic phenomena, and creating more reliable and statistically efficient algorithms.

On the other hand, perhaps the most critical challenge is the theoretical impossibility of causal discovery from purely observational data [Pearl, 2009a]. There are, however, ways to tackle the challenge. For example, specifying a causal DAG using domain knowledge can help mitigate the potential inaccuracy of their assumptions of sufficiency and faithfulness. Another possibility to learn about a certain equivalence class may consider incorporating domain knowledge into structure learning algorithms by using "allow lists" and "deny lists" to determine which edges should or should not be included in a DAG or creating a Bayesian prior to assigning varying levels of probability to certain causal relationships [Castelo and Siebes, 2000, Ness et al., 2017]. This is very much related to using *inductive biases* (such as Occam's razor) [Janzing, 2007] and causal invariances (such as parameter modularity and independence of mechanism) [Hoyer et al., 2008] to learn structure beyond likelihood-based scores and conditional independence constraints. Finally, if data from natural experiments, such as $do(A = a)$, is available, this intervention information can be incorporated into the algorithm [Cooper and Yoo, 1999] to automate this reasoning process.

---

[3]Statistical information deals with the probability of certain variables being observed. In contrast, causal information deals with hypothetical relationships in new situations.

Incorporating the abundant domain knowledge within the causal discovery routine can address the identifiability and faithfulness assumptions (very much in line with the basis or sparsity priors used in equation discovery, cf. §3). By joining forces, both can contribute to resolving pressing scientific issues, ranging from process comprehension to evaluating and upgrading the physics included in physics models.

Many problems in the physical sciences can be framed as causal questions. Yet many researchers in economics and health services, and even many computer scientists in machine learning, have been trained to be reluctant to use the language of causality [Dowd, 2011, Pearl, 2011, Hernán, 2018, Schölkopf et al., 2021]. This is a cognitive barrier to resolve in the future. Besides, the language barrier between the methodological and domain science communities is a significant challenge in the causal discovery endeavour. Bridging this divide by translating domain questions into actionable and precisely stated causal inference tasks seems reasonable. Additionally, hesitance to adopt causal inference can be attributed to the lack of suitable benchmarks to help choose an appropriate method. A benchmarking platform (https://causeme.net) was introduced that covers the causal discovery problem setting. It is necessary to have more of these benchmark platforms and easily accessible databases to facilitate better collaboration between the two communities. To successfully address a causal inference problem, it is essential for the domain scientist and computer scientist to work together; assumptions must be discussed and formalised, data characteristics must be jointly analysed, and conclusions must be assessed from both perspectives.

It is commonly believed that most research questions in science can be interpreted as causal inference problems. Sentences such as "we find that $X$ [increases/ decreases/ lags/ leads/ affects/ drives/ impacts] $Y$" are often found in papers, creating an impression of causality. Nevertheless, scientists should be more explicit and transparent when making assumptions which lead to causal conclusions. This is not only sensible but is also necessary when analysing complex systems like the Earth, the Brain, or the Economy, since the outcomes of such research may have significant economic, environmental, and social implications. The field of Causal Inference provides a comprehensive arsenal of tools grounded in rigorous mathematical principles and a vibrant interdisciplinary milieu to confront the challenge at hand successfully.

# 3. Learning physical laws from data

Distilling mechanistic models of the world is how physicists have successfully understood and explained the natural world. The prototypical process starts with an experiment or observation. A mathematical model is hypothesised, which can predict a new experiment's outcome. The observations will either support or falsify the hypothesis, leading to more experiments and refined models.

> *"For centuries, scientists observed Nature to extract simple laws and equations to explain the world mechanisms, to anticipate and predict behaviours and gain faith in their interventions/actions."*

For many complex systems, we have poor models because certain interaction terms are unknown. Another case is when we know some microscopic interaction laws. Still, the emergent properties at a larger scale do not directly follow, such that predictions at the larger scale need new coarse-scale interactions. To cope with these situations and make sense of the sheer amount of data produced by modern instruments, researchers have been looking into automating the processes involved in model building and creating new insights. Many motivations and inspirations have been adopted to guide the scientific method development (see Fig. 3.1): e.g. a physical law or equation describing the system should be compositional, thus following the "divide and conquer" rule, should be as simple as possible (but not simpler) thus following Occam's razor, further developed and formalised by Solomonoff, eventually focus on the interesting parts of the system and disregard the rest, create understanding, generalise and unify different working theories and models, and the learned representation (and its parameters) should be self-explainable, amenable and intuitive.

This section reviews the state-of-the-art in equation discovery from data. Unlike in traditional law discovery *à la Kepler* where trial-and-error was dominant, modern statistics and machine learning techniques exploit the regularities found in the data to discover plausible, simple and explainable equations, to learn feature representations that describe (typically dynamic) systems. We will first consider the *explicit* discovery of equations that describe observed data, also called *Symbolic Regression*. Second, we look into *implicit* discovery through dimensionality reduction techniques and transfer operators. We finish the section by discussing the main challenges and research opportunities.

- "Divide-and-conquer" - *Julius Caesar*
- "Simplicity rules in Nature" *Occam + Solomonoff*
- "Look at the interesting parts, forget the rest" *Galileo*
- "Nothing in life is to be feared, it is only to be understood." - *Curie*
- "Unify theories by parameterization" - *Einstein*
- "Science gives a partial explanation for life. Insofar as it goes, it is based on fact, experience and experiment." - *Franklin*
- "If you want to master something, teach it!" - *Feynman*



Figure 3.1: Historical Inspirations/Motivations in Law/Equation Discovery.

## 3.1 Explicit equation discovery with symbolic regression

Symbolic Regression (SR) refers to a class of machine learning techniques that aim to discover mathematical relationships and patterns in data. SR aims to find a compact, human-readable mathematical expression that accurately reflects the underlying relationships in the data. This approach is particularly useful in cases where the relationships between variables are complex or unknown, and traditional statistical methods may not provide adequate explanations. SR operates under the assumption that the underlying data-generating mechanism can be described by a sparse and algebraic input-output relationship.

There is a major divide in the method to achieve that. One class of methods use genetic algorithms or other discrete search methods to find mathematical terms, typically represented by a graph of mathematical operations. Another class of methods uses continuous space search methods and solves a relaxation of the discrete problem of finding compact equations. A third and most recent addition to the family of symbolic regression methods uses massive amounts of synthetic data for pretraining a system that can quickly guess a suitable expression at test time.

More formally, we are trying to solve the following optimisation problem:

$$\arg\min_{f} \mathbb{E}_{(x,y)\sim\mathscr{D}}\|f(\boldsymbol{x}) - \boldsymbol{y}\|^2 + \lambda C(f) \tag{3.1}$$

where we attempt to find a low-complexity function (equation) $f$ that best[1] maps the inputs $\boldsymbol{x} \in \mathbb{R}^n$ to their corresponding outputs $\boldsymbol{y} \in \mathbb{R}^m$ in the data distribution $\mathscr{D}$, $C(f)$ refers to a measure of complexity of $f$, and $\lambda$ is the weighting factor. For instance, the complexity measure could be the number of terms in the equation.

A central problem when performing symbolic regression is selecting an appropriate weighting factor. More generically, the question is which level of complexity is right. There is probably no definite answer to this question. Instead, we consider the solution to a symbolic regression problem as a family of Pareto-optimal solutions:

$$f^{(c)} = \arg\min_{f} \mathbb{E}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathscr{D}}\|f(\mathbf{x}) - \mathbf{y}\|^2 \qquad \text{s.t.} \quad C(f) = c \tag{3.2}$$

where $f^{(c)}$ refers to the best fitting expression with complexity $c$. A more complex expression will be able to fit the given data at least as well as a less complex expression. Fig. 3.2 illustrates

---

[1]We use the squared error here for simplicity, but other notions of distance are possible.

a typical Pareto curve along the optimisation objectives: goodness of fit (error reduction) and function complexity. In addition to the training error, which monotonically decreases with function complexity, the illustration also shows a hypothetical test error that shows the overfitting of too complex functions.



Figure 3.2: Illustration of a Pareto curve of solutions to an SR problem.

We will now look more closely into different methods that have been proposed to solve the optimisation problem (Eq. 3.2) in practice.

### 3.1.1 Symbolic regression using discrete search methods

The problem of symbolic regression is, at its core, a search for suitable functions $f$ in (3.2). Since those functions should have low complexity, it is natural to attempt to perform a search for functions. The first attempt to do that was proposed by Cramer [1985] by inventing Genetic Programming, which got popularised and applied through Koza [Koza, 1990, 1994]. The idea is simple: search for computer programs to solve a particular problem by iteratively creating many random programs and selecting the best fit, and create a new pool of candidates by recombination and random modification. This mimics the biological evolution process of nature to create the genetic material of living organisms. Applied to symbolic regression: The functions are represented as a graph of input variables, operators and basic algebraic functions.

In the paper by Schmidt and Lipson [2009b], this approach was refined and applied to the discovery of physical laws. As the method was indeed able to discover Lagrangian and Hamiltonian formulations from data, it stimulated a growing interest in symbolic regression and sparked the development of many methods. These general search methods are also referred to as *evolutionary algorithms*. The approach of Schmidt and Lipson [2009b] is illustrated in Fig. 3.3. The general method for symbolic regression was implemented in a tool called *Eureqa* [Praksova, 2011] that is now only available as an online service [DataRobot Inc, 2023].

There are several publicly available and open-source implementations, such as the PySR [Cranmer, 2020], gplearn [Stephens, 2022], Glyph [Quade et al., 2019] and Operon [Burlacu et al., 2020]. A more detailed overview of genetic algorithm-based methods and their combination with gradient descent can be found in Kommenda et al. [2020].

**Feynman AI.**

An approach that exploits physical knowledge such as units and makes reasonable assumptions for equations in physics is Feynman AI [Udrescu and Tegmark, 2020, Udrescu et al., 2020]. The method augments the genetic algorithm searching for expressions by enforcing fitting physical units, decomposing the problem using symmetries and checking separability. To check for symmetries, a neural network is trained on the data to allow accessing whether the underlying function is symmetric. It is worth noting that a large amount of data is used here. From a set of 100 equations taken from the Feynman Lectures, the method was able to recover all of them whereas Eureqa only solves 71.

| **1** Observational time series $\{x(t), y(t), z(t)\}$ | **2** Partial derivatives $\left\{\frac{\Delta x}{\Delta y}, \frac{\Delta z}{\Delta x}, \frac{\Delta y}{\Delta z}\right\}$ | **3** Candidate symbolic functions | **4** Symbolic partial derivatives | **5** Compare predictive par.der. | **6** Accurate and sparsest solutions |
|---|---|---|---|---|---|
|  |  | $f = x\cos(y)$ <br> $f = 0.7\exp(y/z)$ <br> $f = y^2 - 9.8\cos(x)$ <br> ... | $\frac{\partial f}{\partial y} = y + \sin(x)\frac{\Delta x}{\Delta y}$ <br> $\left.\frac{\partial y}{\partial x}\right|_{f(x,y)} = \frac{\partial f}{\partial x} / \frac{\partial f}{\partial y}$ | $\left.\frac{\Delta y}{\Delta y}\right|_{D_i} = \left.\frac{\partial y}{\partial x}\right|_{f_i}$ | $f = z + 9.8\sin(x)$ <br> $f = y^2 - 9.8\cos(x)$ |

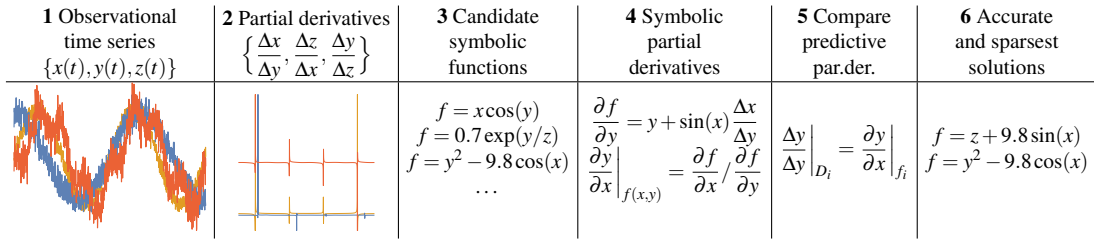Figure 3.3: Schematic view of the symbolic regression method for discovering physical laws in (Schmidt and Lipson, 2009b). Starting from observational data (1), partial derivatives are computed numerically for all pairs of variables (2). A set of candidate symbolic functions $f$ is derived (3), whose symbolic partial derivatives are computed (4) and compared to the predictive ones (5). The process 3-5 is iterated until, finally, a small set of the most accurate and simple equations is returned (6).

### Search with Deep Reinforcement Learning.

A method that uses Deep Reinforcement Learning to search for a suitable solution to the symbolic regression problem is Deep Symbolic Regression (DSR) [Petersen et al., 2021]. The key idea is to treat the search for expressions as an exploration problem in reinforcement learning (RL). The functional expressions are represented as a sequence of tokens corresponding to a depth-first graph traversal and are generated by a recurrent neural network. Numerical constants are fitted using the BFGS optimiser. This generative network is trained on the given dataset using RL to find a highly fitting solution. An interesting contribution is a formulation of a risk-seeking policy gradient that tries to optimise for the best-case scenario (a good solution can be found) rather than the typical average case. The method was able to solve 83% of the standard Nguyen-1 dataset.

### 3.1.2 Sparse linear regression and neural network approach

The symbolic regression problem can also be tackled via traditional regression methods. In contrast to the search in a discrete set of functions, the search is performed in a dense set, typically represented by a real-valued parameterised function. So (3.1) is solved by choosing a large enough function class described by $f_w$ with $\boldsymbol{w} \in \mathbb{R}^p$. The optimisation is then performed over the space of parameter-values $w$. Linear regression is a special case, where the function $f_w(\boldsymbol{x}) = \boldsymbol{w} \cdot \boldsymbol{x}^\top$.

What about the complexity regularisation term $C(f_w)$ in (3.1)? Ideally, the term should count the number of non-zero parameters in $\boldsymbol{w}$, expressed as $|\boldsymbol{w}|_0$ and referred to as $L_0$ norm. However, this term jeopardises the efficient solution of the regression problem because it is non-linear and non-differentiable. One practical alternative is to use the $L_1$ norm instead, i.e. the sum of absolute values, which also leads to sparse solutions (see Fig. 3.4). In the case of linear regression, this is termed LASSO regression [Tibshirani, 1996].

The methods differ in the function class $f_w$, the regularisation term and the optimisation method used.

### SINDy: Sparse identification of dynamical systems.

In some cases, the class of building blocks that might occur as summands in the analytical description of the data are known. Then a rather simple but effective method can be employed that is called *sparse identification of dynamical systems*, SINDy for short. It was proposed in Brunton et al. [2016b] to find differential equations of dynamical systems from observations. For the general symbolic regression problem, the FFX method by McConaghy [2011] was already earlier proposing the same idea. The input data is passed through a predefined library of base functions and interaction terms. Then the resulting high-dimensional representation is fit to the data using sparse linear regression. All relevant terms keep a non-zero weight and constitute the final expression.

Let us unpack this in more detail for a dynamical system of $n$ variables described by the system of ordinary differential equations $\frac{d}{dt}\boldsymbol{x} = \boldsymbol{g}(\boldsymbol{x})$, where $\boldsymbol{x}(t) \in \mathbb{R}^n$. Each component of $\boldsymbol{g}$ can now be

Figure 3.4: $L_1$ regularisation typically leads to sparse solutions. The lines show the isolines of quadratic loss (red) and $|w|$ (blue). Instead of the data-only solution (red star), a sparse solution (green) is found.

substituted by a linear combination of library functions:

$$
\begin{aligned}
\frac{d}{dt}x_1 &= g_1(x_1, x_2, \ldots, x_n) &= w_{11}l_1(x_1, \ldots, x_n) + \ldots + w_{1m}l_m(x_1, \ldots, x_n) \\
\frac{d}{dt}x_2 &= g_2(x_1, x_2, \ldots, x_n) &= w_{21}l_1(x_1, \ldots, x_n) + \ldots + w_{2m}l_m(x_1, \ldots, x_n) \\
&\vdots \\
\frac{d}{dt}x_n &= g_n(x_1, x_2, \ldots, x_n) &= w_{n1}l_1(x_1, \ldots, x_n) + \ldots + w_{nm}l_m(x_1, \ldots, x_n),
\end{aligned}
$$

where $l_1, \ldots, l_m$ is the *predefined finite library of candidate functions* and $w_{ij}$ are the scalar coefficients to be learned following our objective (3.1). To obtain a sparse solution, the $L_1$ regularisation (LASSO) can be used, i.e. $C(f) = |w|_1$, as described above. Alternatively, the (squared) $L_2$ norm of the weights $\|w\|_2^2$ can be used, corresponding to classical ridge regression that permits a closed-form solution. However, an iterative pruning of small weights must be used to obtain a sparse solution.

**Illustration of sparse identification of dynamical systems (SINDy)** [Brunton et al., 2016b] using a synthetic dataset of the well-known Lotka-Volterra system, as shown in Adsuara et al. [2020]. The Lotka-Volterra system models the interaction between prey and its predator in ecology and is given by the following equations:

$$
\begin{aligned}
\frac{d}{dt}x_1 &= \alpha x_1 - \beta x_1 x_2 \\
\frac{d}{dt}x_2 &= -\gamma x_2 + \delta x_1 x_2
\end{aligned}
$$

with the coefficients, $\alpha$ and $\gamma$, being the intrinsic growth/decrease rates of $x_1$ and $x_2$, and $\beta$ and $\delta$ are cross terms taking into account the interaction between species. In our particular case, we will set $\alpha = 3/2$, $\beta = 1$, $\gamma = 3$, and $\delta = 1/2$. We show the results of the identification of these parameters using SINDy in the table below for two levels of additive white Gaussian noise of the signal-to-noise ratio of 5 (high noise level) and 40 dB (low noise level). As usual, the created data was split into train/test data $(75\%, 25\%)$, respectively. The ODE coefficients are recovered sufficiently well to achieve a high correlation coefficient $R$ but are generally underestimated due to the sparsity regularisation.

| Library functions | Learned Coefficients | | | | True Coefficients | |
|---|---|---|---|---|---|---|
| | 40 dB | | 5 dB | | | |
| | $\frac{d}{dt}x_1$ | $\frac{d}{dt}x_2$ | $\frac{d}{dt}x_1$ | $\frac{d}{dt}x_2$ | $\frac{d}{dt}x_1$ | $\frac{d}{dt}x_2$ |
| $x_1$ | 1.3822 | 0 | 1.1404 | 0 | 1.5 | 0 |
| $x_2$ | 0 | -2.9123 | 0 | -2.7946 | 0 | -3 |
| $x_1 x_2$ | -0.9797 | 0.4849 | -0.9520 | 0.4710 | -1 | 0.5 |
| $x_2^3$ | 0 | 0 | 0 | -0.0001 | 0 | 0 |
| R | 0.9999 | | 0.8674 | | | |

For fitting dynamical systems, the temporal derivatives need to be computed. Finite differences are often too sensitive to noise such that kernel regression (aka Gaussian processes), which allow for explicit derivative computation, are preferred [Camps-Valls et al., 2016, Johnson et al., 2018]. A more recent approach is to solve noise estimation and model identification in one joint optimisation procedure [Kaheman et al., 2022]. Intuitively, for every data point, the corruption by noise is estimated. Since this optimisation problem is highly underdetermined, an additional constraint is used, namely that when integrating the estimated dynamical system model a small error should occur. This trick separates noise from the signal and leads to an improved estimation quality.

### Neural network approach: Equation Learner.

Enlarging the function class $f$ is possible using neural networks. Probably the first work in this direction is the Equation Learner (EQL) introduced in Martius and Lampert [2016] that uses a neural network with algebraic base functions and a particular regularisation scheme to solve (3.1) and 3.2. The function $f$ is represented by a neural network, modified only to contain elementary operations that should appear in a potential solution. Figure 3.5(left) shows the architecture of the Equation Learner (EQL) in a simplified form.



Figure 3.5: Equation Learning Architecture and example equation. Left: Illustration of the EQL Architecture, reproduced from (Sahoo et al., 2018), a feed-forward neural network with special *activation* functions (sin, multiplication etc.). Note that each unit type will occur many times. Right: Example system with four inputs $x_{1,2,3,4}$ and one output $y$. Training is only in the $[-1,1]^4$. EQL recovers the equation and extrapolates (Sahoo et al., 2018).

The input variables are mapped with a dense layer to multiple instances of trigonometric functions, identity, multiplication and division, but more base functions, such as squares or exponentials, are possible. The resulting values are again mapped with a matrix to another layer of elementary functions, and so forth, until the last layer corresponds to the output (containing only division operators in the picture).

The network is trained using stochastic gradient descent (e.g. Adam) on the mean squared error loss and $L_1$ regularisation on the weights to induce sparsity:

$$\mathbb{E}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathscr{D}}\|f_W(\boldsymbol{x})-\boldsymbol{y}\|^2+\lambda|\boldsymbol{W}| \tag{3.3}$$

where $f_W$ denotes the neural network with parameters $\boldsymbol{W}$.

Note that the system needs to be differentiable for training, such that a pure complexity term, such as $L_0$ regularisation that would count the number of non-zero weights, does not work out of the

Figure 3.6: Illustration of uncertainty estimates using a mixture of Laplace approximations of learned equations. Top row: toy example $y = 0.8 \cos x - 0.4 + \varepsilon$ where $\varepsilon \sim \mathcal{N}(0, 0.03^2)$ with just 6 datapoints. Bottom row: Atmospheric $CO_2$ concentration at Mauna Loa Observatory (Observatory, 2020) (concentration vs. time, both in arbitrary units). The left panels show the predictive distributions. The panel in the middle shows individual local Laplace approximations with $2\sigma$ (shaded area) for the toy data and a zoomed density for the Mauna dataset. The colour represents the weight and aligns with the Pareto plots on the right side, showing RMSE over the complexity of each equation. Reproduced from (Werner et al., 2022).

box. Although methods have been developed since then [Louizos et al., 2018], the $L_1$ regularisation in (3.3) is effective but creates an undesired trade-off between error and sparsity. Something that we also encountered in the illustrating example when applying SINDy. The EQL method introduces an additional regularisation phase after converging with $L_1$ that clamps all $|W_i| < \varepsilon \ll 1$ to zero and optimises without regularisation. This yields a practical approach to optimising for sparsity without trade-offs. In Kim et al. [2021] an alternative to $L_1$ with $L_{0.5}$ was used.

The reader may wonder how the system is successfully trained with elementary functions such as division or square root. Indeed, a naive application would fail due to exploding values or gradients. In Sahoo et al. [2018] and Werner et al. [2021], suitable parameterisations and training steps are proposed. Choosing different $\lambda$ (Eq. 3.3) will create differently sparse resulting networks. Each represents a particular symbolic expression resembling the Pareto curve illustrated in Fig. 3.2. In Fig. 3.5 (right), a synthetic example system is shown. The training data is only generated in the $[-1, 1]^4$ hypercube. The correct equation was discovered, and perfect extrapolation is possible in this case, see Sahoo et al. [2018] for details.

Instead of manually selecting a particular solution, which might be a good procedure when structural insights are to be obtained when investigating some unknown phenomenon, one can also use several or all solutions along the Pareto curve to estimate uncertainty about the predictions for extrapolation, as proposed in Werner et al. [2022]. Figure 3.6 illustrates this approach. For each found equation, a Laplace approximation allows to approximate the uncertainty due to parameter estimation errors and yields a Gaussian posterior. Combining these using a weight based on the validation error and the complexity yields the estimated density. Note how the uncertainty in extrapolation shows clearly the structure of the discrete set of automatically generated hypotheses.

### 3.1.3 Learning to solve symbolic regression

All methods so far treat every symbolic regression problem in isolation – the search or optimisation algorithm was applied to a new dataset from scratch. We are now looking into the idea of learning to solve a particular problem quickly by using data from a whole class of symbolic regression

instances. Generally, the idea is to approximate the inverse mapping from data to a suitable equation. The Dreamcoder paper by Ellis et al. [2021] showed the first instantiation of this idea. Provided with the language of algebraic expressions (arithmetic operations, variables, base functions) and a *simulator* to generate data for a particular equation instance, the method learns a probabilistic mapping from data to equation terms and a library of common equation building blocks. Given a particular instance of data, the system can relatively quickly guess and verify suitable explaining equations. Developing the idea of pretraining further and specialising it for symbolic regression was done by the following method.

**NeSymReS.**

The approach in Biggio et al. [2021] is to use a high-capacity transformer model pretrained to solve the symbolic regression problem. The method is called *Neural Symbolic Regression that Scales* (NeSymReS). The method uses a large set of symbolic regression problems to approximate the inverse mapping from data to equations. After this pretraining phase, given new data, the inverse mapping can generate likely candidate equations. Intuitively, an experienced data scientist might solve the problem similarly: looking at the data and postulating a particular functional form that might explain it, testing it, and potentially trying a different plausible hypothesis.



Figure 3.7: Overview of the NeSymReS method. Left: step with randomly generated training data. Right: inference of candidate equations for unseen data. Reproduced from (Biggio et al., 2021).

Let us look closer at the method. As visualised in Fig. 3.7, the core is a transformer[2] architecture [Vaswani et al., 2017] that can generate algebraic expressions symbol-by-symbol given a set of data points.

The input is represented as a set of $(\boldsymbol{x}, \boldsymbol{y})$ pairs (1024 in Biggio et al. [2021]), which are processed through a set-encoder. The latter is invariant to permutations of the data points. The output of the transformer are tokens that correspond to the typical symbols of input variables, base functions, operators and constant placeholders, which resemble the skeleton of the predicted function. Importantly, the transformer does not have to guess the right constants, just their location in the expression, as these constant placeholders are fit to the data using non-linear optimisation (here BFGS).

Trained on millions of synthetically generated pairs of random expressions with corresponding data, the transformer does a remarkable job guessing likely equations. Importantly, the prediction is not deterministic but allows sampling of possible functions. Thus, new potential solutions can be generated and validated when new data is presented at test time until a sufficiently good fit is found or the Pareto.

Figure 3.8 shows the accuracy of different SR methods for unseen equations from the Feynman and Nguyen benchmarks. The performance is presented in dependence on wall-clock time. NeSymReS is remarkably fast at finding a well-fitting expression for the data.

As a downside, the method was only shown for three input variables, and it remains to be seen

---

[2]The transformer architecture is the basic building block of many large-scale machine learning systems, such as GPT-3 [Brown et al., 2020].

Figure 3.8: Performance of NeSymReS, DSR, classical SR (using gplearn (Stephens, 2022)), and Gaussian processes on the AI Feynman and Nguyen datasets (equations unseen during training). Reproduced from (Biggio et al., 2021).

how much it can be scaled in this respect. Also, the dataset used to guess an equation at test time cannot be big (currently in the order of 1000 data points) because the set transformer encoder cannot yet handle larger sets well.

### 3.1.4 Comparison

As there are quite a number of methods, we aim to discuss their differences, strengths and weaknesses by comparing them along a set of axes. We start with using domain knowledge, as we seldom face a completely uninformed setting in physics. We continue with aspects of the embedding, scaling, speed and usability, summarised in Table 3.1.

**Using domain knowledge.**
A common form of domain knowledge is the base functions and their approximate frequency of occurrence in describing the system under consideration. In standard symbolic regression with genetic algorithms, the number and kind of base functions are very flexible, and each term can have its individual penalty in terms of complexity. More specific domain knowledge, such as monotonicity, function image constraints and derivative constraints, can also be considered, as presented in Kronberger et al. [2022]. Another recent work is presented in Cornelio et al. [2023] that allows to incorporate axiomatic contraints.

In FFX/SINDy, the library of functions is the prime way to specify domain knowledge. Relative preferences could be implemented by different regularisation strengths.

For the EQL framework, the choice of base functions is a bit more complicated, as the systems need to remain optimisable with gradient descent. In Werner et al. [2021], a suitable relaxation for functions with divergences (in function value or derivative) is proposed, and a way to specify preferences among base functions is analysed. The control of the relative frequency of used terms is possible but less direct than in genetic algorithm-based methods. In NeSymReS, domain knowledge can be embedded by selecting/generating the training set with appropriate synthetic problems, although this was not explicitly demonstrated.

**Scaling.**
Most SR methods are for small-scale problems with a few hundred to a few thousand data points and low-dimensional problems, i.e. 1-10 input and output variables. Classical search methods scale unfavourably with dimensionality as the search space grows exponentially. That is why most SR methods are good at finding relatively small and compact equations for low-dimensional systems

Table 3.1: Comparison of symbolic regression methods. See Sec. 3.1.4 for more details.

| | embeddable | scaling | speed | restriction | domain knowledge |
|---|---|---|---|---|---|
| Genetic Programming [Schmidt and Lipson, 2009b] | ✗ | ✗ | slow | for small systems | base-functions, complexity of terms |
| AI Feynman [Udrescu and Tegmark, 2020] | ✗ | ✗ | slow | for physical systems in canonical form | physic: units, symmetries |
| DSR [Petersen et al., 2021] | ✗ | ✗ | medium | small input dim | training domain |
| FFX [McConaghy, 2011], SINDy [Brunton et al., 2016b] | ✓ | ✓[3] | blazing | needs known library | training domain |
| EQL [Sahoo et al., 2018] | ✓ | ✓ | slow | base functions limited, sometimes less concise | base-functions, complexities |
| NeSymReS [Biggio et al., 2021] | ✗ | ✗ | fast | small input dim | training set |

but fail for both high-dimensional systems or those where larger equations are the most compact solution.

FFX/SINDy can handle large output dimensions easily, using some form of ridge regression. However, it also suffers from large input dimensions as the library bank becomes exponentially large unless a factorisation or other simplifying structure is known, for instance, a strong locality assumption in PDEs.

NeSymReS is also limited to several variables and small data sets. Although, the limitation of the data set size can be lifted by sampling a smaller subset of data points for guessing the skeletons and using all data for the parameter tuning with BFGS.

EQL is the only scalable method, as gradient descent works on all dimensions simultaneously, and machine learning methods are developed to scale. A larger initial neural network should be used for larger systems, but no specific adaptations are required.

**Differentiability and Embeddability.**
An interesting feature is whether the SR system can be used as a module in a larger computational pipeline. For instance, if observations are available as images and the causal variables have to be first extracted from the images before they can be used for a concise SR prediction module. An example using SINDy inside an autoencoder architecture [Champion et al., 2019] learns the coordinate frame and dynamics equations at the same time, as detailed in Sec. 3.2.4 below. EQL is conceptually easy to embed into larger architectures, such as deep networks, as it is end-to-end differentiable. An example is discovering PDEs [Long et al., 2019], or learning an energy function given observations of the dynamics in the context of density functional theory [Lin et al., 2020], discussed in more detail in Sec. 4.5. The other methods are more difficult to embed.

**Speed.**
As shown in Fig. 3.8, symbolic regression methods are best compared in performance per compute-time, because search methods can, in principle, find the global optimum given enough time

---

[3]It scales well to large output sizes; for high-dimensional input strong structural assumptions are required.

(although this time might be longer than the age of the universe), so just the "final" performance is difficult to measure. SINDy is not in this comparison because it requires knowledge of the base functions. However, it would be the fastest method, followed by NeSymReS, DSR and classic GPs. EQL is likely the slowest method on small systems because it requires a long training time. However, the time does not significantly increase with system size and amount of data.

**When to use which method?**
For systems where we have a good idea about the occurring modules of the functional form, FFX and SINDy are probably the method of choice. They are simple and effective. For dynamical systems SINDy is the most specialized method. The other SR methods are good if the functional building blocks are unknown or nested structures are expected. Genetic Programming based methods generally shine on small problem settings. Modern implementations can also fit constants, but they result is sometimes complex nested structures. DSR and NeSymReS are faster than standard SR methods and can yield potentially less complex equations. However, less software is built for them, and they are less easy to use. For high-dimensional data or when high fitting accuracy is required, EQL might be the right choice. Also, when SR should be embedded, only FFX, SINDy and EQL are practically usable.

## 3.2 Implicit equation discovery: dimensionality reduction and transfer operators

This section reviews state-of-the-art methods for recovering implicit feature representations of systems from data. We will review connections among methods and emphasise the role and examples in the broad discipline of physics. Unlike in the previous section, an explicit equation is not discovered but rather an operator that encapsulates the system's characteristics (typically spatio-temporal dynamics). The field is tightly related to dimensionality reduction and feature extraction in machine learning and signal processing [Arenas-García et al., 2013], but also to transfer operators in functional analysis [Khodkar et al., 2019].

### 3.2.1 Reduced-order models

In many domains, the goal is to study system dynamics from model simulations. In this case, equations are encapsulated in the model itself, but large-scale, high-fidelity nonlinear models can be challenging to simulate and require significant computational power. In such cases, reduced order models (ROMs) can simplify analysis and control design by trading off model accuracy for computational complexity reduction. ROM can combine complex component-level simulation models into system-level simulations used for control analysis and design.

Two main classes of techniques for building ROMs: model-based and data-driven. Model-based methods rely on a mathematical or physical understanding of the underlying model and are designed for specific PDE-based models. In contrast, data-driven methods use input-output data from the original high-fidelity first-principles model to construct a ROM that accurately represents the underlying system. Data-driven ROMs can be either static or dynamic models. Static ROMs can be developed using techniques such as curve fitting and lookup tables (LUT), while dynamic ROMs can be developed using deep learning techniques such as LSTM, feedforward neural nets, and neural ODEs. The obtained ROM ideally contains the essential physical mechanisms of the original system while exhibiting simpler dynamics that can enhance interpretability. In addition to the simpler physics, the ROM would be much cheaper from the computational point of view than the numerical integration of the governing equations, which can help obtain more efficient control and optimisation techniques.

Developing a ROM typically requires two steps. The first step is to find a set of coordinates where the original dynamics can be expressed in a compact form. This is typically done in terms

of *modes*, associated with coherent features in the original system [Taira et al., 2017, Rowley and Dawson, 2017]. In the second step, it is necessary to find a set of differential equations governing the temporal evolution of the amplitudes of the aforementioned modes. These equations, which constitute a dynamical system, enable shedding light on the physics of the (reduced) phenomenon under study. A widely-used approach to perform the modal decomposition is the so-called proper-orthogonal decomposition (POD) [Lumley, 1967], which is also known as principal component analysis (PCA) in statistics and empirical orthogonal function (EOF) analysis in meteorology [Arenas-García et al., 2013], and is closely connected with the singular-value decomposition (SVD) [Brunton and Kutz, 2022].

> **Proper-orthogonal decomposition (POD), aka PCA or EOF** [Lumley, 1967, Arenas-García et al., 2013] decomposes a dataset or a high-dimensional (spatio-temporal) field into a set of orthogonal basis functions called modes or eigenfunctions, which capture the dominant features of the data. The first few modes explain most of the variability in the data, while the later modes explain smaller and smaller amounts of variability. By truncating the number of modes, one can obtain a low-dimensional representation of the data that preserves the essential features of the original system.
>
> Given a set of data snapshots $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$, where $m$ is the dimension of the data and $n$ is the number of snapshots, we seek to decompose the data into a set of $r$ orthogonal modes $\{\mathbf{u}_i\}_{i=1}^r$, such that $\mathbf{X} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, where $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \ldots, \mathbf{u}_r] \in \mathbb{R}^{m \times r}$, $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$ is a diagonal matrix containing the singular values, and $\mathbf{V}^\top \in \mathbb{R}^{r \times n}$ is the matrix of temporal coefficients. The modes $\{\mathbf{u}_i\}_{i=1}^r$ can be computed by performing a singular value decomposition (SVD) of the data matrix $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ where $\sigma_i$ is the $i$-th singular value, and $\mathbf{v}_i$ is the $i$-th right singular vector. The modes are then given by $\mathbf{u}_i = \frac{1}{\sqrt{\sigma_i}}\mathbf{X}\mathbf{v}_i$.

The POD/PCA framework enables decomposing spatio-temporal data (e.g. flow velocities, weather or climate variables which depend on the spatial coordinates and time) into a set of spatial modes (which only depend on the spatial coordinates), multiplied by their temporal coefficients (which define their change of amplitude with time). POD/PCA ensures that components are orthogonal and optimality with respect to the variance explained by a reduced number of modes, cf. Fig. 3.9 for an example of PCA on a toy spatio-temporal data flow. Other alternative multivariate methods, like partial least squares (PLS) or canonical correlation analysis (CCA) seek projections that maximise covariance or correlation, respectively [Arenas-García et al., 2013]. Still, all these projection methods are linear and thus cannot cope with nonlinear spatio-temporal feature relations and complex dynamics. This can be addressed with kernel machines [Arenas-García et al., 2013]. Oblique and nonlinear transformations can also be learned by embedding Varimax in Reproducing Kernel Hilbert Spaces (RKHS) explicitly [Bueso et al., 2020, Diego Bueso, 2020]. Other ways to obtain non-linear transformations from observation space to ROM space, e.g. using neural networks, will be discussed below.

### 3.2.2  Transfer operators for learning nonlinear dynamics

Transfer operators are related to the abovementioned methods and allow the characterisation and modelling of complex dynamic systems. These operators' eigenfunctions can decompose a system given by an ergodic Markov process into fast and slow dynamics and identify modes of the stationary measure called metastable sets.

The Koopman operator is a linear operator that describes the dynamics of a system by lifting the state variables into an infinite-dimensional Hilbert space. Thus, it enables us to effectively linearise complex temporal trajectories and hence is a compelling approach in dynamical systems research [Khodkar et al., 2019]. Its application is expanding in both theoretical [Khodkar et al., 2019, Brunton et al., 2017], and practical domains from molecular dynamics and fluid dynamics, atmospheric sciences, and control theory [Schmid, 2010, Brunton et al., 2016a, Klus et al., 2018].

The advantages of the Koopman operator are numerous. First, it is a powerful tool for analysing

and predicting the behaviour of a system over time. By lifting the state variables into a higher dimensional space, the Koopman operator can identify patterns in a system's behaviour that may otherwise be difficult to detect. This can be especially useful for uncovering hidden dynamical structures in chaotic systems.

Second, the Koopman operator allows us to develop data-driven models of dynamical systems. Using the operator's eigenfunctions as basis functions, it is possible to develop models of dynamical systems operating on these summarised coordinates, without solving or even understanding the underlying equations of motion. This makes the Koopman operator an attractive tool for model-based control and optimisation. Third, the Koopman operator is useful to discern key properties of highly nonlinear dynamical systems. In short, with the expansion of original state variables to infinite dimensions, the operator can uncover subtle nonlinear behaviour in a system that would otherwise be difficult or impossible to detect [Khodkar et al., 2019].

**Koopman operator** [Khodkar et al., 2019] The Koopman operator is a linear operator that describes the evolution of an observable function of a dynamical system. Let $\mathscr{M}$ be a manifold of dimension $n$ and $f : \mathscr{M} \to \mathbb{R}$ be a real-valued function. The Koopman operator, denoted by $\mathscr{K}$, is defined as an infinite-dimensional linear operator that acts on the space of observable functions $f$ such that for any $f \in L^2(\mathscr{M})$,

$$\mathscr{K} f(\mathbf{x}) = f(T(\mathbf{x})),$$

where $T : \mathscr{M} \to \mathscr{M}$ is the evolution operator that maps each point $\mathbf{x} \in \mathscr{M}$ to its next iterate in time. Now, let $f : \mathscr{M} \to \mathbb{R}$ be a bounded, measurable observable and $\mathscr{K}$ be the Koopman operator associated with a dynamical system. Then, there exists a sequence of eigenfunctions $\psi_j : \mathscr{M} \to \mathbb{C}$ and a corresponding sequence of eigenvalues $\lambda_j \in \mathbb{C}$ such that

$$\mathscr{K} \psi_j(\mathbf{x}) = \lambda_j \psi_j(\mathbf{x}).$$

The Koopman operator preserves the linear structure of the space of observables, provides a linear representation of nonlinear dynamics, and its eigenfunctions provide a useful basis for approximating dynamical systems.

There are, however, some drawbacks associated with the Koopman operator. First, it is intrinsically an infinite-dimensional operator, and although there are efficient finite-dimensional approximations available, the accurate computation of their eigenvalues and eigenfunctions can be computationally expensive [Kaiser et al., 2020]. This can be a problem for real-time applications, such as model-based control [Kaiser et al., 2021a]. Second, related to this drawback, the operator's eigenfunctions are often difficult to interpret, hampering the capacity of the learned representations to explain the underlying system's dynamics. Finally, the Koopman operator assumes that a locally-linear behaviour can represent nonlinearities sufficiently accurately. Thus, the Koopman operator is an important theoretical and applied research tool for understanding and predicting the behaviour of complex dynamical systems. Its data-driven approach to model-based control and optimisation has opened up new possibilities for real-time applications [Kaiser et al., 2021a]. Moreover, its ability to uncover subtle nonlinear behaviour in chaotic systems has made it invaluable for studying chaotic dynamical systems [Takeishi et al., 2017].

More specifically, given a space in which the dynamics is linear, a successful approach to approximate transfer operators (such as the Koopman operator) from the data is called dynamic-mode decomposition (DMD) [Schmid, 2010]. It is also possible to obtain ROMs using DMD [Schmid, 2010], which is also based on concepts from linear algebra and assumes that the system's state can be advanced in time via a linear operator $\mathbf{A}$. While the POD modes are orthogonal in space, the DMD ones are orthogonal in time, and each mode is associated with a particular frequency and a growth rate. Therefore, DMD may help to identify temporal patterns in the data more clearly than POD. In contrast, POD may lead to a more compact low-order representation of the original system due to its optimality property. See an illustrative example in Fig. 3.9.

**Dynamic-mode decomposition (DMD)** [Schmid, 2010] is a technique used to approximate the normal modes and eigenvalues of a linear system. Additionally, these modes can be associated with a damped or driven sinusoidal behaviour in time. DMD is useful for identifying a system's frequency and decay/growth rate. Let us define a dynamical process formulated as $\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, t, \mu)$, where $\mathbf{x}$ defines a measurement, $t$ is a time, $\mu$ is a parametric dependence, and $f$ indicates an unspecified system but from which we obtain many data. Therefore, the complex dynamical system $f$ can be approximated as follows $\frac{d\mathbf{x}}{dt} \approx \mathbf{Ax}$, where $x \in \mathbb{R}^n$, $n \gg 1$ and $\mathbf{A}$ defines a linear dynamical system. Then its general solution is the 'exponential solution' defined as $\mathbf{x} = \mathbf{v}e^{\lambda t}$, where $\mathbf{v}$ and $\lambda$ are eigenvectors and eigenvalues of the linear system $\mathbf{A}$. The problem of finding the eigenvectors $\mathbf{v}$ and the eigenvalues $\lambda$ is a eigenvalue problem defined as $\lambda \mathbf{v} = \mathbf{Av}$.

Yet, we are interested in obtaining $\mathbf{A}$, not its eigendecomposition. This is what the so-called 'exact DMD' does. DMD uses observations/measurements $x_j = \mathbf{x}(t_j)$, defined at a time point $j$ to construct two matrices: the first concatenating the data from the first snapshot to $(m-1)$-th snapshot, and the second with the shifted-by-1-time-step samples, $\mathbf{X}$ and $\mathbf{Y}$, respectively. The goal is thus building a linear dynamical system $A$ fitted with $\frac{d\mathbf{x}}{dt} = \mathbf{Ax}$, and thus *learn* the linear dynamical system $\mathbf{A}$ that takes the data $\mathbf{x}$ from current state $(j-1)$ to future state $(j)$, that is $\mathbf{Y} = \mathbf{AX}$. The linear dynamical system $A$ can be extracted using a pseudo-inverse $\mathbf{X}^\dagger$ of $\mathbf{X}$, that is $\mathbf{A} = \mathbf{YX}^\dagger$. Intuitively, the linear dynamical system $\mathbf{A}$ performs a least-square fitting from the current state $\mathbf{X}$ to the future state $\mathbf{Y}$.

Over the last decades, different numerical methods have been introduced: Ulam's method [Ulam, 1960], extended dynamic-mode decomposition (EDMD) [Williams et al., 2015a,b, Klus et al., 2016], and the variational approach of conformation dynamics (VAC) [Noé and Nüske, 2013, Nüske et al., 2014]. The advantage of purely data-driven methods is that they can be applied to simulation and observational data. Hence, information about the underlying system itself is not required. An overview and comparison of such methods can be found in [Klus et al., 2018]. Applications and variants of these methods are also described in [Rowley et al., 2009b, Tu et al., 2014, McGibbon and Pande, 2015], while kernel-based reformulations of the methods above have been proposed before in [Williams et al., 2015b, Schwantes and Pande, 2015]. Note that the framework of higher-order dynamic-mode decomposition (HODMD) [Le Clainche and Vega, 2017] enables relaxing the linear assumption by including several temporal snapshots to build the operator by exploiting Takens' delay-embedding theorem [Takens, 1981]. The HODMD approach requires additional hyper-parameter tuning, but it has led to very insightful results, for instance, in the context of complex turbulent flows, where this method has enabled identifying the coherent structures responsible for the concentration of pollutants in cities [Lazpita et al., 2022]. Another relevant application of HODMD includes cardiovascular flows [Groun et al., 2022].

### 3.2.3 Dynamic modes in neural-network latent spaces

Despite the interesting properties of POD and DMD, their inherent linearity typically leads to the requirement of very large numbers of modes to reconstruct most of the variance of the original signal, for example, in three-dimensional turbulent flows [Baars and Tinney, 2014]. Neural networks, especially autoencoders (AEs), have been proposed to obtain a reduced-order nonlinear representation of the original data. AEs exploit non-linear activation functions to produce significantly more compact representations in the latent space than those with, e.g. POD [Hinton and Salakhutdinov, 2006]. Figure 3.10[a-b] shows the use of AEs for learning (dynamic) feature representations.

AEs have been used in fluid mechanics to obtain compact modal decompositions of the flow around a two-dimensional cylinder [Murata et al., 2020] and in more complex turbulent flows, e.g. the flow in a simplified urban environment [Eivazi et al., 2022]. Interestingly, when restricting neural networks to linear activation functions, one recovers the POD modes, as shown by Milano and Koumoutsakos [2002] with a multilayer perceptron (MLP) in turbulent channel flow. Shallow NNs have been used for flow reconstruction, in this case from sparse measurements, as illustrated by Erichson et al. [2020] for several flow cases. A more general illustration of the potential of

| (a) Space | (b) Time | (c) Data | (d) PCA | (e) DMD |
|:---:|:---:|:---:|:---:|:---:|
| $x_1(s,\cdot)$ | $x_1(\cdot,t)$ | $x(s,t=10)$ | Mode 1 | Mode 1 |
| $x_2(s,\cdot)$ | $x_2(\cdot,t)$ | | Mode 2 | Mode 2 |

Figure 3.9: Comparison between DMD and PCA with synthetic spatio-temporal data. The signal under analysis $x(s,t)$ (c) is the sum of two generative signals (a,b): $x_1(s)$ is a Gaussian that decays exponentially, and $x_2(s,t)$ is a square that oscillates at a lower frequency. The projections onto the two top components of PCA (d) and DMD (e) show that DMD extracts cleaner spatial coherence patterns from the data.

AEs based on convolutional neural networks (CNNs) was presented by Lee and Carlberg [2020], and an application to spectral submanifolds was developed by Cenedese et al. [2022]. The reader is referred to Refs. [Benner et al., 2015, Carlberg et al., 2017] for a survey of classical methods applicable to linear subspaces.

Despite the superior compression performance of AEs compared with POD, the former does not have two very interesting properties of the latter, namely the optimality and orthogonality of the resulting modes. These are important properties due to their connection with interpretable and parsimonious ROMs. Regarding optimality, Fukami et al. [2020] proposed an interesting approach based on hierarchical autoencoders (HAEs). They first trained a CNN-based AE fixing the dimension of the latent space to just one, obtaining one latent vector. Then, they trained another CNN-AE with a latent dimension of two and fixed the first latent variable to the one obtained in the previous NN, thus obtaining a second latent vector. Through this recursive strategy, they obtained a sequence of latent vectors exhibiting progressively less contribution to the reconstruction of the original signal, allowing them to establish a ranking in the resulting modes. This approach was tested in the flow around a two-dimensional cylinder, although it is important to note that the resulting modes were not orthogonal. This was addressed by Eivazi et al. [2022], who used $\beta$-variational autoencoders ($\beta$-VAEs), which enable introducing stochasticity in the latent space to impose orthogonality in the resulting AE modes, a phenomenon that was explained in Rolinek et al. [2019] among a connection of $\beta$-VAEs to PCA. Also in the case of the $\beta$-VAEs, the modes can also be ranked in terms of their contribution to the reconstruction.

### 3.2.4 Equation discovery in latent representations

Perhaps the biggest challenge in data-driven model discovery is balancing model efficiency with descriptive capabilities. Parsimonious models with the fewest terms required to capture essential interactions promote interpretability and generalizability. However, obtaining parsimonious models is linked to the coordinate system in which the dynamics are measured. The previous methods based on dimensionality reduction, e.g. ROM, DMD or AE, extract expressive components without simultaneously discovering coordinates. In [Champion et al., 2019], AEs were trained for data reconstruction and to recover a parsimonious dynamical system model through sparse regression using SINDy (see Sec. 3.1.2). See Fig. 3.10[c]. The joint goal of discovering models and coordinates

Figure 3.10: Approximating Koopman operator with explicit nonlinear mappings using autoencoders. (a) An autoencoder neural network learns a mapping $\Phi$ compressed latent representation $\mathbf{y}$ from input data $\mathbf{x}$ by minimising the reconstruction error $L$. (b) One can incorporate the Koopman linear op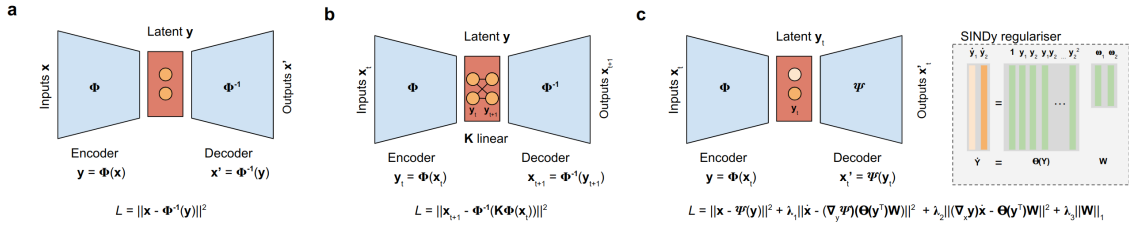erator $\mathbf{K}$ operating in the latent representation $\mathbf{y}_t$, which can then be used for prediction $\mathbf{x}_{t+1}$ from the transformed $\mathbf{y}_{t+1}$ (Brunton et al., 2017). (c) Those representations are not necessarily physically consistent. This can be addressed by enforcing equation dictionaries using SINDy in the loss for the simultaneous discovery of coordinates and parsimonious dynamics (Champion et al., 2019). The loss now accounts for the reconstruction of the input data, as in a regular autoencoder, and the temporal dynamics (gradients $\nabla_x, \nabla_y$) of $\mathbf{x}$ and $\mathbf{y}$ projected onto SINDy bases.

is critical for understanding many modern systems. Using SINDy as an explicit equation discovery regulariser in the latent space balances simplifying coordinate transformations and nonlinear dynamics to identify coordinate transformations where only a few nonlinear terms are present.

### 3.2.5   Discovering fundamental variables

Despite advances in equation discovery (either through implicit or explicit representations), the main core problem is identifying state variables. The discovery typically refers to the identification of the governing equations, not the identification of the physical forces or variables. The vast majority of data-driven models of discovery rely, however, on pre-existing knowledge of the state variables, e.g., the position and velocity of a rigid body object. This relies on deep domain knowledge and strong assumptions. In addition, such assumptions cannot work properly for new physical systems or when those state variables cannot be measured. The work [Chen et al., 2022] proposed a principle for determining the number and identity of state variables in a system from high-dimensional data and demonstrated high effectiveness using video recordings of physical systems. The algorithm discovered the intrinsic dimension of the observed dynamics and could identify candidate sets of state variables without prior knowledge of the underlying physics. Alternatively, other studies sought to identify the fundamental state variables via manifold learning in ambient RKHS (termed *Diffusion* maps, e.g., [Kemeth et al., 2022, Thiem et al., 2020], see also Section 4.1.2).

In short, the field of *variable discovery* is filled with many opportunities in the physical, biological and chemical sciences [Pukrittayakamee et al., 2009, Kemeth et al., 2022, Schütt et al., 2017], as well as many challenges [Chen et al., 2022]. Finally, and interestingly, we want to emphasise that variable discovery is intimately related to revealing latent confounders in the field of causal inference [Stegle et al., 2010, Monti et al., 2020, Diaz et al., 2023].

## 3.3   Perspectives

Let us indulge ourselves with a brief overview of the main challenges (both conceptual and technical) and the opportunities for future research in the field of equation discovery for the physical sciences.

### 3.3.1   Challenges

The field of equation discovery from data is very prolific and is situated at the intersection of many communities: statistics, machine learning, computational fluid dynamics, Bayesian inference, dynamical systems and control theory, functional analysis and causal inference. The field has occupied scientists for centuries at all levels. The quest for optimal and automated solutions has tradition-

ally considered moving in one or several subspaces in the sparsity-extrapolation-generalisation space, i.e., models should be simple, generalisable/robust, and capable of extrapolating outside the sample space (Fig. 3.11). These are very ambitious goals, implying both theoretical and practical challenges.



Figure 3.11: On the quest for the optimal model in the sparsity-extrapolation-generalisation space.

**Theoretical challenges.**
From a more theoretical perspective, both the *identifiability* of the system's equations [Fajardo-Fontiveros et al., 2023, Antonelli et al., 2022a] and the role (or preference) for *sparsity* (simplicity) have been questioned [Fuentes et al., 2021, Guimerà et al., 2020]. In addition, there is a long-standing debate on the *evaluation* of the obtained solution, where many criteria can be adopted. Is it only about invariances and robustness in space and time? Is Nature always simple and compositional, such that compactness and sparsity rule in natural systems? The issue of model's (i.e. hypotheses) intercomparison and evaluation also speaks to the more elusive question of how to reconcile solutions offered by different equation discovery methods.

Another important theoretical challenge is related to the fact that, very often, one (1) assumes that all involved state variables are given/observed, which resembles the sufficiency assumption in causal inference, and (2) selects a subset of representative states, assumes a particular basis to express the solution, or can operate on a manifold subspace [Champion et al., 2019, Brunton and Kutz, 2022, Daniels and Nemenman, 2015, Waltz and Buchanan, 2009]. Both are strong assumptions that challenge the process of discovering equations and raise many questions. How do we choose the right variables to include in the equation discovery method? How much does the solution change when a variable is omitted or added? Here, *identifiability* issues and *latent confounders* play a substantial role. And foremost, what if we cannot measure the underlying state variables? Can we identify them automatically? Several methods have arguably proposed to discover the latent variables [Chen et al., 2022], and other efforts exploit the link between RKHS techniques and a plausible master equation underlying the spatiotemporal evolution of the data probability density function [Coifman and Lafon, 2006] to automatically learn a set of fundamental coordinates of the system, cf. Sec. 3.2.5. The inferred subspace has shown effectiveness in identifying the dynamics of (partial) differential equations [Kemeth et al., 2022, Thiem et al., 2020], see examples in Sec. 4.1.2. However, like other explicit methods discussed in this section, it requires prior guidelines on the differential equation to model and incorporates a range of heuristics during its processing pipeline [Kemeth et al., 2022].

**Practical challenges.**
Important practical challenges are related to the *model development* and *data characteristics*: (1) *high dimensionality*, (2) *nonlinear relationships* and (3) *risk of overfitting*. In many cases, the number of variables and parameters can be very large, making it difficult to find the most relevant features and relationships. This is the scenario where non-identifiability arises. Another challenge

is that real-world systems often exhibit highly nonlinear relationships between input and output variables. This makes finding a good representation of the underlying dynamics difficult, and the search space for equations can become very large. While several nonlinear methods that capture complex dynamics exist (kernels, neural networks), performance evaluation and hyperparameter tuning are still important challenges. Finally, symbolic regression approaches can also suffer from overfitting, where the model fits the training data well but performs poorly on new data, which speaks to the trade-off between the model's accuracy and complexity.

When working with spatiotemporal data, it is relatively unclear how to incorporate information about time-lagged relations and interventional data. Both the explicit and the implicit approximations show particular challenges, though: On the one hand, symbolic regression techniques (such as SINDy) face significant problems in defining the basis functions, working in high-dimensional problems, and the impact of (even a limited) amount of noise. Other methods, like those based on AI Feynman, even if they incorporate sensible criteria to guide the equation discovery (like compositionally, reversibility or physical unit consistency), almost predict completely different equations when changing constants in the true equation [Suseela et al., 2022]. On the other hand, implicit methods (like DMD or Koopman operators) do not provide an explicit equation but a latent feature representation to explain system dynamics. These methods struggle with nonstationarities, nonlinearities, and gaps noise, which remain unresolved problems in the literature [Wu et al., 2021]. Note that similar challenges to those in causal discovery remain here, cf. Sec. 2.2. Besides, DMD methods fail to generalise outside the training data and violate basic physical laws. To alleviate this, integrating domain knowledge (such as symmetries, invariances and conservation laws) in DMD has been recently introduced as an effective, robust approach [Baddoo et al., 2023].

### 3.3.2 Opportunities

While symbolic regression presents challenges, it offers exciting research opportunities for the physical sciences. Three main opportunities can be identified: (1) *model interpretability*, (2) *model compression and evaluation*, and (3) *model selection*. Equation discovery is a step forward in the system's understanding. The field leverages fully interpretable models, and unlike causal models, equation discovery (symbolic regression) models are directly applicable predictive models. The discovered equations from data can provide insights into the underlying dynamics of a system, thus helping researchers better understand how different factors interact and contribute to a particular outcome. Even with implicit latent representations, interpretability can be accomplished with interventional analysis. Another interesting opportunity is model compression, as symbolic regression models can also be used to compress large datasets into simple equations that capture the system's essential features; this can make it easier to analyse and visualise the data and make it more computationally efficient to work with. Another practical opportunity of symbolic regression models is that they offer a reduced set of possible solutions typically ranked in amenable Pareto fronts, which of course, trigger difficulties in choosing the right model but also fruitful scientific discussions about the plausibility of identified relations. This can save researchers time and effort and help identify unexpected patterns and relationships in the data.

**Model interpretability and intervention analysis.**

The discovered explicit models are interpretable in nature. However, when complexity cannot be traded for accuracy or whenever an implicit feature representation is learned, intervention (or sensitivity) analysis offers opportunities for interpretability. For example, one can (1) *ignore or simplify the problem* by performing small perturbations away from real-world dynamics, which might help identify the proper relationship between variables; (2) *intervene on exogenous variables* (e.g. wind or solar irradiation, mixing coefficients, initial conditions in climate sciences, or targeted, direct brain stimulation in neurosciences) which is equivalent to collecting more data; (3) *create*

*a library of trajectories* under different conditions and select those which match the desired intervention; and (4) *intervene in the learned latent space* and decode the intervention back to input space, thus allowing us to generate interventions that follow the system's natural trajectories. Further analysis methods, e.g. for studying interventions of the learned ODE, might be fruitful, as they can access long-term dynamical properties, such as how distortion in the eigenvalues affects the system's stability as the phase space changes.

**Model compression and evaluation.**
Scientists frequently use metrics to evaluate new ideas or distinguish between competing hypotheses. As we have seen before, a governing equation should be simple but not simpler, accurate for prediction, robust under distortions and changes, and invariant in space and time. Equation discovery offers a *direct* way to learn plausible models and an *indirect* way to contrast and evaluate derived models. For this, one typically assesses (1) the *predictive accuracy* when answering how well the (simplest) hypothesis explains the data; (2) the model's *invariance* under distributional shift to account for the causal mechanisms; and (3) the *robustness* under interventions to study how the proposed process (or descriptive equation –representation) is consistent with interventions on the model dynamics, such as deactivation of components or targeted modification of exogenous variables. The latter is the rarest form of validation due to its high computational cost and difficulty in experimental design.

**Model selection.**
Enforcing sparsity in model selection can lead to unrealistically too simple models. That is why methods that can provide solutions along the Pareto line (Fig. 3.2) are needed to capture complex relationships and offer subsets of plausible model solutions. Alternative regularisation schemes will likely be important alongside profound estimates of uncertainty and extrapolation indicators.

# 4. Case studies in the physical sciences

This section gives concrete examples of applying different data-driven causal and equation discovery in important fields of the physical sciences: neuroscience, Earth and climate sciences, and fluid and mechanical dynamics, cf. Table 4.1.

Table 4.1: Case studies presented and the main methods used in this section.

|  | Neuroscience | Earth & climate | Fluid dynamics |
|---|---|---|---|
| Causal discovery | Causal connectivity (DCM, GC, TE, SCM) | Carbon-water interactions Climate model comparison (CCM, PCMCI) | – |
| Equation discovery | Learning trajectories (kFDA, GP, Variational Bayes RNN, Diffusion Maps) | Ocean Mesoscale closures (RVM, DMD, SINDy) | Turbulence understanding Vortex shedding (SINDy, Genetic Programming) |

## 4.1 Neuroscientific applications of physics-based machine learning

### 4.1.1 Overview of parsimonious models for neural population dynamics

Neuroscientific modelling falls within the remit of the field known as computational or theoretical neuroscience, which studies the transmission of information in the nervous system at multiple spatiotemporal scales (ranging from neuronal to whole-brain levels) in relation to perception, cognition, and behaviour (e.g., [Gerstner et al., 2014, Koch, 2004]). Thus, an ongoing challenge in computational neuroscience is to link biophysically detailed models operating at microscopic levels with meso/macroscopic theories of cortical processing [Rabinovich and Varona, 2018]. This enterprise is often addressed by deducting low-dimensional systems of partial differential equations or maps describing *coarse-grained* variables derived from collective neural responses. Such different neurobiologically plausible simplifications are commonly termed ensemble, population, neural-mass, or simply *firing-rate* models (see, for instance, [Byrne et al., 2021, Brunel, 2000]).

These synthesis efforts are intimately connected with empirically discovering a reduced dynamical system generating the observed neuronal activity. However, neurocomputational modelling

traditionally focused on analytical, deductive approaches mapping realistic cortical networks to *tissue-level* descriptions, as opposed to data-driven model discovery, reviewed in Sec. 3. Thus, ensemble models are typically principles-based, often hinged on assumptions about dynamical interactions arising within homogeneous pools of neurons (e.g., [Amari, 1977, Byrne et al., 2020, Tabas et al., 2019, Wilson and Cowan, 2021]).

Early neural ensemble models stemmed from applying statistical mechanical principles to the interaction of homogeneous pools of (excitatory and inhibitory) populations [Wilson and Cowan, 1972, Potthast, 2013]. Later, physics formalisms like the Fokker-Planck approach for describing the spatiotemporal evolution of the probability distribution of neuronal activity enabled theoretical neuroscientists to take a more holistic approach to identify mean-field approximations of networks of spiking neurons (e.g., [Gerstner et al., 2014, Mattia and Del Giudice, 2002, Brunel, 2000]). These and other nonlinear dynamical systems tools [Rabinovich et al., 2008] provided closed-from, exact solutions for the collective behaviour of neural populations [Montbrió et al., 2015, Mattia et al., 2019, Byrne et al., 2020] capable of an extensive dynamical repertoire, although strongly dependent on universal theoretical assumptions, given their deductive nature (see Sec. 1). Alternatively, a Laplacian assumption on this probability distribution resulted in neural mass descriptions [Marreiros et al., 2010], recently proposed as building blocks for whole-brain models with translational applications [Schirner et al., 2022].

Overall, these chiefly deductive approaches rendered compact models of differential equations based on *a priori* assumptions about neural and synaptic variables, fostering the interpretability of high-complex neuronal networks. By contrast, inferential approaches in neuroscience have been typically utilised to empirically identify neural dynamics underlying cognition and behaviour, as discussed next.

## 4.1.2   Empirical reconstruction of neuronal trajectories

There is an increasing focus on applying classic and deep machine learning approaches to reconstruct attracting and transient dynamics of cerebral cortex responses [Duncker and Sahani, 2021]. A neural trajectory $T$ $(n \times d)$ is often defined as the sequence of $n$ neural response vectors $x(t)$ embedded in a $d$-dimensional state-space (the *ambient* space), spanned by neural ensemble activity or proxies thereof (e.g., firing-rates, electromagnetic potentials), their lags and nonlinear transformations [Galgali et al., 2023, Duncker and Sahani, 2021, Balaguer-Ballester et al., 2011].

Traditionally, standard dimensionality reduction techniques (e.g., PCA, Multi-dimensional Scaling, Discriminant analysis etc., see Sec. 3.2) were directly applied for the visualisation of high-dimensional neural trajectories, showcasing coarse-grained aspects of firing-rate dynamics concerning, for example, cognitive decisions or motor functions [Cunningham and Yu, 2014, Hyman et al., 2012]. More recently, Gaussian processes provided a flexible approach to derive a low-dimensional manifold representing the dynamical systems generating the observed activity. They just require a reasonable hypothesis on temporal correlations between observations (the prior covariance function [Rasmussen and Williams, 2006]). For instance, Gaussian process-based factor analysis (GPFA) [Yu et al., 2009], and other latent-variable methods [Aoi et al., 2020], provide such low-dimensional subspace while simultaneously approximating the probability of spiking -without the compelling need for probability density estimation [Yu et al., 2009, Gokcen et al., 2022]. These and related approaches can identify latent neural trajectory manifolds in prefrontal and motor cortices underlying decision-making [Rutten et al., 2020, Aoi et al., 2020, Duncker and Sahani, 2018]. Specifically, recent GPFA variants were able to discern between competing models for the contribution of upstream areas to recurrent dynamics supporting decision-making in the monkey prefrontal cortex [Galgali et al., 2023].

In general, new developments in deep auto-encoders such as Latent Factor Analysis via Dynamical Systems (LFADS) enabled the successful reconstruction of single-trial spiking activity

[Pandarinath et al., 2018]. Moreover, related approaches such as preferential subspace identification (PSI) incorporate behavioural labels to inform the dimensionality reduction algorithm; effectively discerning *behaviourally relevant* dynamics from the general neuro-dynamical landscape [Sani et al., 2021]. Remarkably, recent self-supervised methods, also guided by behavioural observations, leverage contrastive learning and nonlinear ICA to produce consistent sub-spaces across experimental sessions and subjects [Schneider et al., 2023]. In addition, like previous approaches, they can operate in a supervised fashion, fostering decoding with respect to competing methods in a range of calcium and electrophysiological recordings in different species [Schneider et al., 2023].

Covariance (kernel) function methods can also recreate salient facets of cortical dynamics like attracting sets. To this end, they leverage delay-embedding techniques in RKHS spanned by neuronal correlations and their temporal structure [Balaguer-Ballester et al., 2011, 2014, Lapish et al., 2015, Balaguer-Ballester et al., 2020] for identifying compact manifolds mapping animal's choice with attracting sets of ensemble trajectories [Lapish et al., 2015, Balaguer-Ballester et al., 2020]. Figure 4.1 presents an illustrative example of these RKHS techniques for recreating neuronal trajectories underlying the effect of dopamine at the circuit level. This approach facilitated the evaluation of mechanistic theories of dopamine modulation during decision-making, which was challenging given the limitations of direct experimental manipulations [Lapish et al., 2015]. The figure shows the flow field of trajectories derived from the activity of neuronal constellations in the rodent anterior cingulate cortex. Interestingly, the dynamic landscape depicted during working memory tasks in an optimal RHKS can be approximately described as transients connecting multiple attracting sets mapping spatial choices (Fig. 4.1a). This robust multi-stable scenario is completely disrupted by high doses of amphetamine (a well-known trigger of dopamine release, Figure 4.1b), while it is enhanced by low doses (see [Lapish et al., 2015]), in line with long-standing theoretical predictions of biophysical models [Durstewitz and Seamans, 2008].

Further facets of brain dynamics, such as chaotic attractors, have been recently addressed with recurrent, piecewise-linear architectures, amenable to optimisation via back-propagation variants or Bayesian variational inference [Durstewitz, 2017, Brenner et al., 2022]. In these approaches, tractability is promoted by leveraging units' linearisation for approximating trajectory inference in a parsimonious, transparent fashion. Empowered by these characteristics, such linearised recurrent networks could infer pathological whole-brain dynamics from functional magnetic resonance imaging (fMRI) recordings [Koppe et al., 2019]. Moreover, recent developments of these methods embody biophysically-inspired computations such as dendritic processing, fostering their reconstruction capabilities of nonlinear dynamical systems [Brenner et al., 2022].

More broadly, classic and deep architectures can facilitate the Bayesian inference of the optimal range of biophysically realistic model parameters. This is a challenging task given the potentially high sensitivity of realistic networks to different parametrisations [Barrett et al., 2019]. Along these lines, approximate Bayesian computation (ABC) has been combined with connectionist approaches to identify parameters in models operating at multiple spatial scales, ranging from microscopic-level Hodking-Huxley-type single neurons [Lueckmann et al., 2017] to macroscopic, cognitive-level decision-making models [Boelts et al., 2022].

For instance, Sequential Neural Posterior Estimation (an ABC method) alternates between deep learners and variational Bayes approaches for parameter approximation. First, a standard (non-biophysical plausible) classifier is used to constrain the range of initial parameters $\theta$ generated from the prior $p(\theta)$ by predicting their suitability. Subsequently, a deep learner operating on multivariate data $x$ and parameters sampled from such constrained prior $\hat{p}(\theta)$ estimates the likelihood $\hat{p}(x|\theta)$, enabling a progressively more refined Bayesian computation of the posterior over neuronal and synaptic parameters $p(\theta|x)$ [Gonçalves et al., 2020]. These approaches were able, e.g., to discern between neuronal model configurations, essentially indistinguishable in observed activity, the ones metabolically optimal in the pyloric network in crustacean [Deistler et al., 2022]; or to infer reaction
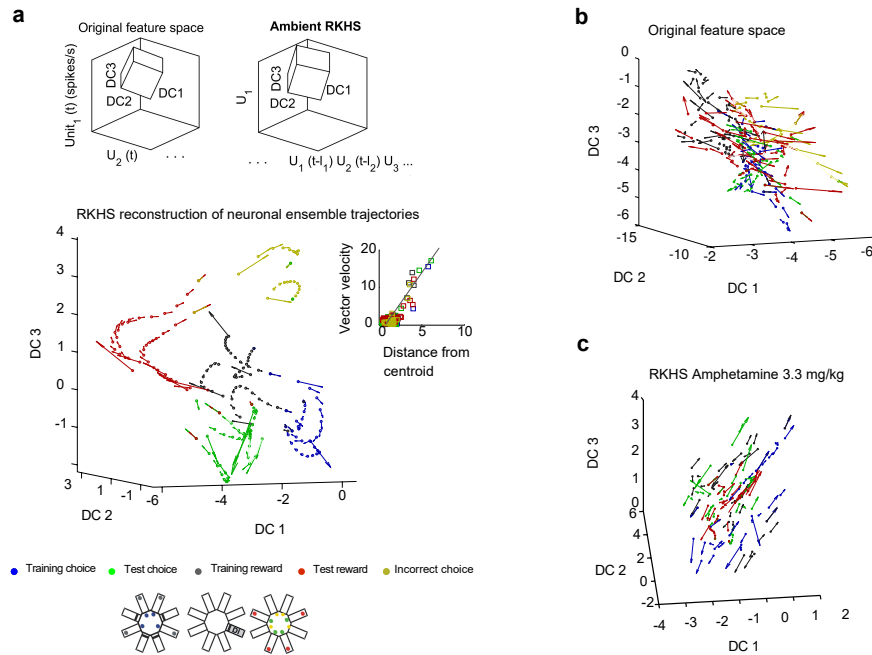
Figure 4.1: Example of the reconstruction of neural trajectories in an RKHS derived from rodent anterior cingulate cortex (ACC) multi-array recordings, taken from (Lapish et al., 2015). **a**. The flow field stems from projecting an ACC ensemble firing rates onto the three main coordinates of a discriminant subspace (DC1-DC3, computed here by kernel-fisher discriminant analysis orthogonalised for a faithful representation, details in, e.g., (Balaguer-Ballester et al., 2011, 2020, 2014, Lapish et al., 2015)). The DC1-DC3 is embedded into an optimal ambient high-dimensional RKHS (top schematics), spanned by neuronal ensemble firing rate and its higher-order correlations (up to 3rd order in this example) operating on a delay-coordinate map. The colour code (bottom) corresponds to rat choices in this experiment (the schematic of the radial-arm maze used in this experiment is taken from (Balaguer-Ballester et al., 2011, 2014)), which occupy distinct regions of the subspace. The flow field indicates faster shifts (large vector lengths) during the transition points, while it slows downs nearer the centroids of the clusters, suggesting an attracting-like dynamic landscape. The inset quantifies this uneven distribution of flow field speeds as a function of the distance to the centroid, further supporting this observation on flow convergence. **b**. This ordered phase-space structure cannot be achieved in the original feature space or with delay-coordinate maps. **c**. It is also destroyed when the animal receives a high dose of amphetamine, even in an optimal ambient space. Figure adapted from (Lapish et al., 2015) and from (Balaguer-Ballester et al., 2011) with publishers' permission (the Society for Neuroscience and the Public Library of Science).

times and choices in classic (drift-diffusion-type), descriptive models of decision-making [Boelts et al., 2022].

Alternatively, neural trajectory reconstruction of Hodking-Huxley ensembles has been recently tackled with nonlinear manifold learning in RKHS [Thiem et al., 2020, Kemeth et al., 2022]. These methods, like common dimensionality reduction-based techniques (see Sec. 3.2), identify first an optimally reduced set of coordinates from an original higher-dimensional ambient space embedding the time series $y$. However, by contrast with other approaches, components spanning the low-dimensional representation are typically lead eigenvectors of a discretised Laplace operator governing the spatiotemporal evolution of the underlying $p(y(x,t))$, where $x$ is homologous to a spatial coordinate. Thus, the reduced subspace spanned by the main non-redundant eigenvectors is often termed a Diffusion Map [Coifman and Lafon, 2006] (see also Sec. 3.2.5) given by the set of $i^{th}$ emergent coordinates $\{\phi(x)_i\}$.

In a subsequent stage, a fully connected (deep/shallow) architecture learns the dynamical system on the diffusion map, in other words, estimates the function $f$ that maps the temporal flow of the time series $y$ to its derivatives w.r.t. emerging coordinates, for instance, in a one-dimensional

diffusion map $\phi_1$,

$$\frac{\partial y(\phi_1, t)}{\partial t} \approx f\left(y, \frac{\partial y}{\partial \phi_1}, \frac{\partial^2 y}{\partial \phi_1^2}, \ldots, \frac{\partial^n y}{\phi_1^n}, \gamma\right), \tag{4.1}$$

where $\gamma$ is a set of parameters which enables the learned map to reproduce bifurcations.

Recently, diffusion maps and related approaches were also utilised for the one-dimensional reduction of neural recordings at multiple spatial scales, ranging from single-unit to macroscopic levels imaging data [Brennan et al., 2023, Nieh et al., 2021]. Such drastic simplification facilitates, for instance, the interpretability of decision-making models, while preserving a competitive prediction accuracy [Brennan et al., 2023].

Interestingly, when the argument of $f$ contains no derivatives and incorporates additive Gaussian noise, the system's dynamics reduces to the well-known Langevin stochastic differential equation, stemming from a Fokker-Planck process governing the temporal evolution of a probability distribution [Genkin and Engel, 2020]. Thus, by approximating a discretised Fokker-Plank operator, it is possible to empirically infer parameters of the stochastic process leveraging conventional likelihood estimation techniques [Genkin and Engel, 2020, Genkin et al., 2021]. Langevin dynamics fit many natural phenomena and is of special interest in high-level decision-making models in neuroscience. This formalism was recently used in [Genkin et al., 2021] for model discovery based on a one-dimensional Langevin equation and an additional stochastic spike generator,

$$\begin{aligned} \frac{dx(t)}{dt} &\approx D \cdot F(x) + \sqrt{2D} \cdot \psi(t; 0, 1)), \\ y &\sim Poisson(x; \lambda(t)), \end{aligned} \tag{4.2}$$

where $x(t)$ is a latent trajectory, $D$ is the diffusion constant, $\psi$ is white normal noise, $F$ is the deterministic map to be inferred, and $\lambda(t)$ is the parameter of an inhomogeneous Poisson process generating the observed spike train time series $y(t)$. This approach is capable of identifying a parsimonious model (that is, a set of $\{D, F(x), \lambda(t)\}$) from the spike train time series. Thus, it was used to discern between competing models of perceptual decision-making by comparing probability distributions underpinning such alternative parameter sets via standard Kullback-Lieber divergence [Genkin et al., 2021].

In short, neuroscientific studies typically conceive the discovery of biophysical laws as the inference of deterministic dynamics embedded in essentially stochastic neural processes. This goal has been interpreted either as empirically reconstructing attracting and transient components of neural activity (disentangled from coupled noise e.g., [Balaguer-Ballester et al., 2011, Rutten et al., 2020, Kemeth et al., 2022]); or as identifying parameters of parsimonious, *a priori* model shapes [Genkin et al., 2021, Lueckmann et al., 2017]. These approaches provide valuable insights on the latent neural dynamical landscape.

However, a liking theme in such inferential methods is that, despite their advances in tractability (e.g., [Brenner et al., 2022, Genkin et al., 2021]) and interpretability (e.g., [Balaguer-Ballester et al., 2020, Kemeth et al., 2022]), they are often not designed to empirically discern a unique set of differential equations, as it is popular in other areas of physics (e.g., [Brunton et al., 2016b, Guan et al., 2021, Loiseau, 2020], see examples in Sec. 4). This is at odds with the chief goal for deductive approaches, outlined in Section 4.1.1. Key reasons for this shortcoming might be found in the high-noise levels arising in intrinsically stochastically-dominated neural processes, in which most hidden variables are not experimentally accessible (especially in *in-vivo*) [Genkin et al., 2021, Balaguer-Ballester et al., 2011]. This challenging scenario hinders the direct application of approaches common in other fields and poses an intriguing question for future research endeavours.

## 4.2 Learning causally interacting brain regions

### 4.2.1 Causality in the connected brain

The debate on the causal role of brain connectivity has a long-standing tradition (see e.g. [Razi and Friston, 2016]). The classic view of functional segregation (mapping functions to physical brain regions) veered to connectionism, that is, brain functions result from interactions between neurocomputational units [Razi and Friston, 2016]. Consistently, the focus gradually shifted from functional segregation (the study of regionally-specific brain activation) to functional integration (the study of the connectivity between cortical areas [Razi and Friston, 2016]).

Historically, connectivity studies establish the distinction between structural (the anatomical location of white matter, axonal tracts), functional and *effective* connectivity (and sometimes with normative connectivity, in contrast to individual-specific connectomes) [Siddiqi et al., 2022]. This classification is relevant to determine the type of causality questions that can -or *cannot*- be addressed [Reid et al., 2019, Ross, 2015]. Typically, functional connectivity methods estimate statistical dependencies such as spatiotemporal correlations or coherence measures between brain ensembles. In contrast, a subset of these methods, commonly termed effective connectivity approaches, refer to the quantification of directed interactions between brain circuits (e.g., [Siddiqi et al., 2022, Reid et al., 2019, Tognoli and Kelso, 2014, Balaguer-Ballester et al., 2018]). In this arena, the quest for demonstrating causality relationships in neuroscience has attracted much attention over the recent decades [Siddiqi et al., 2022, Ross, 2015], and its plausibility has been widely debated [Weichwald and Peters, 2021, Barack et al., 2022, Barnett et al., 2018, Ross, 2015, Haufe et al., 2014]. For instance, in neuroimaging, multiple issues such as confounding factors [Woolgar et al., 2014, Todd et al., 2013] and varying temporal delays (intrinsic to, for instance, fMRI) challenge estimates of network information flow directionality (e.g., [Siddiqi et al., 2022, Weichwald and Peters, 2021, Haufe et al., 2014, Lohmann et al., 2012, Davis et al., 2014] among many others).

Methodologically, a key characteristic of causal approaches -in difference with conventional probabilistic modelling- is the need for predicting how the system reacts under interventions [Weichwald and Peters, 2021]; in other words, for defining counterfactual models (see Sec. 2.3.8). Problematically, a large amount of interventional data is necessary to falsify the wide range of causal hypotheses in a high-dimensional system like the cerebral cortex [Weichwald and Peters, 2021]. For instance, *targeted* brain interventions via intracranial electrical stimulation (iES) in conscious patients is typically a robust approach for testing causality [Siddiqi et al., 2022], but large-scale datasets using this experimental protocol are scarce, given ethical and experimental limitations of invasive techniques [Weichwald and Peters, 2021]. However, comprehensive, high-quality interventional data would be fundamental to falsify as many competing causal scenarios as possible. This is especially important in cognitive neuroscience given the lack of experimental access to some fundamental variables, which increases the number of plausible causal models underlying observable behaviour [Weichwald and Peters, 2021].

This shortage of comprehensive targeted lesion/stimulation datasets, and the improvement of whole-brain registration techniques, led to the development of analytical methods (or adaptations of existing ones) to better understand causality in cortical circuits. Most notably, Granger Causality and related approaches (GC, originated in the field of Economics [Granger, 1969]), Structural Causal Modelling (SCM [Pearl, 2009c, Bongers et al., 2018]), and Dynamic Causal Modelling (DCM [Friston et al., 2003]) have been extensively used, as will be discussed next.

### 4.2.2 Causal methods in neuroscience

GC and an extension of this concept, Transfer Entropy (TE), are perhaps the most common *model-free* methods for assessing causal relations in neuroscience. These two generalist approaches estimate the direction of causality between interacting neural populations by analysing the time

series derived from brain responses [Ding et al., 2006, Friston et al., 2013, Barnett et al., 2018]. They are regular statistical tools for studying orchestrated interactions between brain regions via magneto/electroencephalography (M/EEG) and fMRI recordings (e.g., [Friston et al., 2013, Stokes and Purdon, 2017, Bassett and Sporns, 2017]). At microscopic levels, they have also been applied to detect synaptic connections between neurons [Sheikhattar et al., 2018]. Specifically, GC is based on the assumption that time series prediction leveraging its past values significantly improves by inputting historical values from another, causally connected time series (see details in Sec. 2.1.3). Thus, the presence of causal relationships is detected by testing the hypothesis that one time series autocorrelations have predictive power for the other time series [Ding et al., 2006].

TE expands this idea to accommodate broader types of nonlinear temporal interactions by computing the amount of information that one time series *transfers* to another. Similarly to GC, it conjectures that the current value of one time series can be better estimated by conditioning the predictive probability to past values of both itself and another time series, inferring causality direction [Barnett et al., 2009]. Alternatively, SCM and its recent variants are Bayesian approaches for assessing plausible causal graphs in brain networks. They have been applied, for instance, to foster interpretability in behavioural decoding approaches [Weichwald et al., 2015]. However, their use in cognitive neuroscience is still challenging, given the key difficulties discussed in Section 4.2.1, and the indirect nature of most neuroimaging measurements (reviewed in [Weichwald and Peters, 2021]).

Accompanying these model-free approaches, perhaps the most standard model-based technique for connectivity inference between brain regions is Dynamic Causal Modelling (DCM). DCM is a Bayesian method incorporating different degrees of *a priori* biological plausibility for understanding mechanisms underlying neuroimaging data (see e.g.,[Penny et al., 2004, Stephan et al., 2010, Marreiros et al., 2010, Cooray et al., 2015]). It has been employed to study neural pathways of effective connectivity in e.g., motor control, attention, learning, decision-making, emotion, and other higher cognitive functions [Friston et al., 2019]; and even to model EEG seizure activity dynamics in epilepsy [Cooray et al., 2015].

In general, whole-brain modelling methods like DCM or other more recent models [Cabral et al., 2022] can provide a more nuanced understanding of the underlying mechanisms of brain function than model-free approaches [Friston, 2005]. However, the need for large datasets w.r.t. the complexity of the range of alternative models hampers the interpretation of the estimated connectivity [Siddiqi et al., 2022, Penny et al., 2004, Stephan et al., 2010, Marreiros et al., 2010]. Indeed, classic DCMs [Penny et al., 2010] have been criticised for the difficulty in falsifying their model selection approach [Lohmann et al., 2012] and perhaps for this reason, they were not extensively tested in clinical settings [Penny et al., 2004]. Specifically [Lohmann et al., 2012] suggested the ambiguity of DCM inference in generating a unique optimal connectivity map due to, e.g., known challenges in model fitting and selection in such a large space of possible architectures [Lohmann et al., 2012, Chicharro and Panzeri, 2014].

These caveats of DCM as a robust approach for causality assessment led to the development of variants such as spectral DCMs, the canonical microcircuit DCM -introducing higher degrees of laminar-specific, biophysical detail towards more informative priors for E/MEG modelling-, or the stochastic dynamic causal model, sDCM (see a review in [Friston et al., 2019]). sDCM incorporates random processes to the basic DCM equations, enhancing its fitting capability to hemodynamic responses and hence alleviating excessive dominance of priors in Bayes model selection [Bernal-Casas et al., 2013]. In a classic DCM for fMRI data, the neural state $\boldsymbol{x}(t) \in \{1, n\}$ (for $n$ interacting brain regions) corresponding to a single task-based input $u(t)$, is determined using the simple first-order differential equation $\frac{d\boldsymbol{x}(t)}{dt} = (\boldsymbol{A} + u(t) \cdot \boldsymbol{B}) \boldsymbol{x}(t) + u(t) \cdot \boldsymbol{c}$; where the matrix $\boldsymbol{A}$ encodes (endogenous) connections between brain regions, $\boldsymbol{B}$ the strength in which inputs modulate each connection (*modulatory* inputs) and $\boldsymbol{c}$ the gain of the *driving* inputs to each region.

sDCM expands this approach by adding intrinsic $\boldsymbol{\beta}(t)$ and extrinsic $\boldsymbol{\gamma}(t)$ stochastic fluctuations to account for the incomplete observability of both states and inputs to brain areas relevant to the cognitive task:

$$\frac{d\boldsymbol{x}(t)}{dt} = (\boldsymbol{A} + \boldsymbol{v}(t) \cdot \boldsymbol{B}) \, \boldsymbol{x}(t) + \boldsymbol{v}(t) \cdot \boldsymbol{c} + \boldsymbol{\beta}(t),$$
$$\boldsymbol{v}(t) = u(t) + \boldsymbol{\gamma}(t),$$
$$\boldsymbol{y}(t) = g(\boldsymbol{\theta}) * \boldsymbol{x}(t) + \boldsymbol{\epsilon}(t),$$

(4.3)

where $\boldsymbol{v}(t)$ is a *hidden* input cause masked by fluctuations (univariate here for simplicity), and the last equation represents a hemodynamic model (present in all DCMs variants) of non-neural parameters $\{\boldsymbol{\theta}, \boldsymbol{\epsilon}(t)\}$. Finally, their convolution with the neural state $\boldsymbol{x}(t)$, yields the observed fMRI blood-oxygenation level-dependent (BOLD) response $\boldsymbol{y}(t)$ in relevant brain areas (termed regions of interest, ROI) [Friston et al., 2019].

Figure 4.2 summarises an illustrative reliability study from [Bernal-Casas et al., 2013], specially designed for assessing sDCM robustness. In this example, a large sample of participants (*n*=180) was recruited from three different geographical locations. fMRI recordings were obtained from healthy subjects from the same age range while performing a classic 2-Back working memory task (to recall numbers shown two trials before). This N-Back task activates the dorsolateral prefrontal cortex-hippocampal formation (DLFC-HF) network connectivity, which is abnormal in schizophrenia patients [Meyer-Lindenberg et al., 2005]. Bernal-Casas et al. [2013] used sDCM to identify the DLPFC-HF effective connectivity and compared the consistency of the models in this multi-centre setting (Figure 4.2). Three *a priori* likely mechanisms to explain the BOLD responses to this task were implemented in three different families of models: with only driving inputs to the two regions ($B \equiv 0$, Figure 4.2a, left), only connectivity modulation ($C \equiv 0$, Figure 4.2a, centre) and both mechanisms combined (Figure 4.2a, right). Noticeably, the random effects Bayes Model selection process strongly favours a specific connectivity model belonging to the driving inputs family (Figure 4.2b) over all the rest, consistently for the three independent locations. Specifically, the DLPFC-HF connectivity parameters were statistically indistinguishable across datasets (Figure 4.2b), supporting the reliability of sDCM results.

In line with these whole-brain analyses, other studies showcased the consistency of causality methods at microscopic levels. As a representative example, [Chen et al., 2023] recently proposed the effectiveness of GC in inferring information directionality in zebrafish motor circuits from single-cell calcium imaging signals. Causally strong, interventional data was inaccessible in this setting. Despite this, results were in full agreement with the known physiology of this species. In addition, and besides these standard methods (GC, TE and DCM), recent approaches have also addressed the causality robustness question from different angles, for instance, by focusing on changes in information *reversibility* as a sign of aberrant resting-state brain dynamics -which could subserve as a biomarker of Alzheimer's disease [Cruzat et al., 2023].

These and other even more indirect causality measures (like standard statistical approaches [Reid et al., 2019]) have provided useful insights when operating on neurophysiological recordings with high temporal precision. The ideal recording modalities are thus those capable of directly recording local (electrical) field potentials (LFP), such as intracranial electroencephalography (iEEG) or neuronal-level techniques (like in Chen et al. [2023]). However, when indirect causality measures are estimated from other modalities -especially from functional imaging- the multiple confounders discussed earlier rank them in a weak position in a causality scale when compared with approaches based on interventions [Siddiqi et al., 2022]. Therefore, their capability for providing a reliable indication of causality interactions is highly disputed [Mehler and Kording, 2018, Siddiqi et al., 2022, Lohmann et al., 2012].

Nevertheless, limitations for assessing causal relationships in neuroimaging do not preclude
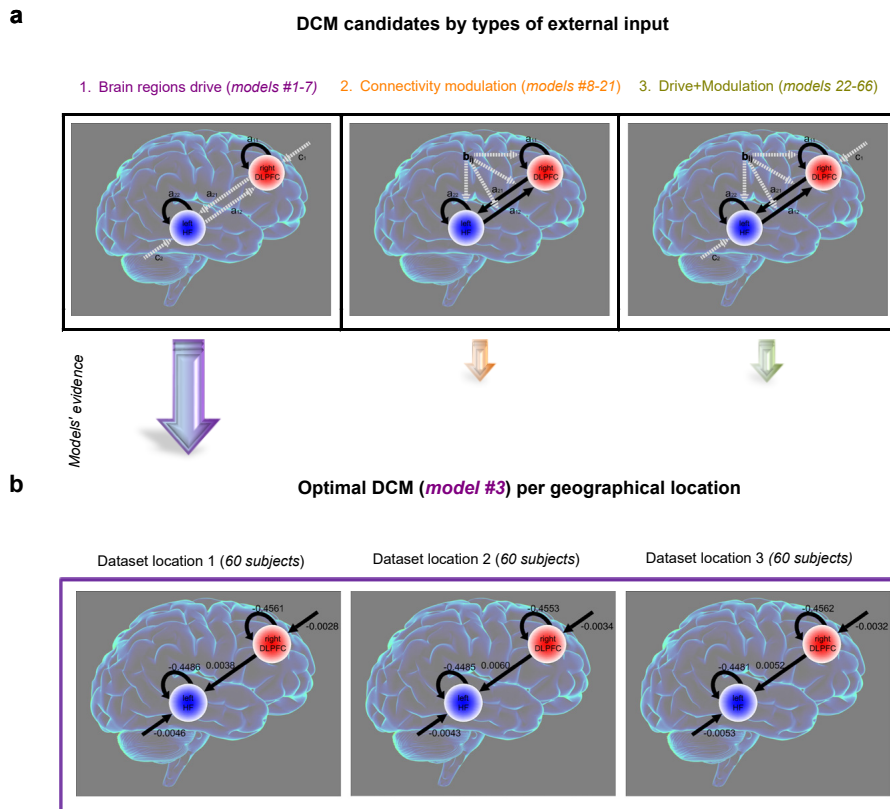
Figure 4.2: Example of consistency assessment for stochastic DCM, taken from (Bernal-Casas et al., 2013). **a**. Connectivity hypotheses associated with the performance of a 2-Back working memory task by healthy human participants (Bernal-Casas et al., 2013). Left: input fluctuations to the two ROIs (right DLPFC and left HF; 7 distinct model combinations). Centre: input variance modulates the connections themselves (14 models). Right: combinations of both mechanisms (44 models). Connectivity hypotheses are tested on independent datasets collected from three locations (Bonn, Berlin and Mannheim, 60 subjects each). Model evidence (log-likelihood marginalised over model's free parameters (Penny et al., 2004)) is much stronger for a model of the first family (#3, random effects Bayes factor 98%+ in favour of this model). **b**. Remarkably, connectivity parameters do not differ across sites (Friedman non-parametric test, $p > 0.2$), and interaction between sites and model parameters were not found ($p > 0.8$), supporting the model's robustness. Figure adapted from (Bernal-Casas et al., 2013) with the publisher's permission (Elsevier).

indirect analytical approaches to constrain the universe of plausible causal graphs for a specific scenario [Reid et al., 2019, Siddiqi et al., 2022]. Thus, there is a reasonable consensus in considering them as valuable contributors to strengthen causality claims, provided they are combined with more direct measures of causality based on interventional data [Reid et al., 2019, Siddiqi et al., 2022]. Indeed, the ideal scenario from a causality perspective occurs when its inference is consistent throughout different approaches; in other words, when different methods having complementary views provide synergistic evidence [Siddiqi et al., 2022]. For instance, converging evidence between a causal fMRI/EEG model, a targeted lesion, and the stimulation of a specific cortical circuit would score high on a causality scale than either of these methods alone; since the counterfactual could be established by focal stimulation [Neumann et al., 2023] combined with real-time neuroimaging recordings [Siddiqi et al., 2022], enriching the conclusions of the lesion study.

This optimal coalescence of multiple causal approaches for effective connectivity inference was termed Convergent Causal Mapping and is the recommended approach for designing new experiments [Siddiqi et al., 2022]. Thus, from this concerted perspective, studies considering a single approach in isolation -especially if it is not based on interventions- should not make strong

translational claims, like suggesting direct therapeutic applications [Siddiqi et al., 2022, Neumann et al., 2023]. In addition, future works should consider testing robustness to different environments [Weichwald and Peters, 2021] (like in the example shown in Figure 4.2 [Bernal-Casas et al., 2013]) for further reinforcing the credibility of the inferred causal flow.

## 4.3  Learning causal graphs of carbon and water fluxes

### 4.3.1  Introduction

The Earth is a highly complex, dynamic, and networked system where very different physical, chemical and biological processes interact in and across several spheres. Land and atmosphere are tightly coupled systems interacting at different spatial and temporal scales [Diaz et al., 2022]. The main challenge to quantifying such relations globally comes from the lack of sufficient in-situ measurements and the fact that some of these variables are latent and not directly observable with remote sensing systems. One can, for example, measure SM but not GPP directly. As an alternative, many studies have relied on model simulations to investigate SM-precipitation [Koster et al., 2006], GPP-SM [Green et al., 2019] and ET-SM relations [Milly, 1992, Jung et al., 2010], to name just a few. However, assuming a model implies assuming the knowledge of the causal mechanisms and relations governing the system. This is not necessarily a correct assumption, especially in model misspecification, non-linearities and non-stationarities. Discovering such relations from data is of paramount relevance in these cases. In the following, we review the performance of two standard methods of causal discovery from time series data to learn the relationships between environmental factors and carbon and heat and energy fluxes at the local (site, flux tower) level and the global (planetary, product-derived) level. At the local level, we exploit data acquire by eddy-covariance instruments estimating fluxes exchange. At the global level, we exploit Earth observation data from satellite observations.

### 4.3.2  Clustering of biosphere-atmosphere causal graphs at the site level

The atmosphere and terrestrial ecosystems constitute another closely interconnected complex system where processes interact across a range of temporal and spatial scales. Further, causal relations also depend on vegetation types, climatic regions, and the season. Fortunately, measurement campaigns of the past decades have resulted in good coverage of measurement sites, available in the FLUXNET database [Baldocchi, 2014], a collection of long-term global observations of biosphere-atmosphere fluxes measured via the eddy covariance method. Runge et al. [2023] discuss a similar case study in-depth.

Here we review the study of Krich et al. [2021] that analysed causal networks for different seasons at eddy covariance flux tower observations in the FLUXNET network and how they depend on meteorological conditions. Figure 4.3 explains the methodological setup. From a selection of 119 FLUXNET sites (Fig. 4.3(a)) daily time series data of the following variables were considered (see Fig. 4.3(b) for one site): short-wave downward radiation (or global radiation, Rg), air temperature (T), net ecosystem exchange (NEE) (inverted), vapour pressure deficit (VPD), sensible heat (H), latent heat flux (LE), gross primary productivity (GPP), precipitation (P), and soil water content (SWC). For details on data processing, we refer to Krich et al. [2021].

Causal networks were then estimated with PCMCI [Runge et al., 2019b] (time lags from 0 to 5 days) in sliding windows of 3 months to capture the temporal evolution of biosphere-atmosphere interactions. Based on findings in Krich et al. [2020], a smoothed seasonal mean was subtracted to remove the common driver influence of the seasonal cycle. This results in 10.038 networks for the different months and sites (an example network is shown in Fig. 4.3(c)). Node colours indicate the level of autocorrelation (auto-MCI-partial correlation [Runge et al., 2019b]), and link colours the cross-link strength (cross-MCI); time lags are indicated by small labels, and straight edges

Figure 4.3: Clustering causal graphs of local measurement stations of biosphere-atmosphere interactions (adapted from Krich et al. (2021)). See the main text for explanations.

are contemporaneous. Since the strongest and most consistent links are contemporaneous, further analysis focused on these 15 links.

A previous study [Krich et al., 2020] discussed individual networks in more detail; the scope of Krich et al. [2021] was to apply a dimension reduction, here t-distributed stochastic neighbour embedding (t-SNE [Van der Maaten and Hinton, 2008]) which considers each of the causal graphs as an observation in a high-dimensional space of the contemporaneous MCI partial correlation values (Fig. 4.3(d)). t-SNE allows projecting this high-dimensional space onto two dimensions (Fig. 4.3(e, left)) that are the dominant features of transitions between different states of biosphere-atmosphere interactions. The coloured clusters in Fig. 4.3(e, left) are based on the OPTICS approach [Ankerst et al., 1999], and the four corners indicate the four archetypes of network connectivity and the networks' underlying meteorological conditions (averages taken over the sliding windows in Fig. 4.3(b)). Finally, Fig. 4.3(e, right) shows the convex hulls of clusters and their average network.

Each point of the low-dimensional embedding represents a specific ecosystem's biosphere-atmosphere interactions at a specific time and allows us to investigate their behaviour. A main finding of Krich et al. [2021] was that ecosystems from different climate zones or vegetation types have similar biosphere-atmosphere interactions if their meteorological conditions are similar. For example, temperate and high-latitude ecosystems feature similar networks to tropical forests during peak productivity. During droughts, both ecosystems behave more like typical Mediterranean

ecosystems during their dry season. Such meta-analyses of causal networks allow for another perspective on understanding ecosystems, including an analysis of anomalous changes in network structure as indicators of ecosystem shifts (see Sec. 2.3.4).

### 4.3.3 Causal relations at global scale

As an alternative to Granger causality, the work [Sugihara et al., 2012] presented the convergent cross-mapping (CCM) method, which may deal with the issues of non-stationary and nonlinear processes and deterministic relations in dynamic systems with weak to moderate cause-effect variable coupling. CCM assesses the reconstruction of a variable's state space using time embeddings to determine if $X \rightarrow Y$. This method has been extended to account for causal relations operating at different time lags and applied to various research areas. However, it is sensitive to noise levels, hyperparameter selection, and false detections in strong, unidirectional variable coupling cases. To address these issues, the robust CCM (RCCM) [Diaz et al., 2022] alternatively relies on bootstrap resampling through time and the derivation of more stringent cross-map skill scores. The method also exploited the information-geometric causal inference (IGCI) method in [Mønster et al., 2017] to infer weak and strong causal relationships between variables and estimate the embedding dimension to derive global maps of causal relations.

Let us exemplify the RCCM method to discover interactions of three key variables in the carbon cycle: moisture, photosynthesis and air temperature (Tair). For that, we use data compiled in the Earth System Data Lab (ESDL), which contains harmonised products with a spatial resolution of $0.25°$ and a temporal resolution of 8 days, spanning over 11 years from 2001 to 2011. The RCCM method is applied in each grid cell, which allows us to infer spatial patterns of causal relations between several key variables of the carbon and water cycles.



Figure 4.4: Applying RCCM in (Diaz et al., 2022) to discover causal relations between GPP, Tair and SM. GPP drives Tair in cold ecosystems; Tair controls SM in water-limited areas; GPP dominates SM. Croplands were masked to avoid interference from human activity.

Fig. 4.4 shows GPP drives Tair mostly in cold ecosystems due to changes in land surface albedo. Results show GPP is an important forcing of local temperature in many areas. Recent studies have found temperature is an important factor of GPP, driven by radiative factors in cold climates

and turbulent energy fluxes in warmer, drier ecosystems. SM and Tair are closely linked, limiting evaporation and raising Tair under dry conditions. This could explain the significant impact of Tair in high latitudes. GPP is mainly influenced by Tair in water-limited regions, especially in high northern latitudes where cold temperatures limit photosynthesis and plant growth. GPP and ET are tightly related as carbon assimilation in plants is linked with water losses through transpiration [Field et al., 1995]. Low water availability reduces GPP and ET, causing increased air and surface temperatures and a drier atmosphere. SM being stronger than GPP is mostly seen in transitional wet/dry climates [Koster et al., 2004]. No strong forcings in tropical rainforest areas indicate GPP is mostly driven by solar radiation and affected by high VPD values [Madani et al., 2020].

## 4.4 Causal climate model intercomparison

As introduced in Sec. 2.3.7, causal inference can help to assess the output of physical models and evaluate and compare them against observations at the level of causal dependencies [Eyring et al., 2020, 2019, Nowack et al., 2020, Pérez-Suay and Camps-Valls, 2019]. Climate models [Eyring et al., 2016] provide short-term predictions and future climate projections under given anthropogenic emission scenarios and are the basis for climate-related decision-making. As models can only provide an approximation to the real system, it is essential to evaluate them against observations. Such climate model evaluation is largely based on means, climatologies, or spectral properties [Stocker et al., 2013, Eyring et al., 2020]. Here the problem of equifinality may occur: even though a particular model might well fit descriptive statistics of the data, the model might not well simulate the causal physical mechanisms that produce this statistic, given that multiple model formulations and parameterisations, even when wrong, can fit the observations equally well. The issue is that such models would lead to erroneous future projections–a causal problem of out-of-distribution prediction. Causal model evaluation [Runge et al., 2019a] can evaluate the ability of models to simulate the causal interdependencies of its subprocesses in a process-based model evaluation framework [Maraun et al., 2017].

Here we briefly summarise one approach in this direction [Nowack et al., 2020]. The author aimed to compare causal networks among regional subprocesses in sea-level pressure between observations and climate models of the CMIP ensemble [Eyring et al., 2016]. Figure 4.5**a-d** illustrates the method's steps. First, the regional subprocesses were constructed from gridded climate time series (daily-mean sea level pressure from the NCEP-NCAR reanalysis [Kalnay et al., 1996]) using Varimax principal component analysis (PCA) to obtain a set of regionally confined climate modes of variability (Fig. 4.5**b**). The Varimax-PCA weights were then applied to the pressure data from each climate model (the regional weights' cores are indicated in red). Each component is associated with a time series (3-day averaged) and is one of the causal network nodes. Then the causal discovery method PCMCI [Runge et al., 2019b] was applied to these time series to reconstruct the lagged time series graph among these nodes, which constitute characteristic causal fingerprints (Fig. 4.5**c,d**) for the observational data as well as the individual models. Node colours indicate the level of autocorrelation (auto-MCI-partial correlation [Runge et al., 2019b]), and link colours the cross-link strength (cross-MCI); time lags are indicated by small labels. Only the around 200 most significant links are shown.

These causal fingerprints can then be used for model evaluation and intercomparison. Figure 4.5**e** depicts a comparison of models among each other, that is, the matrix of average F1-scores for pair-wise network comparisons between ensemble members of 20 climate models (labelled following CMIP-nomenclature in capital letters) for simulations spanning approximately the historical period from 1948 to 2017 and two surrogate models (Random, Independent). The rows show the models taken as references in each case, and the columns indicate the models compared to these references. Higher scores imply a better agreement between networks, i.e., that the two models are more similar regarding their causal fingerprint. One can see that causal fingerprints

Figure 4.5: Causal climate model evaluation (adapted from Nowack et al. (2020)). See the main text for explanations.

from different ensemble members of the same model (diagonal in Fig. 4.5e) are more consistent than networks estimated from two different models (off-diagonal). The blocks are consistent with different models sharing a common development background. In Figure 4.5f, the models' causal fingerprints are each compared to the fingerprint of the observational data (ordered F1-scores). The result is a continuum of more- and less-similar models (but models have significantly different causal fingerprints). The networks can be further investigated to analyse which regional interactions the models differ more from observations.

Causal model evaluation can provide important information to model developers on where their models can be improved. Furthermore, Nowack et al. [2020] show that more realistic fingerprints also affect projected changes in land surface precipitation. Hence, causal model analyses could

Figure 4.6: Results for learning the density functional for Lennard-Jones fluids. (a) shows the equation of state $P(\rho)$ (pressure) for different interaction strengths $\varepsilon$ comparing Monte-Carlo simulations (MC) with the Machine learning (ML) results. (b) density profile for $\varepsilon = 1.25$, $\mu = \ln(1.15)$ inside the training region, but $V$ is not in the training data. (c) density profile at a hard wall for $\varepsilon = 1.9$, $\mu = \ln(1.9)$ (outside the training region $\varepsilon \in [0.5, 1.5]$). Dark solid lines are simulation profiles, and blue dashed lines are ML results. Insets in (b) and (c) show $\Delta\rho = \rho^{\mathrm{mc}} - \rho^{\mathrm{ml}}$.

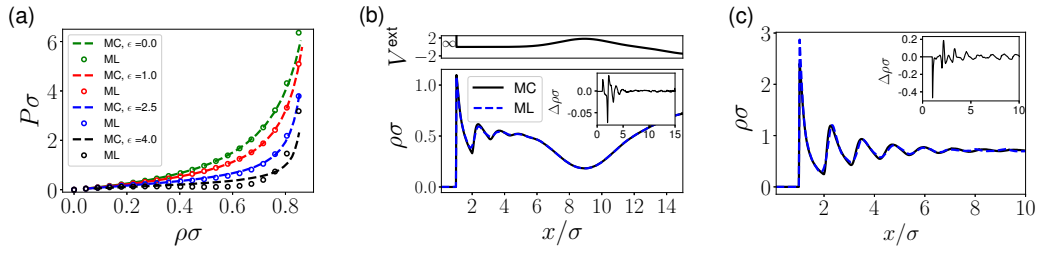be used to constrain climate change projections. The assumption is that the underlying physical processes (e.g., large-scale circulation) lead to dynamical coupling mechanisms captured in the causal fingerprints. One may now argue that high modelling skill on historical data is also relevant for modelling future changes if the physical processes remain important under future climate change.

## 4.5 Learning density functionals

Being able to describe many-body systems is exciting and important for many applications. Density functional theory Evans [1979] (DFT) is an approach to creating a description for classical and quantum many-body systems in equilibrium. The aim is to find a unique (free) energy functional that gives rise to the particle density profile. The analytical form of the (free) energy functional is generally unknown, except for a handful of particular model systems. One way to treat more complex systems is to perform computer simulations and learn the energy functional via machine learning. The first attempts in classical DFT used a convolutional network [Lin and Oettel, 2019], which does not allow much theoretical insight.

In Lin et al. [2020], the above-mentioned symbolic regression method, EQL [Sahoo et al., 2018], was adapted to represent part of the energy function. This is an interesting application, as the problem contains known parts of the computational pipeline that we do not want to replace and other parts that should be replaced via the data-driven approach. The fact that EQL can be embedded into any differentiable computational structure is crucial here.

The problem can be formulated as a self-consistency equation:

$$\rho(x) = \exp\left( \mu - \left.\frac{\delta F(\rho(x))}{\delta\rho}\right|_{\rho=\rho^{eq}} - V \right), \tag{4.4}$$

where $\rho$ is the particle density, $F$ is the external free energy functional that needs to be learned, and $\mu$ and $V$ are chemical and external potential, respectively. Notice that the derivative of $F$ (which is represented by an EQL network) is used in the equation. An analytical description for $F$ can be obtained using symbolic regression on simulation data. In Lin et al. [2020], for the case of hard rod particles and Lennard-Jones fluids, solutions were found that extrapolate well to unseen situations (different external potential or mean density), as shown in Fig. 4.6. It is a promising approach to gain more theoretical insights when applied to less studied systems.

## 4.6 Discovering governing equations in boundary-layer transition to turbulence

A classical approach to discovering the governing equations of a reduced-order model (ROM) describing a particular phenomenon, for which the governing partial differential equations (PDEs)

are known, is to perform Galerkin projection [Wang et al., 2012, Noack et al., 2003]. In Galerkin projection, a set of orthogonal basis modes (obtained, for instance, via POD) are used to develop a ROM of the system from data. Then, the governing PDEs are projected onto these modes, transforming the PDEs into a system of ordinary differential equations (ODEs) governing the dynamics of the temporal coefficients associated with those modes [Brunton and Kutz, 2022].

For incompressible fluid flows, the spatiotemporal velocity vector $\boldsymbol{u}(\boldsymbol{x},t)$ (where $\boldsymbol{x}$ are the spatial coordinates and $t$ time) can be expressed as follows after performing POD:

$$\boldsymbol{u}(\boldsymbol{x},t) \simeq \boldsymbol{u}_0(\boldsymbol{x}) + \sum_{k=1}^{r} a_k(t)\boldsymbol{u}_k(\boldsymbol{x}), \tag{4.5}$$

where $\boldsymbol{u}_0$ is the mean flow, $\boldsymbol{u}_k$ are the spatial modes, $a_k(t)$ are the temporal coefficients and $r$ is the number of retained modes in the ROM. The expansion (4.5) is then substituted into the governing PDEs, *i.e.* the incompressible Navier–Stokes equations, taking advantage that the POD modes are linear combinations of the instantaneous flow realisations (thus satisfying the boundary conditions) and are solenoidal (*, i.e.* divergence-free, due to the incompressibility condition). It is then possible to take an inner product in space with $\boldsymbol{u}_i(\boldsymbol{x})$. Since the POD modes are orthogonal, a set of ODEs can be obtained for the time derivatives of the temporal coefficients $\mathrm{d}a_i(t)/\mathrm{d}t$ as a function of the spatial modes and also $a_i(t)$. Despite being a widely-used method, it has the limitations of requiring knowledge of the underlying PDEs, and it also may exhibit convergence problems in more challenging scenarios. As discussed in Sec. 3.2, another alternative to produce a ROM for physical systems in a purely data-driven way is dynamic-mode decomposition (DMD), in which the obtained modes are orthogonal in time [Kutz et al., 2016]; note that this approach, in its compressive version [Bai et al., 2020], shares similarities with the eigensystem-realisation algorithm (ERA) [Juang and Pappa, 1985]. In this sense, DMD and its connections with the Koopman operator [Rowley et al., 2009a] were exploited by Eivazi *et al.* [Eivazi et al., 2021] to reproduce the dynamics of a near-wall model of turbulence [Moehlis et al., 2004] by using external forcing to reproduce the nonlinear behaviour of the system [Khodkar et al., 2019, Brunton et al., 2017].

In addition to the approaches mentioned above, other techniques enable learning the equations of motion just from data, as discussed in the early work by Crutchfield *et al.* [Crutchfield and McNamara, 1987]. These approaches typically rely on a library of candidate functions to build the resulting governing equation and solve an optimisation problem to obtain the expression that best represents the data. Note that it is essential to use any knowledge on the physical properties of the analysed data to inform the library (e.g. whether non-linearities, periodicities, etc. are present in the system that produced the data under study), as well as to define the relevant state variables, sampling rate, the initial set of parameters defining relevant trajectories, etc. Embedding prior physical information into the obtained model is crucial for the success of these approaches, and failing to do so may lead to rate and even incorrect models [Antonelli et al., 2022b]. Furthermore, these approaches typically suffer from the curse of dimensionality [Gelß et al., 2019], making it even more important to make the right choices in the library of candidate functions to ensure convergence; thus, being able to define the best basis functions to reduce the dimensionality of the system while retaining the most relevant physics is also is critical. Generally, only after solving the optimisation problem is it possible to assess which terms in the library are necessary and which ones may be combined, a fact that complicates *a-priori* equation discovery.

SINDy has been successfully applied to boundary-value problems [Shea et al., 2021] using forcing functions and performs well even with noise. Furthermore, SINDy has produced very successful results in a wide variety of fluid-mechanics problems, ranging from thermal convection [Loiseau, 2020], chaotic electroconvection [Guan et al., 2021], the so-called "fluidic-pinball" problem [Deng et al., 2020], turbulent wakes [Callaham et al., 2021b] and ROM development [Callaham et al.,

2022]. Interestingly, SINDy has also been successfully combined with autoencoders to discover low-dimensional latent spaces [Champion et al., 2019], benefiting from the non-linear data-compression capabilities of the latter and the interpretability of the former. This is certainly a promising direction to discover hidden complex relations in fluid-flow data and other high-dimensional physical systems, which requires further investigation, particularly when obtaining deeper insight into the interpretation of the latent variables.

Besides the methods above based on discovering nonlinear dynamical systems, other strategies exist to obtain equations from data. For instance, gene-expression programming (GEP), a branch of evolutionary computing [Koza, 1992], is based on having a population of candidate functions to build the solution that best approximates the data and progressively improving this population by the survival of the fittest. The main advantage of this approach is that it leads to closed-form equations, even for data where the governing equation is unknown. In principle, it leads to interpretable solutions (although, in some cases, the resulting equations are so convoluted that interpretability is complicated). GEP has been used to model turbulence [Weatheritt and Sandberg, 2016], particularly in the context of the so-called Reynolds-averaged Navier–Stokes (RANS) equations. In short, the RANS equations are obtained after decomposing the instantaneous velocity into a mean and a fluctuation component (Reynolds averaging [Reynolds, 1895]), and although this simplifies the flow-simulation process (RANS approaches are widely used in industry), the so-called closure problem emerges [Pope, 2000, Tennekes and Lumley, 1972]. This problem is associated with the unknown impact of turbulent fluctuations on the mean flow. All the existing models for these stresses are empirical, which precludes RANS simulations from producing accurate results for arbitrary flow cases. In this context, Weatheritt and Sandberg [Weatheritt and Sandberg, 2017] used GEP to obtain general expressions for these turbulent stresses in various cases, including turbulent ducts [Vinuesa et al., 2018], which are challenging for RANS models due to the presence of secondary flows. They achieved quite successful RANS models for the secondary flows. GEP effectively obtained more general expressions for the turbulent stresses than those in the classical literature [Boussinesq, 1923], a critical step for RANS models to produce a reasonable performance for complex flows [Spalart, 2000].

Finally, we conclude with a technique that is not aimed at discovering equations from data but rather focuses on identifying the dominant terms in the equations for various geometrical regions in the domain under study, given the available data. This method is based on data-driven balance models [Callaham et al., 2021a]. It can help improve our system's physical interpretation by understanding the most relevant terms defining various mechanisms in the data, particularly in non-asymptotic cases where the negligible terms are not obvious. Using unsupervised learning, the authors sought clusters of points in the domain with negligible covariance in directions that represent terms with a negligible contribution to the physics, a condition equivalent to stating that the equation is satisfied by a few dominant terms within the cluster. In particular, they used Gaussian-mixture models (GMMs) [Bishop, 2006] to cluster the data. Then they obtained a sparse approximation in the direction of maximum variance using sparse principal-component analysis (SPCA) [Zou et al., 2012]. Callaham *et al.* [Callaham et al., 2021a] show the applicability of this framework to a wide range of problems, including turbulence transition, combustion, nonlinear optics, geophysical fluids, and neuroscience. In all these cases, they obtained relevant insight into the governing equations, which can help uncover novel and unexpected physical relations. In particular, their application to the case of transition to turbulence is very illustrative, as shown in Figure 4.7. This figure shows that starting from high-fidelity turbulence data, the RANS equations [Pope, 2000, Tennekes and Lumley, 1972] mentioned above are obtained and their terms analysed. Visualisation of the clustering in equation space reveals some interesting relations, such as the high covariance of the so-called viscous and Reynolds shear-stress-gradient terms, $\nu \nabla^2 \overline{u}$ and $\overline{(u'v')}_y$ respectively, which identify the viscous sublayer in the domain. Note that subscripts here denote partial derivatives with

respect to the corresponding spatial variable, the overbar indicates time averaging, and the prime is used for fluctuating quantities. The inertial sublayer, which would correspond to the turbulent region in this boundary-layer flow, would be dominated by the convection of the mean flow $\bar{u}\,\bar{u}_x$ and $\overline{(u'v')}_y$, which are again correctly identified through a strong covariance and highlighted in the corresponding region of the domain.
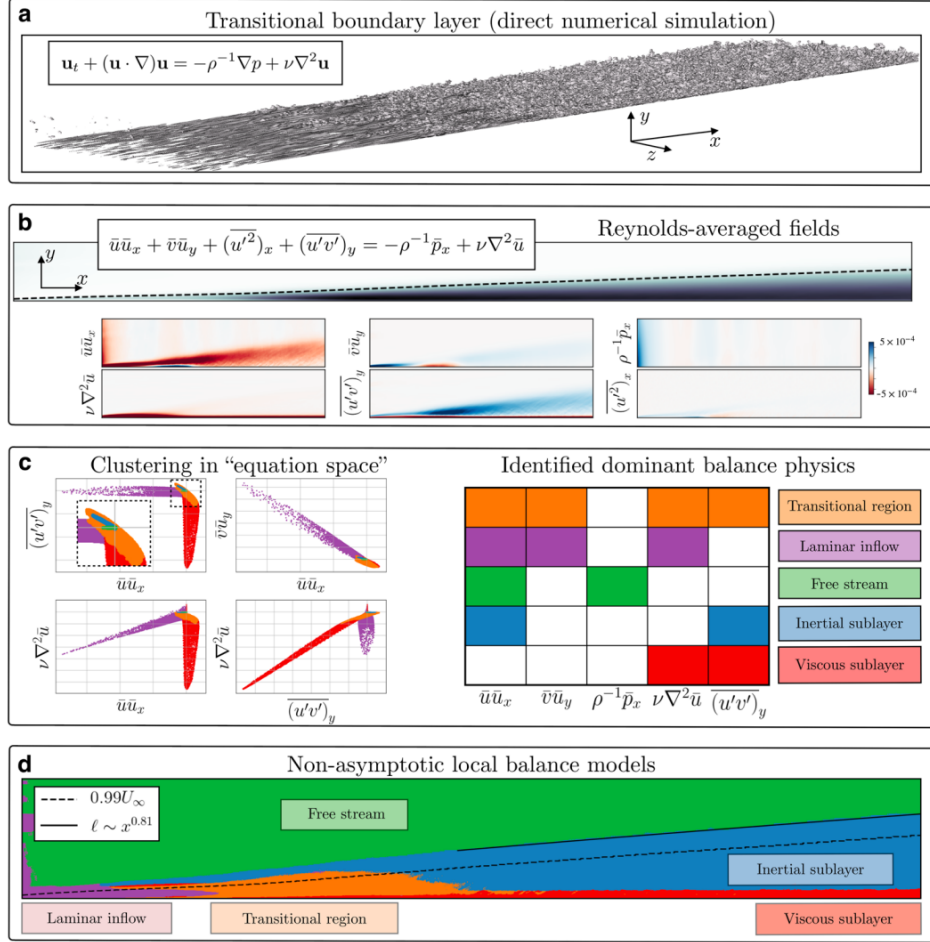


Figure 4.7: Data-driven balance model by Callaham *et al.* (Callaham et al., 2021a) applied to a boundary layer undergoing transition to turbulence. a) Instantaneous data from high-fidelity simulations (Lee and Zaki, 2018) and b) terms in the RANS equations obtained from the turbulence statistics. c) Covariance of the various terms grouped into clusters, labelled based on their physical meaning. d) Representation of the various clusters in the flow field, together with various boundary-layer quantities. Figure reproduced from Ref. (Callaham et al., 2021a) with permission of the publisher (Springer Nature).

## 4.7 Learning reduced-order models for vortex shedding behind an obstacle

In this section, we illustrate the possibility of learning ROMs in the case of flow around an obstacle, focusing on the wake. One possibility is to perform a modal decomposition, for instance, based on POD, and then carry out Galerkin projection of the governing Navier–Stokes equations onto the POD modes, as discussed in Sec. 4.6. This would lead to differential equations governing the temporal evolution of the POD coefficients associated with the spatial modes. This approach may exhibit two main problems in the case of turbulent flows, namely the possible numerical challenges of performing Galerkin projection and the need for many modes to reconstruct a significant fraction of the flow energy. As stated above, autoencoders can provide a compressed version of the original data by exploiting non-linearities, thus exhibiting the great potential to express high-dimensional

turbulence data in a few non-linear modes. As shown by Eivazi *et al.* [Eivazi et al., 2022], it is possible to learn a reduced representation of the original data where the latent vectors expressed in physical space exhibit orthogonality. This is achieved by promoting the learning of a latent space with disentangled latent vectors, which also enables learning parsimonious latent spaces. This is done by regularising the loss function, where the associated hyperparameter $\beta$ gives the name to the $\beta$-VAE framework discussed in Sec. 3.2. Larger values of $\beta$ give more weight to the term in the loss responsible for learning statistically-independent latent variables, therefore, when $\beta = 0$ one obtains the standard reconstruction loss function. In contrast, larger values of $\beta$ lead to higher orthogonality of the learned modes. At the same time, larger $\beta$ values will yield a worse reconstruction for the set number of latent vectors in the model. Based on this trade-off, it is possible to obtain a good balance between reconstruction and orthogonality. Eivazi *et al.* [Eivazi et al., 2022] illustrated this on the turbulent flow around two wall-mounted obstacles and showed that with only 5 AE modes, it is possible to reconstruct around 90% of the turbulent kinetic energy (TKE) with over 99% orthogonality of the modes. In comparison, 5 POD modes only reconstruct around 30% of the TKE. This is very interesting because the $\beta$-VAE, which, unlike other AE-based methods, produces orthogonal modes, yields a reduced representation that can be interpreted from a physical point of view. The first AE and POD modes are very similar, identifying the shear layers around the obstacles and the wake shedding. However, the AE modes exhibit a broader range of scales, incorporating additional higher-frequency turbulent fluctuations into the basic identified features (similar to those in the POD results). Consequently, there is great potential for this type of method to shed light on the physics of complex turbulent flows, in particular when using novel data-driven methods, such as transformers [Vaswani et al., 2017, Yousif et al., 2022], to predict the dynamics of the latent space.

Another linear approach discussed above to obtain low-dimensional representations of the flow is dynamic-mode decomposition (DMD), which is based on building a linear operator connecting the instantaneous snapshots in time. Unlike POD, the DMD modes are orthogonal in time, *i.e.* they are associated with a single frequency, which helps identify temporal features in fluid flows. HODMD enables establishing more complex relationships among snapshots, and although it requires additional hyper-parameter tuning, it can help to identify more detailed patterns in the flow. Martínez-Sánchez *et al.* [Martínez-Sánchez et al., 2023] used HODMD to study the turbulent flow in a simplified urban environment, emphasising the structures behind a wall-mounted obstacle. In this type of flow, a number of flow features emerge around the obstacle [Monnier et al., 2010], where a very important feature is the so-called arch vortex. This vortex, where the legs exhibit wall-normal and the spanwise roof vorticities, is responsible for the high concentration of pollutants in urban environments; therefore, understanding its formation mechanisms can have important implications for urban sustainability. In this context, HODMD enabled identifying two types of modes, namely vortex-generator and vortex-breaker features. The former is associated with low frequency, whereas the latter exhibits higher frequency, and both play important dynamic roles in flow physics. Another extension of the HODMD method also applied to the flow around a wall-mounted obstacle, was proposed by Amor *et al.* [Amor et al., 2023]. This study featured the so-called on-the-fly version of HODMD. The data is analysed dynamically as the simulation is run, without storing massive amounts of data for post-processing. Furthermore, more refined criteria for convergence of the modal decomposition were proposed, thus yielding a more effective way to analyse the data. Consequently, this on-the-fly approach reduces up to 80% in memory requirements compared with the traditional offline method. This is a big advantage when applied to large-scale numerical databases.

Causality maps, discussed in Sec. 2, have been used to study the dynamic interactions present in turbulent flows, focusing on the physical roles of various features. In particular, Lozano-Durán *et al.* [Lozano-Durán et al., 2020] studied the time series of the first Fourier modes in a turbulent
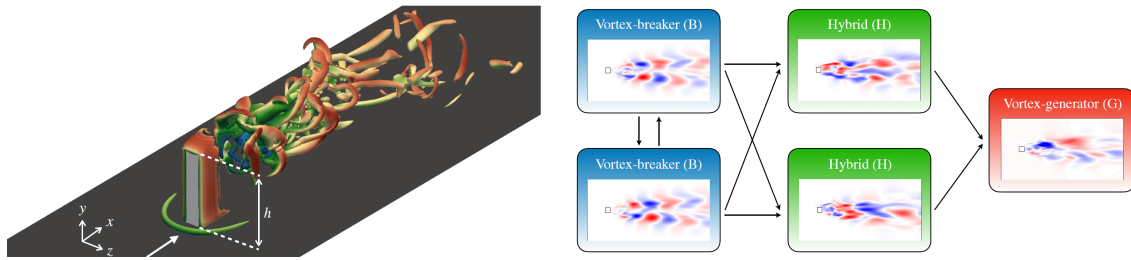
Figure 4.8: (Left) Instantaneous snapshot of the flow around a wall-mounted square cylinder, where the vortex clusters are identified with the $Q$ criterion (Hunt et al., 1998). The structures are coloured by their streamwise velocity, ranging from (dark blue) negative to (dark red) positive values. (Right) Schematic representation of the causal relations among modes, where two vortex-breaker (B), two hybrids (H), and one vortex-generator (G) modes are shown. Figure adapted from Ref. (Martínez-Sánchez et al., 2022).

channel. They found the following strong causal relations among modes: i) wall-normal modes causing streamwise modes, a phenomenon very closely connected with the well-known lift-up mechanism [Orr, 1907, Landahl and Landahlt, 1975] in near-wall turbulence; ii) wall-normal modes causing spanwise modes, which is associated with the roll generation, also connected with the lift-up process and the incompressibility of the flow; iii) streamwise modes causing spanwise ones, and spanwise modes causing wall-normal ones; both phenomena are connected with the mean-flow instability, including spanwise meandering and breakdown of the streaks [Swearingen and Blackwelder, 1987, Waleffe, 1995]. These causal relations were also identified [Martínez-Sánchez et al., 2022] in other simplified models of near-wall turbulence, such as the nine-equation model by Moehlis *et al.* [Moehlis et al., 2004], a fact that confirms the robustness of the causality framework utilised to study turbulence phenomena. Regarding the flow around a wall-mounted obstacle, the various modes discussed above and their connection with the arch vortex were assessed by Martínez-Sánchez *et al.* [Martínez-Sánchez et al., 2022] also using causality analysis. As can be observed in Figure 4.8 (left), the flow under consideration exhibits large-scale separation at the sharp edges of the obstacle and very prominent vortical structures in the wake. Figure 4.8 (right) exhibits the vortex-generator and breaker modes discussed above (associated with low and high frequencies, respectively), as well as an additional type of mode of intermediate frequency, denoted as hybrid mode. Clear causal relations are identified between the vortex-breaker and hybrid modes, closely connected with developing vortex-generator modes. This is of great interest because these causal relations define a sequence of events required for the production of the arch vortex (and the subsequent accumulation of pollutants in urban environments); thus, being able to control and inhibit this sequence of events may lead to novel sustainability solutions in cities (as well as to a deeper physical understanding of these complex turbulent flows).

## 4.8   Uncovering new physical understanding in wall-bounded turbulence

Turbulent flow is one of the most elusive areas of study within fluid mechanics. The wide range of spatial and temporal scales present in turbulence, and the highly non-linear behaviour that characterises it, significantly complicate the possibilities of gaining a deep physical understanding of the main mechanisms within turbulence; this becomes even more complicated in the case of wall-bounded turbulence, which is ubiquitous in science and engineering. Turbulence is characterised by coherent structures, three-dimensional regions that instantaneously satisfy certain physical properties. Note that this term sometimes refers to the features extracted by modal analysis. Still, we will consider the above definition in this work's context. A very important coherent structure in wall-bounded turbulence is the near-wall streak, extensively studied in the 1960s by Kline *et al.* [Kline et al., 1967]. As reported by Kim *et al.* [Kim et al., 1971], the near-wall production of turbulence is very closely connected with the dynamics of these streaks. Another important
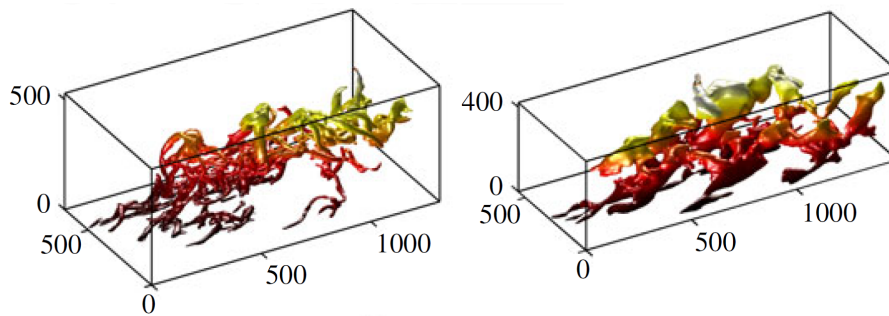
Figure 4.9: Coherent structures in a turbulent channel flow. We show (left) a vortex cluster and (right) an intense Reynolds-stress event. Figure adapted from Ref. (Lozano-Durán et al., 2012) with permission from the publisher (Cambridge University Press).

quantity in wall turbulence is the Reynolds shear stress, briefly introduced in Sec. 4.6. This quantity is essentially a correlation between stream-wise ($u'$) and wall-normal ($v'$) fluctuations and is responsible for the wall-normal momentum transport. Studying the coherent structures most relevant to the development of the Reynolds stresses is a critical goal for reaching a deeper understanding of turbulence. Several studies in the 1970s [Wallace et al., 1972, Lu and Willmarth, 1973] focused on the quadrant analysis to carry out this task; in this analysis, different near-wall events are classified in terms of the sign of their fluctuations, where the most dominant events are the so-called sweeps ($u' > 0$, $v' < 0$) and ejections ($u' < 0$, $v' > 0$). More recently, del Álamo *et al.* [del Álamo et al., 2006] have studied vortex clusters in turbulent channels, and Lozano-Durán *et al.* [Lozano-Durán et al., 2012] have analysed extreme Reynolds-stress events in the same flow case. The latter is defined as the three-dimensional connected regions satisfying:

$$|u'v'| > Hu_{\mathrm{rms}}v_{\mathrm{rms}}, \tag{4.6}$$

where the subscript 'rms' denotes root-mean-squared quantities, and $H$ is an empirical threshold denoted as a hyperbolic hole. In Figure 4.9, we show both types of coherent structure in a turbulent channel flow. Additional insight into the role of both types of structures can be obtained by tracking their evolution in time, such that the various structure interactions (advection, merges, splits, and dissipation) can be assessed [Lozano-Durán and Jiménez, 2014]. Convolutional neural networks (CNNs) have been used to predict the temporal evolution of the structures in turbulent channels [Schmekel et al., 2022], an approach that enables a deeper understanding of their dynamic behaviour. In turbulence, there is a direct cascade of energy from the larger, energy-containing structures towards the smaller dissipative ones; however, there is also an energy path in the opposite direction [Cardesa et al., 2017]. This picture, observed in homogeneous isotropic turbulence, becomes even more complicated in the case of wall-bounded turbulence [Jiménez, 2012]. Each wall-normal location has a different energy cascade because the wall segregates the flow by introducing wall-normal inhomogeneity. A comprehensive review of coherent structures in turbulence was provided by Jiménez [Jiménez, 2018b], who highlighted the potential and challenges of this perspective on turbulence. A very interesting open question raised in this work is the multi-scale organisation and interaction among the various individual structures and how they can dynamically produce the underlying physics of the flow.

Despite the extensive body of work on coherent structures in turbulence, there are still a number of open questions regarding the objective identification of the structures which play the most important role in the dynamics of turbulent flows. This fundamental question has implications for the theoretical knowledge of the physics of turbulence and the potential of flow-control strategies. If it is possible to identify these structures and they can be suppressed, there may be potential for novel and effective drag-reduction techniques. The vortex clusters and Reynolds-stress structures were

defined based on historical reasons and physical intuition. Although they play an important role in the flow, it is unclear whether these are the most relevant motions. A new type of structure that maximises the momentum transfer in a turbulent channel was identified by Jiménez [Jiménez, 2016], and he reported significant differences between these and the Reynolds stresses. An extension of this idea was implemented in two-dimensional decaying turbulence by removing subregions of the domain and assessing their relative influence in the future evolution of the flow [Jiménez, 2018a]. The idea is to quantify the "significance" of the various regions, and the result confirmed the initial physical intuition regarding this case: the most significant regions were vortices. The least significant ones exhibited high strain. In this direction, Cremades *et al.* [Cremades et al., 2023] proposed an approach to exploit the explainability of neural networks to assess the relevance of the coherent structures in turbulent flows. In this study, the SHapley Additive exPlanations (SHAP) framework [Winter, 2002, Lun-Chau et al., 2022] was used on the coherent structures identified in a turbulent channel; more concretely, the intense Reynolds stresses were first identified, and then a CNN was used to predict the location of those structures in the next time step [Schmekel et al., 2022]. The SHAP technique allows for identifying the impact of each of the features in the input (in this case, the three-dimensional Reynolds-stress events) on the prediction of the next step, thus enabling an assessment of their relevance to the future evolution of the flow. This framework could be used to find new ways of objectively identifying coherent regions in the flow. Another way to gain insight into the detailed mechanisms of turbulence via neural networks is to perform flow estimation, e.g. from the quantities measured at the wall to the turbulent fluctuations above [Guastoni et al., 2021, Güemes et al., 2021]. After training a neural network to make this prediction, detailed knowledge of the connection between the scales at the wall and the ones above can be gained through neural-network interpretability [Vinuesa and Sirmacek, 2021]. This approach allows us to discover a symbolic equation that can reproduce the predictive capabilities of the network. This can be achieved through the methodology developed by Cranmer *et al.* [Cranmer et al., 2020], which relies on symbolic regression (e.g. based on genetic programming) to obtain the equation relating input and output; see Sec. 4.6 for a related discussion. By analysing such an equation, it is possible to identify the characteristics of the scales relevant to this wall-normal interaction in wall-bounded turbulent flows.

## 4.9    Discovery of ocean mesoscale closures

The closure problem, described above in RANS equations, apply to many ocean and atmosphere modelling. In climate modelling, we must resolve (spatial) scales from meters to thousand kilometres. However, due to computational limitations, we need to truncate the spatial spectrum at a given scale - equivalent to the grid spacing of the numerical climate model. Therefore, all processes occurring below the spatial scales need to be approximated - this is the so-called parameterisation or closure problem for subgrid processes. The closure problem, described above in RANS equations, apply to many ocean and atmosphere modelling.

While RANS separates terms into time-averaged and fluctuating components, the most common approach is based on Large Eddy Simulation (LES), in which the filtering separates into a resolved scale and a sub-grid scale. The LES decomposition is based on the self-similarity of small-scale turbulent structures. The resolved scales are defined using a convolution integral with associated physical width, usually the grid cell size. Commonly used filters are box filters, normalised Gaussian, or a combination of both filters. Applying the filtering to the governing equations of the fluid (momentum and buoyancy) gives rise to a set of equations for the resolved scale, with a term - coined subgrid scale forcing - which depends on the fine scale. For the momentum equation, this

term subgrid term would be expressed as

$$\mathbf{S} = \begin{pmatrix} S_x \\ S_y \end{pmatrix} = (\overline{\mathbf{u}} \cdot \overline{\nabla})\overline{\mathbf{u}} - \overline{(\mathbf{u} \cdot \nabla)\mathbf{u}}, \tag{4.7}$$

where $\nabla$ is the horizontal 2D gradient operator, and the horizontal velocity $\mathbf{u} = (u, v)$, and the overline denotes the filtered (hence resolved) velocity on the grid.

Therefore $\mathbf{S}$ must be approximated with only resolved scales $\overline{\mathbf{u}}$ since the total variable $\mathbf{u}$ is not available to the model. Typically turbulence subgrid closures in a fluid are ad-hoc, such as Smagorinsky-type closures, in which the form of the closure is based on some physical argument that is assumed to hold across scale and regimes. This is rarely the case.

For ocean and atmosphere problems, the closure idea can also be boiled down to finding an expression for multiscale interaction that only depends on the resolved scale of the fluid. Similarly to traditional fluid problems, closures or parameterisations in the ocean and atmosphere modelling is often empirical and a source of error in simulations. Instead, equation discovery algorithms, as discussed in previous sections, can be used to uncover relationships between variables. For the closure problem, these algorithms can be applied to derive equations that describe the behaviour of the subgrid scales using resolved variables based on simulated data. The goal is to find the simplest (in some sense) mathematical relationship that accurately captures the behaviour of the subgrid-scale model, which can then be used for prediction. The main advantage of equation discovery algorithms is that they can uncover relationships that may not be immediately apparent, reducing the need for expert knowledge and human intuition in model building, as typically done for parameterisation in ocean and atmospheric modelling.

Zanna and Bolton [2020] used Relevance Vector Machine (RVM), a sparse Bayesian regression algorithm, to find closure models for momentum and buoyancy subgrid forcing. The RVM algorithm finds the most relevant input features - functions of the resolved scales - that will describe the subgrid-scale model. The RVM starts with many basis functions and iteratively removes irrelevant basis functions, arriving at a compact set of basis functions that best represent the data. Compared to other methods, the RVM algorithm has the advantage of handling noisy and redundant data and high-dimensional input spaces. Finally, it provides a probabilistic output, which can be a useful measure of uncertainty.

Below is a closure found by Zanna and Bolton [2020], using data from an ocean primitive equation model

$$\mathbf{S} \approx \kappa_{BT} \overline{\nabla} \cdot \begin{pmatrix} \zeta^2 - \zeta D & \zeta \tilde{D} \\ \zeta \tilde{D} & \zeta^2 + \zeta D \end{pmatrix}, \tag{4.8}$$

where $\zeta = \overline{v}_x - \overline{u}_y, D = \overline{u}_y + \overline{v}_x, \tilde{D} = \overline{u}_x - \overline{v}_y$, the short-hands $()_{x,y} \equiv \frac{\partial}{\partial x,y}$ are used for spatial derivatives, $\zeta$ is the relative vorticity, and $D$ and $\tilde{D}$ are the shearing and stretching deformation of the flow field, respectively. The authors were able to relate the found expression to energy transfer across scales, which mimics the impact of unresolved scales on large-scale energetics.

However, sparse linear regression entails trade-offs between the size and expressiveness of the feature library and the complexity and cost of sparse regression, as discussed in Zanna and Bolton [2020] and above. If we wish to include a deep library of functions, the number of different expressions needed will grow exponentially and might be limited by accurately taking derivatives of functions. Finally, many expressions might be highly correlated, preventing convergence [Hastie et al., 2015].

As discussed above, genetic programming (GP) [Koza, 1994] is an alternative approach. GP algorithms, unlike sparse regression, do not require a defined library of functions. [Ross et al., 2023] used GP with some modifications, including building spatial derivatives in spectral space and combining them with sparse regression to find robust expressions in turbulent datasets generated
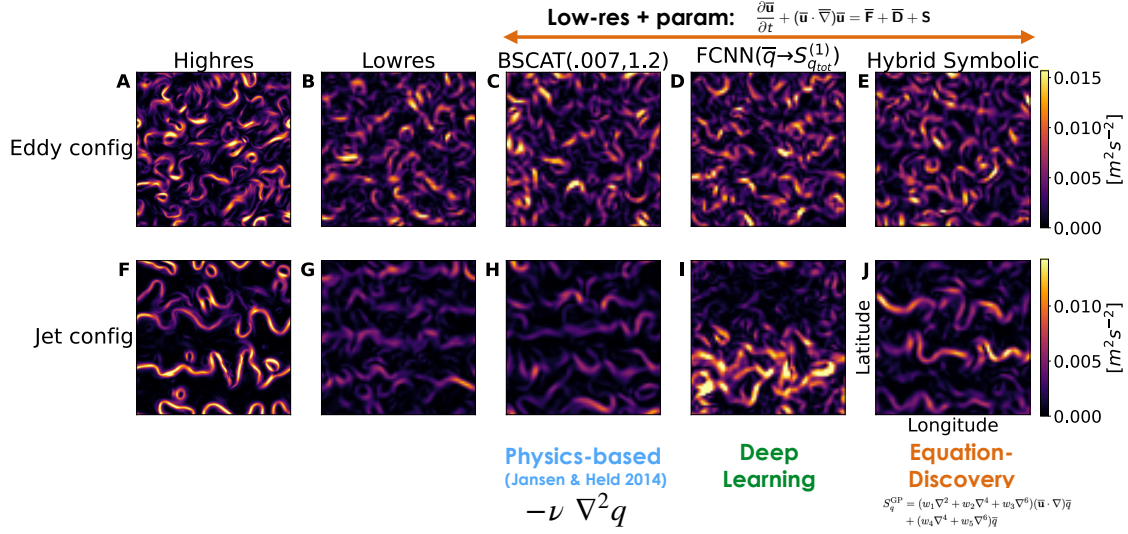
Figure 4.10: Snapshot of potential vorticity in two different simulations: Top = Eddy (mostly isotropic turbulence), Bottom = Jet (some elongated sharp features mixed with isotropic turbulence features). A, F: High-Resolution simulations; B, G: Coarse Resolution; C, H: Coarse Resolution with physics-based parameterisations; D, I: Coarse Resolution with a Neural Network-based Dramatisation; E, J: Coarse Resolution with parameterisation discovered with symbolic regression. The data-driven parameterisations are trained on eddy configuration, and only the equation-discovery lead to robust generalisation in different regimes without retraining. (Ross et al., 2023)

by idealised simulations. Focusing on results from [Ross et al., 2023], they look for the missing subgrid forcing for potential vorticity, $q$ - a variable that combines momentum and buoyancy effects in geophysical flows, and related to $\nabla \times \mathbf{u}$. In the first few iterations, the algorithm discovered quadratic expressions proportional to $(\overline{\mathbf{u}} \cdot \nabla)\overline{q}$, similarly to previous theoretical studies [Meneveau and Katz, 2000, Anstey and Zanna, 2017]. Often these expressions cannot be used as standalone parameterisations implemented in coarse-resolution models due to numerical stability constraints. The next few iterations of the GP-sparse regression algorithm led to eddy-viscosity models that dissipate energy at small scales, $\nabla^4\overline{q}$ and redistribute energy to larger scales, i.e. kinetic energy backscatter $\nabla^6\overline{q}$ [Jansen and Held, 2014]. Additional terms, which are cubic in model variables and contain a double-advection operator, $(\overline{\mathbf{u}} \cdot \nabla)^2$, can ensure dissipation of enstrophy [Marshall and Adcroft, 2010], helping with model stability. In addition, there were additional terms that we were not discovered previously. In summary, our discovered closure contains elements of existing subgrid parameterisations, which have pros and cons when used as standalone ones but, when combined, could capture all necessary properties for stable implementation and accurate representation of momentum, energy and enstrophy fluxes missing at coarse resolution.

To test our discovered closure, we implement it in a coarse-resolution simulation (see Fig. 4.10). The goal is to improve the physics of the coarse-resolution model (panel B) relative to the high-resolution model simulation (panel A). To this end, we run the coarse-resolution simulation with a physics-based (empirical parameterisation; panel C), with data-driven parameterisation learned using a convolutional neural network (panel D), and with the equation-discovery parameterisation (panel E). All simulations are improving the flow, and some aspects of the statistics are also improved. However, generalisation is vastly different without retraining the data-driven driven parameterisations or tuning the physics-based parameterisations. We test our parameterisations in the same model in which we changed the rotation rate in order to form jets and less isotropic turbulence (panels F for high resolution and G for low resolution without any parameterisation). The physics-based parameterisation has little impact on the flow (panel H), but the implementation

of the deep learning parameterisation has a very detrimental effect on the flow, most trying to make the flow more isotropic (panel I). On the other hand, the implementation of the equation discovery-based parameterisation substantially improves the flow (panel J)- reinforcing the need to discover relationships from data that encapsulate the necessary laws of physics to mimic the scale interactions which are internal to the fluid and not dependent on the configuration of the simulations.

This symbolic parameterisation includes up to the seventh spatial derivative of $\overline{q}$, which may be unrealistic to implement into a climate model. However, it might be more realistic than a fully non-local approach, such as the convolutional neural network parameterisations considered in other studies or extremely local physics-based parameterisations (such as anti-viscosity). Most importantly, the sparse model can generalise well without retraining, while the neural network-based parameterisations perform poorly.

# 5. Concluding remarks

The fields of causal and equation discovery have emerged in recent years as important research areas that apply artificial intelligence and machine learning to analyse complex systems [Peters et al., 2017b, Schmidt et al., 2011, Runge et al., 2023]. The fields respectively aim to identify causal relationships and discover equations that can be used to predict the behaviour of the system, including the effects of interventions.

In this paper, we have reviewed the state-of-the-art in both fields and discussed their respective approaches and techniques. Causal discovery aims to discover the qualitative cause-and-effect relationships between the variables in a system. In order to achieve this task using non-experimental data, causal discovery employs certain enabling assumptions (see Sec. 2). Among these assumptions is that the data-generating process can be described by a structural causal model and the corresponding causal graph. Methods for causal discovery are manifold and can be partitioned into constraint-based, score-based, asymmetry-based, and context-based methods. Data from the physical world typically comes in the form of time series with autocorrelation and potentially non-stationary behaviour. Autocorrelation and potential non-stationarity pose statistical challenges for many causal discovery methods as they are typically designed for i.i.d. data. In addition, the true functional relationship between variables can be highly non-linear, and the variables can be high-dimensional, both increasing model complexity and affecting the efficiency of causal discovery methods. All these challenges are compounded by the fact that the data acquired from real-world processes are often far from ideal, with problems such as missing data and inherent selection bias that might lead to the observed data not being representative of the process underlying it. These and many other challenges are avenues for future research in causal discovery and many of its sister fields, such as Bayesian networks and conditional independence testing, to name a few.

In equation discovery, the focus is on understanding the structure of a system by discovering equations, state variables and laws that can be used to predict (and, more importantly, to understand) its behaviour (see Sec. 3). The main techniques used in this field are symbolic regression, evolutionary algorithms, and deep learning. These methods offer the potential to discover both linear and nonlinear equations but suffer from the need for large datasets and the difficulty of finding accurate equations in complex systems. More relevant challenges have to do with identifiability issues and

the impossibility of evaluating the generality of the equations or even the criteria to select the most general ones. Broader (and perhaps more philosophical) questions need to be addressed, such as compressibility or sparsity, confronted with expressive power, the role of physical units and modularity, to name a few.

Thus, causality studies and equation inference approaches have synergistic goals. Both fields have made significant advances in recent years and offer considerable promise for further research. In particular, techniques from both fields can be combined to create hybrid models capable of uncovering causal relationships and equations. Additionally, developing more efficient algorithms and better methods for dealing with the challenge of overfitting could lead to further progress in both fields. More specifically, the question remains as to which current approaches provide stronger guarantees for the uniqueness/equifinality of the discovered equation or inferred causal graph. This key aspect in the inferential discovery of physical models should receive more attention in future research enterprises.

A wide range of case studies in many areas of interest in the physical sciences (neurosciences, Earth and climate sciences, fluid mechanics) has illustrated the performance of causal discovery and equation discovery algorithms (see Sec. 4). We noted that specific methods and techniques reside in particular fields and do not permeate to others, mainly because of the needed assumptions and data characteristics. Yet, as has been the case for centuries, there is a lack of transdisciplinary in science. Directly stemming from this review, it is evident that analysing complex systems require an inter/trans-disciplinary approach that combines method and domain expertise. Techniques from artificial intelligence, machine learning, and control theory can be combined to better understand a system's behaviour and make accurate predictions. As such, future research in causal and equation discovery should consider the potential benefits of a more integrative and fused approach to analysing complex systems. This is perhaps especially important for recent developments designed for answering critical questions in specific areas but which, given their fundamental nature, have a wider appeal. For instance, the progress in the empirical inference of transfer operators in fluid mechanics or chemical reaction pathways has unexplored implications in understanding the metastable dynamics of neuronal network responses.

Overall, the fields of causal discovery and equation discovery are rapidly advancing, and there is a growing synergy between them. Despite the remaining challenges, researchers have made great strides in uncovering the underlying structure of complex systems. With continued research and development, we can look forward to further advances in both fields and unlocking complex systems' mysteries.

# Acknowledgements & Contributions

# Bibliography

G. Abbati, P. Wenk, M. A. Osborne, A. Krause, B. Schölkopf, and S. Bauer. AReS and MaRS adversarial and MMD-minimizing regression for SDEs. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1–10. PMLR, 09–15 Jun 2019.

A. Abhyankar. Linear and nonlinear Granger causality: Evidence from the UK stock index futures market. *The Journal of Futures Markets (1986-1998)*, 18(5):519, 1998.

J. E. Adsuara, A. Pérez-Suay, Á. Moreno-Martínez, G. Camps-Valls, G. Kraemer, M. Reichstein, and M. Mahecha. Discovering differential equations from Earth observation data. In *IGARSS 2020 - 2020 IEEE International Geoscience and Remote Sensing Symposium*, pages 3999–4002, 2020.

S. Amari. Dynamics of pattern formation in lateral-inhibition type neural fields. *Biol. Cybern.*, 27 (2):77–87, 8 1977.

C. Amor, P. Schlatter, R. Vinuesa, and S. Le Clainche. Higher-order dynamic mode decomposition on-the-fly: A low-order algorithm for complex fluid flows. *J. Comput. Phys.*, 475:111849, 2023.

N. Ancona, D. Marinazzo, and S. Stramaglia. Radial basis function approach to nonlinear Granger causality of time series. *Physical Review E*, 70(5):056221, 2004.

B. Andersen and T. Fagerhaug. *Root cause analysis: simplified tools and techniques*. Quality Press, 2006.

M. Ankerst, M. M. Breunig, H.-P. Kriegel, and J. Sander. OPTICS: Ordering points to identify the clustering structure. *ACM Sigmod record*, 28(2):49–60, 1999.

J. A. Anstey and L. Zanna. A deformation-based parametrization of ocean mesoscale eddy Reynolds stresses. *Ocean Modelling*, 112:99–111, 2017.

G. Antonelli, S. Chiaverini, and P. Di Lillo. On data-driven identification: Is automatically discovering equations of motion from data a Chimera? *Nonlinear Dynamics*, pages 1–12, 2022a.

G. Antonelli, S. Chiaverini, and P. Di Lillo. On data-driven identification: Is automatically discovering equations of motion from data a Chimera? *Nonlinear Dyn*, 2022b.

M. C. Aoi, V. Mante, and J. W. Pillow. Prefrontal cortex exhibits multidimensional dynamic encoding during decision-making. *Nat. Neurosci.*, 23(11):1410–1420, Nov. 2020.

J. Arenas-García, K. Petersen, G. Camps-Valls, and L. Hansen. Kernel multivariate analysis framework for supervised subspace learning: A tutorial on linear and kernel multivariate methods. *IEEE Signal Processing Magazine*, 30(4):16–29, 2013.

C. K. Assaad, E. Devijver, and E. Gaussier. Discovery of extended summary graphs in time series. In J. Cussens and K. Zhang, editors, *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pages 96–106. Pmlr, 01–05 Aug 2022a.

C. K. Assaad, E. Devijver, and E. Gaussier. Survey and evaluation of causal discovery methods for time series. *Journal of Artificial Intelligence Research*, 73:767–819, 2022b.

W. J. Baars and C. Tinney. Proper orthogonal decomposition-based spectral higher-order stochastic estimation. *Phys. Fluids*, 26:055112, 2014.

P. J. Baddoo, B. Herrmann, B. J. McKeon, J. Nathan Kutz, and S. L. Brunton. Physics-informed dynamic mode decomposition. *Proceedings of the Royal Society A*, 479(2271):20220576, 2023.

Z. Bai, E. Kaiser, J. Proctor, J. Kutz, and S. Brunton. Dynamic mode decomposition for compressive system identification. *Aiaa J.*, 58:561, 2020.

E. Balaguer-Ballester, C. C. Lapish, J. K. Seamans, and D. Durstewitz. Attracting dynamics of frontal cortex ensembles during memory-guided decision-making. *PLoS Computational Biology*, 7(5):e1002057, 2011.

E. Balaguer-Ballester, A. Tabas-Diaz, and M. Budka. Can we identify non-stationary dynamics of trial-to-trial variability? *PLoS ONE*, 9(4):1–13, 04 2014.

E. Balaguer-Ballester, R. Moreno-Bote, G. Deco, and D. Durstewitz. Editorial: Metastable dynamics of neural ensembles. *Frontiers in Systems Neuroscience*, 11, 2018.

E. Balaguer-Ballester, R. Nogueira, J. M. Abofalia, R. Moreno-Bote, and M. V. Sanchez-Vives. Representation of foreseeable choice outcomes in orbitofrontal cortex triplet-wise interactions. *PLOS Computational Biology*, 16(6):1–30, 06 2020.

D. Baldocchi. Measuring fluxes of trace gases and energy between ecosystems and the atmosphere – the state and future of the eddy covariance method. *Global Change Biology*, 20(12):3600–3609, 2014.

D. L. Barack, E. K. Miller, C. I. Moore, A. M. Packer, L. Pessoa, L. N. Ross, and N. C. Rust. A call for more clarity around causality in neuroscience. *Trends in neurosciences*, 2022.

L. Barnett, A. B. Barrett, and A. K. Seth. Granger causality and transfer entropy are equivalent for Gaussian variables. *Physical review letters*, 103(23):238701, 2009.

L. Barnett, A. B. Barrett, and A. K. Seth. Misunderstandings regarding the application of Granger causality in neuroscience. *Proceedings of the National Academy of Sciences*, 115(29):E6676–e6677, 2018.

A. B. Barrett, L. Barnett, and A. K. Seth. Multivariate Granger causality and generalized variance. *Phys. Rev. E*, 81(4):41907, 2010.

D. G. Barrett, A. S. Morcos, and J. H. Macke. Analyzing biological and artificial neural networks: Challenges with opportunities for synergy? *Current Opinion in Neurobiology*, 55:55–64, 2019. Machine Learning, Big Data, and Neuroscience.

D. S. Bassett and O. Sporns. Network neuroscience. *Nature neuroscience*, 20(3):353–364, 2017.

D. Bell, J. Kay, and J. Malley. A non-parametric approach to non-linear causality testing. *Economics Letters*, 51(1):7–18, 1996.

K. Bello, B. Aragam, and P. Ravikumar. DAGMA: Learning DAGs via M-matrices and a log-determinant acyclicity characterization. In *Advances in Neural Information Processing Systems*, 2022.

A. Belloni, V. Chernozhukov, and C. Hansen. Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, 81(2):608–650, 2014.

P. Benner, S. Gugercin, and K. Willcox. A survey of projection-based model reduction methods for parametric dynamical systems. *SIAM Rev.*, 57(4):483–531, 2015.

R. Berk, L. Brown, A. Buja, K. Zhang, and L. Zhao. Valid post-selection inference. *The Annals of Statistics*, 41(2):802 – 837, 2013.

D. Bernal-Casas, E. Balaguer-Ballester, M. F. Gerchen, S. Iglesias, H. Walter, A. Heinz, A. Meyer-Lindenberg, K. E. Stephan, and P. Kirsch. Multi-site reproducibility of prefrontal-hippocampal connectivity estimates by stochastic DCM. *Neuroimage*, 82:555–563, 2013.

W. Bialek, I. Nemenman, and N. Tishby. Predictability, complexity, and learning. *Neural computation*, 13(11):2409–2463, 2001.

L. Biggio, T. Bendinelli, A. Neitz, A. Lucchi, and G. Parascandolo. Neural symbolic regression that scales. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 936–945. Pmlr, 18–24 Jul 2021.

C. Bishop. Pattern recognition and machine learning. *Springer New York*, 2006.

S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, and D.-U. Hwang. Complex networks: Structure and dynamics. *Physics reports*, 424(4-5):175–308, 2006.

M. A. Boden. Creativity and artificial intelligence: A contradiction in terms. *The philosophy of creativity: New essays*, pages 224–46, 2014.

J. Boelts, J.-M. Lueckmann, R. Gao, and J. H. Macke. Flexible and efficient simulation-based inference for models of decision-making. *eLife*, 11:e77220, July 2022.

G. Boffetta, M. Cencini, M. Falcioni, and A. Vulpiani. Predictability: A way to characterize complexity. *Physics reports*, 356(6):367–474, 2002.

K. A. Bollen. *Structural Equations with Latent Variables*. John Wiley & Sons, New York, NY, USA, 1989.

J. Bongard and H. Lipson. Automated reverse engineering of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences of the United States of America*, 104(24): 9943–9948, 2007.

S. Bongers and J. M. Mooij. From random differential equations to structural causal models: The stochastic case. *arXiv preprint arXiv:1803.08784*, 2018.

S. Bongers, T. Blom, and J. M. Mooij. Causal modeling of dynamical systems, 2018.

S. Bongers, P. Forré, J. Peters, and J. M. Mooij. Foundations of structural causal models with cycles and latent variables. *The Annals of Statistics*, 49(5):2885 – 2915, 2021.

T. Bouezmarni, J. V. Rombouts, and A. Taamouti. Nonparametric copula-based test for conditional independence with applications to Granger causality. *Journal of Business & Economic Statistics*, 30(2):275–287, 2012.

J. V. Boussinesq. *Théorie analytique de la chaleur: mise en harmonie avec la thermodynamique et avec la théorie mécanique de la lumière T. 2, Refroidissement et échauffement par rayonnement conductibilité des tiges, lames et masses cristallines courants de convection théorie mécanique de la lumière*. Gauthier-Villars, 1923.

P. Bradley, S. Gold, and S. Silverman. Constructive induction from incomplete data: A comparative study. *Machine Learning*, 42(1):7–48, 2001.

C. Brennan, A. Aggarwal, R. Pei, D. Sussillo, and A. Proekt. One dimensional approximations of neuronal dynamics reveal computational strategy. *PLOS Computational Biology*, 19(1):1–27, 01 2023. doi: 10.1371/journal.pcbi.1010784. URL https://doi.org/10.1371/journal.pcbi.1010784.

M. Brenner, F. Hess, J. M. Mikhaeil, L. F. Bereska, Z. Monfared, P.-C. Kuo, and D. Durstewitz. Tractable dendritic RNNs for reconstructing nonlinear dynamical systems. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 2292–2320. Pmlr, 17–23 Jul 2022.

S. L. Bressler and A. K. Seth. Wiener-Granger causality: A well established methodology. *Neuroimage*, 58(2):323–329, 2011.

P. Brouillard, S. Lachapelle, A. Lacoste, S. Lacoste-Julien, and A. Drouin. Differentiable causal discovery from interventional data. *Advances in Neural Information Processing Systems*, 33: 21865–21877, 2020.

G. Brown, A. Pocock, M.-J. Zhao, and M. Luján. Conditional likelihood maximisation: A unifying framework for information theoretic feature selection. *The journal of machine learning research*, 13:27–66, 2012.

T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors,

*Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.

N. Brunel. Dynamics of sparsely connected networks of excitatory and inhibitory spiking neurons. *J. Comput. Neurosci.*, 8(3):183–208, May 2000.

S. Brunton, B. Brunton, J. Proctor, E. Kaiser, and J. Kutz. Chaos as an intermittently forced linear system. *Nat. Commun.*, 8:19, 2017.

S. L. Brunton and J. N. Kutz. *Data-driven science and engineering: Machine learning, dynamical systems, and control.* Cambridge University Press, 2022.

S. L. Brunton, B. W. Brunton, J. L. Proctor, and J. N. Kutz. Koopman invariant subspaces and finite linear representations of nonlinear dynamical systems for control. *PLoS ONE*, 11, 2016a.

S. L. Brunton, J. L. Proctor, and J. N. Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the national academy of sciences*, 113(15):3932–3937, 2016b.

D. Bueso, M. Piles, and G. Camps-Valls. Nonlinear PCA for spatio-temporal analysis of Earth observation data. *IEEE Transactions on Geoscience and Remote Sensing*, 58(8):5752–5763, Aug. 2020.

E. Bullmore and O. Sporns. Complex brain networks: Graph theoretical analysis of structural and functional systems. *Nature reviews neuroscience*, 10(3):186–198, 2009.

B. Burlacu, G. Kronberger, and M. Kommenda. Operon C++: An Efficient Genetic Programming Framework for Symbolic Regression. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference Companion*, Gecco '20, page 1562–1570, New York, NY, USA, 2020. Association for Computing Machinery.

H. Butterfield. *The origins of modern science*, volume 90507. Simon and Schuster, 1965.

A. Byrne, R. D. O'Dea, M. Forrester, J. Ross, and S. Coombes. Next-generation neural mass and field modeling. *Journal of Neurophysiology*, 123(2):726–742, 2020. Pmid: 31774370.

A. Byrne, J. Ross, R. Nicks, and S. Coombes. Mean-field models for EEG/MEG: From oscillations to waves. *Brain Topography*, 35, 05 2021.

J. Cabral, F. Castaldo, J. Vohryzek, V. Litvak, C. Bick, R. Lambiotte, K. Friston, M. L. Kringelbach, and G. Deco. Metastable oscillatory modes emerge from synchronization in the brain spacetime connectome. *Communications Physics*, 5(1):184, 2022.

J. L. Callaham, J. V. Koch, B. W. Brunton, J. N. Kutz, and S. L. Brunton. Learning dominant physical processes with data-driven balance models. *Nat. Commun.*, 12(1):1–10, 2021a.

J. L. Callaham, G. Rigas, J.-C. Loiseau, and S. L. Brunton. An empirical mean-field model of symmetry-breaking in a turbulent wake. *Sci. Adv.*, 8:eabm4786, 2021b.

J. L. Callaham, S. L. Brunton, and J.-C. Loiseau. On the role of nonlinear correlations in reduced-order modeling. *J. Fluid Mech.*, 938:A1, 2022.

G. Camps-Valls, J. Verrelst, J. Munoz-Mari, V. Laparra, F. Mateo-Jimenez, and J. Gomez-Dans. A survey on Gaussian processes for Earth-observation data analysis: A comprehensive investigation. *IEEE Geoscience and Remote Sensing Magazine*, 4(2):58–78, 2016. doi: 10.1109/MGRS.2015.2510084.

J. I. Cardesa, A. Vela-Martín, and J. Jiménez. The turbulent cascade in five dimensions. *Science*, 357:782–784, 2017.

K. Carlberg, M. Barone, and H. Antil. Galerkin v. least-squares Petrov–Galerkin projection in nonlinear model reduction. *J. Comput. Phys.*, 330:693–734, 2017.

R. Castelo and A. Siebes. Priors on network structures. Biasing the search for Bayesian networks. *International Journal of Approximate Reasoning*, 24(1):39–57, 2000.

M. Cenedese, J. Axås, B. Bäuerlein, K. Avila, and G. Haller. Data-driven modeling and prediction of nonlinearizable dynamics via spectral submanifolds. *Nat. Commun.*, 13:872, 2022.

K. Champion, B. Lusch, J. N. Kutz, and S. L. Brunton. Data-driven discovery of coordinates and governing equations. *Proceedings of the National Academy of Sciences*, 116(45):22445–22451, 2019.

V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3):1–58, 2009.

C. Chatfield. *The analysis of time series: Theory and practice*. Springer, 2013.

B. Chen, K. Huang, S. Raghupathi, I. Chandratreya, Q. Du, and H. Lipson. Automated discovery of fundamental variables hidden in experimental data. *Nature Computational Science*, 2(7): 433–442, 2022.

X. Chen, F. Ginoux, M. Carbo-Tano, T. Mora, A. M. Walczak, and C. Wyart. Granger causality analysis for calcium transients in neuronal networks, challenges and improvements. *eLife*, 12: e81279, Feb. 2023.

Y. Chen, G. Rangarajan, J. Feng, and M. Ding. Analyzing multiple nonlinear time series with extended Granger causality. *Physics letters A*, 324(1):26–35, 2004.

D. Chicharro. On the spectral formulation of Granger causality. *Biological cybernetics*, 105(5): 331–347, 2011.

D. Chicharro and S. Panzeri. Algorithms of causal inference for the analysis of effective connectivity among brain regions. *Frontiers in Neuroinformatics*, 8, 2014.

D. M. Chickering. Learning equivalence classes of Bayesian-network structures. *Journal of Machine Learning Research*, 2:445–498, Mar. 2002a.

D. M. Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002b.

D. M. Chickering and C. Meek. Selective greedy equivalence search: Finding optimal Bayesian networks using a polynomial number of score evaluations. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, pages 211–219, 2015.

M. Chickering. Statistically efficient greedy equivalence search. In J. Peters and D. Sontag, editors, *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 124 of *Proceedings of Machine Learning Research*, pages 241–249. Pmlr, 03–06 Aug 2020.

R. Christiansen, M. Baumann, T. Kuemmerle, M. D. Mahecha, and J. Peters. Toward causal inference for spatio-temporal data: Conflict and forest loss in Colombia. *Journal of the American Statistical Association*, 117(538):591–601, 2022.

T. Claassen and I. G. Bucur. Greedy equivalence search in the presence of latent confounders. In J. Cussens and K. Zhang, editors, *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pages 443–452. Pmlr, 01–05 Aug 2022.

R. R. Coifman and S. Lafon. Diffusion maps. *Applied and Computational Harmonic Analysis*, 21 (1):5–30, 2006. Special Issue: Diffusion Maps and Wavelets.

D. Colombo, M. H. Maathuis, et al. Order-independent constraint-based causal structure learning. *J. Mach. Learn. Res.*, 15(1):3741–3782, 2014.

G. F. Cooper and C. Yoo. Causal discovery from a mixture of experimental and observational data. In *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence*, UAI'99, page 116–125, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.

G. K. Cooray, B. Sengupta, P. Douglas, and K. Friston. Dynamic causal modelling of electrographic seizure activity using Bayesian belief updating. *Neuroimage*, July 2015.

N. Copernicus, M.-P. Lerner, A. P. Segonds, J.-P. Verdet, C. Luna, D. Savoie, and M. Toulmonde. *De revolutionibus orbium coelestium*, volume 1. Johnson Reprint Corporation, 1965.

C. Cornelio, S. Dash, V. Austel, T. R. Josephson, J. Goncalves, K. L. Clarkson, N. Megiddo, B. E. Khadir, and L. Horesh. Combining data and theory for derivable scientific discovery with AI-Descartes. *Nature Communications*, 14:1777, 2023.

J. Correa, S. Lee, and E. Bareinboim. Nested counterfactual identification from arbitrary surrogate experiments. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, volume 34, pages 6856–6867, 2021.

N. L. Cramer. A representation for the adaptive generation of simple sequential programs. In J. J. Grefenstette, editor, *Proceedings of an International Conference on Genetic Algorithms and the Applications*, pages 183–187, Carnegie-Mellon University, Pittsburgh, PA, USA, 1985.

M. Cranmer. PySR: Fast & Parallelized Symbolic Regression in Python/Julia, Sept. 2020.

M. Cranmer, A. Sanchez-Gonzalez, P. Battaglia, R. Xu, K. Cranmer, D. Spergel, and S. Ho. Discovering symbolic models from deep learning with inductive biases. In *Advances in Neural Information Processing Systems (NeurIPS 2020)*, 2020.

A. Cremades, S. Hoyas, P. Quintero, M. Lellep, M. Linkmann, and R. Vinuesa. Explaining wall-bounded turbulence through deep learning. *Preprint arXiv:2302.01250*, 2023.

J. Crutchfield and B. McNamara. Equations of motion from a data series. *Complex Syst.*, 1:417–452, 1987.

J. Cruzat, R. Herzog, P. Prado, Y. Sanz-Perl, R. Gonzalez-Gomez, S. Moguilner, M. L. Kringelbach, G. Deco, E. Tagliazucchi, and A. Ibañez. Temporal irreversibility of large-scale brain dynamics in Alzheimer's disease. *Journal of Neuroscience*, 43(9):1643–1656, 2023.

J. P. Cunningham and B. M. Yu. Dimensionality reduction for large-scale neural recordings. *Nat. Neurosci.*, 17(11):1500–1509, Nov. 2014.

R. Dahlhaus and M. Eichler. Causality and graphical models in time series analysis. *Oxford Stat. Sci. Ser*, 27, 01 2003.

B. C. Daniels and I. Nemenman. Automated adaptive inference of phenomenological dynamical models. *Nature communications*, 6:8133, 2015.

P. Daniušis, D. Janzing, J. Mooij, J. Zscheischler, B. Steudel, K. Zhang, and B. Schölkopf. Inferring deterministic causal relations. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, Uai'10, page 143–150, Arlington, Virginia, USA, 2010. AUAI Press. ISBN 9780974903965.

C. Darwin. *On the Origin of Species*. John Murray, London, 1859.

DataRobot Inc. Eureqa as part of DataRobot's service. https://www.datarobot.com/nutonian/, 2023.

T. Davis, K. F. LaRocque, J. A. Mumford, K. A. Norman, A. D. Wagner, and R. A. Poldrack. What do differences between multi-voxel and univariate analysis mean? How subject-, voxel-, and trial-level variance impact fMRI analysis. *NeuroImage*, 97:271–283, 2014.

M. Deistler, J. H. Macke, and P. J. Gonçalves. Energy-efficient network activity from disparate circuit parameters. *Proceedings of the National Academy of Sciences*, 119(44):e2207632119, 2022.

J. C. del Álamo, J. Jiménez, P. Zandonade, and R. D. Moser. Self-similar vortex clusters in the turbulent logarithmic region. *J. Fluid Mech.*, 561:329–358, 2006.

N. Deng, B. R. Noack, M. Morzynski, and L. R. Pastur. Low-order model for successive bifurcations of the fluidic pinball. *J. Fluid Mech.*, 884:A37, 2020.

D. C. Dennett. Real patterns. *The journal of Philosophy*, 88(1):27–51, 1991.

G. Di Capua, M. Kretschmer, J. Runge, A. Alessandri, R. Donner, B. van Den Hurk, R. Vellore, R. Krishnan, and D. Coumou. Long-lead statistical forecasts of the indian summer monsoon rainfall based on causal precursors. *Weather and Forecasting*, 34(5):1377–1394, 2019.

A. Diaz, J. Johnson, G. Varando, and G. Camps-Valls. Learning latent functions for causal discovery. *Submitted*, 2023.

E. Diaz, J. Adsuara, A. Moreno-Martinez, M. Piles, and G. Camps-Valls. Inferring causal relations from observational long-term carbon and water fluxes records. *Scientific Reports*, 12:1610, 2022.

A. B. Dickerson. *Kant on representation and objectivity*. Cambridge University Press, 2003.

V. Didelez. Asymmetric separation for local independence graphs. In *23rd Annual Conference on Uncertainty in Artifical Intelligence*, 2006.

V. Didelez. Graphical models for marked point processes based on local independence. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 70(1):245–264, 2008.

G. C.-V. Diego Bueso, Maria Piles. Explicit Granger causality in kernel Hilbert spaces. *Physical Review E*, 102:062201, 2020.

C. Diks and V. Panchenko. A new statistic and practical guidelines for nonparametric Granger causality testing. *Journal of Economic Dynamics and Control*, 30(9-10):1647–1669, 2006.

M. Ding, Y. Chen, and S. L. Bressler. Granger causality: Basic theory and application to neuroscience. *Handbook of time series analysis: recent theoretical developments and applications*, pages 437–460, 2006.

J. Donges, Y. Zou, N. Marwan, and J. Kurths. The backbone of the climate network. *Epl*, 87:48007, 2009a.

J. F. Donges, Y. Zou, N. Marwan, and J. Kurths. Complex networks in climate dynamics. *The European Physical Journal Special Topics*, 174(1):157–179, 2009b.

B. E. Dowd. Separated at birth: Statisticians, social scientists, and causality in health services research. *Health Services Research*, 46(2):397–420, 2011.

C. J. Ducasse. Whewell's philosophy of scientific discovery. II. *The Philosophical Review*, 60(2): 213–234, 1951.

L. Duncker and M. Sahani. Temporal alignment and latent Gaussian process factor inference in population spike trains. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

L. Duncker and M. Sahani. Dynamics on the manifold: Identifying computational dynamical activity from neural population recordings. *Current Opinion in Neurobiology*, 70:163–170, 2021. Computational Neuroscience.

D. Durstewitz. A state space approach for piecewise-linear recurrent neural networks for identifying computational dynamics from neural measurements. *PLOS Computational Biology*, 13(6):1–33, 06 2017.

D. Durstewitz and J. K. Seamans. The dual-state theory of prefrontal cortex dopamine function with relevance to catechol-o-methyltransferase genotypes and schizophrenia. *Biological Psychiatry*, 64(9):739–749, 2008. Neurodevelopment and the Transition from Schizophrenia Prodrome to Schizophrenia.

S. Dzeroski and L. Todorovski. *Inductive Logic Programming: Techniques and Applications*. Springer Science+Business Media, 2007.

S. Džeroski and L. Todorovski. Reliable induction of recursive production rules. *Machine Learning*, 20(3):229–256, 1995.

I. Ebert-Uphoff and Y. Deng. Causal discovery in the geosciences—Using synthetic data to learn how to interpret results. *Computers & Geosciences*, 99:50–60, 2017.

H. Eivazi, L. Guastoni, P. Schlatter, H. Azizpour, and R. Vinuesa. Recurrent neural networks and Koopman-based frameworks for temporal predictions in a low-order model of turbulence. *Int. J. Heat Fluid Flow*, 90:108816, 2021.

H. Eivazi, S. Le Clainche, S. Hoyas, and R. Vinuesa. Towards extraction of orthogonal and parsimonious non-linear modes from turbulent flows. *Expert Syst. Appl.*, 202:117038, 2022.

K. Ellis, C. Wong, M. Nye, M. Sablé-Meyer, L. Morales, L. Hewitt, L. Cary, A. Solar-Lezama, and J. B. Tenenbaum. DreamCoder: Bootstrapping inductive program synthesis with wake-sleep library learning. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, Pldi 2021, pages 835–850, New York, NY, USA, June 2021. Association for Computing Machinery. ISBN 978-1-4503-8391-2.

D. Entner and P. O. Hoyer. On causal discovery from time series data using FCI. In P. Myllymäki, T. Roos, and T. Jaakkola, editors, *Proceedings of the 5th European Workshop on Probabilistic Graphical Models*, pages 121–128, Helsinki, FI, 2010. Helsinki Institute for Information Technology HIIT.

N. B. Erichson, L. Mathelin, Z. Yao, S. L. Brunton, M. W. Mahoney, and J. N. Kutz. Shallow neural networks for fluid flow reconstruction with limited sensors. *Proc. R. Soc. Lond. A*, 476: 20200097, 2020.

J. Evans and A. Rzhetsky. Machine science. *Science*, 329(5990):399–400, 2010.

R. Evans. The nature of the liquid-vapour interface and other topics in the statistical mechanics of non-uniform, classical fluids. *Advances in Physics*, 28(2):143–200, 1979.

V. Eyring, S. Bony, G. A. Meehl, C. A. Senior, B. Stevens, R. J. Stouffer, and K. E. Taylor. Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization. *Geoscientific Model Development*, 9(5):1937–1958, 2016.

V. Eyring, P. M. Cox, G. M. Flato, P. J. Gleckler, G. Abramowitz, P. Caldwell, W. D. Collins, B. K. Gier, A. D. Hall, F. M. Hoffman, et al. Taking climate model evaluation to the next level. *Nature Climate Change*, 9(2):102–110, 2019.

V. Eyring, L. Bock, A. Lauer, M. Righi, M. Schlund, B. Andela, E. Arnone, O. Bellprat, B. Brötz, L.-P. Caron, et al. Earth System Model Evaluation Tool (ESMValTool) v2. 0–An extended set of large-scale diagnostics for quasi-operational and comprehensive evaluation of earth system models in CMIP. *Geoscientific Model Development*, 13(7):3383–3438, 2020.

L. Faes, S. Erla, and G. Nollo. Measuring connectivity in linear multivariate processes: definitions, interpretation, and practical analysis. *Computational and mathematical methods in medicine*, 2012, 2012.

O. Fajardo-Fontiveros, I. Reichardt, H. R. De Los Ríos, J. Duch, M. Sales-Pardo, and R. Guimerà. Fundamental limits to learning closed-form mathematical models from data. *Nature Communications*, 14(1):1043, 2023.

B. Falkenhainer and R. Michalski. The structure mapping engine: Algorithm and examples. *Artificial Intelligence*, 32(1):1–63, 1986.

E. Feigenbaum, B. Buchanan, and J. Lederberg. The DENDRAL Project. *AI Magazine*, 2:37–46, 1971.

P. K. Feyerabend. *Problems of Empiricism*, volume 2. Cambridge University Press, 1981.

P. Fiedor. Networks in financial markets based on the mutual information rate. *Physical Review E*, 89(5):052801, 2014.

C. B. Field, R. B. Jackson, and H. A. Mooney. Stomatal responses to increased CO2: Implications from the plant to the global scale. *Plant, Cell & Environment*, 18(10):1214–1225, 1995.

R. A. Fisher. *The Design of Experiments*. Hafner Press, 1935.

P. Forré and J. M. Mooij. Constraint-based causal discovery for non-linear structural causal models with cycles and latent confounders. In A. Globerson and R. Silva, editors, *Proceedings of the 34th Conference on Uncertainty in Artificial Intelligence (UAI-18)*. AUAI Press, 2018.

S. Fortunato, C. T. Bergstrom, K. Börner, J. A. Evans, D. Helbing, S. Milojević, A. M. Petersen, F. Radicchi, R. Sinatra, B. Uzzi, et al. Science of science. *Science*, 359(6379):eaao0185, 2018.

L. C. Freeman. A set of measures of centrality based on betweenness. *Sociometry*, pages 35–41, 1977.

K. Friston. A theory of cortical responses. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 360(1456):815–836, 2005.

K. Friston, R. Moran, and A. K. Seth. Analysing connectivity with Granger causality and dynamic causal modelling. *Current opinion in neurobiology*, 23(2):172–178, 2013.

K. Friston, K. H. Preller, C. Mathys, H. Cagnan, J. Heinzle, A. Razi, and P. Zeidman. Dynamic causal modelling revisited. *NeuroImage*, 199:730–744, 2019.

K. J. Friston, L. Harrison, and W. Penny. Dynamic causal modelling. *Neuroimage*, 19(4):1273–1302, 2003.

K. J. Friston, T. Parr, P. Zeidman, A. Razi, G. Flandin, J. Daunizeau, O. J. Hulme, A. J. Billig, V. Litvak, R. J. Moran, et al. Dynamic causal modelling of COVID-19. *Wellcome open research*, 5, 2020.

K. J. Friston, G. Flandin, and A. Razi. Dynamic causal modelling of cOVID-19 and its mitigations. *Scientific reports*, 12(1):12419, 2022.

R. Fuentes, R. Nayek, P. Gardner, N. Dervilis, T. Rogers, K. Worden, and E. Cross. Equation discovery for nonlinear dynamical systems: A Bayesian viewpoint. *Mechanical Systems and Signal Processing*, 154:107528, 2021.

K. Fukami, T. Nakamura, and K. Fukagata. Convolutional neural network based hierarchical autoencoder for nonlinear mode decomposition of fluid field data. *Phys. Fluids*, 32:095110, 2020.

A. Gain and I. Shpitser. Structure learning under missing data. In *International conference on probabilistic graphical models*, pages 121–132. Pmlr, 2018.

A. R. Galgali, M. Sahani, and V. Mante. Residual dynamics resolves recurrent contributions to neural computation. *Nature Neuroscience*, 26(2):326–338, Feb. 2023.

T. Gao, D. Bhattacharjya, E. Nelson, M. Liu, and Y. Yu. IDYNO: Learning nonparametric DAGs from interventional dynamic data. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 6988–7001. Pmlr, 17–23 Jul 2022.

D. Geiger, T. Verma, and J. Pearl. Identifying independence in Bayesian networks. *Networks*, 20 (5):507–534, 1990.

P. Gelß, S. Klus, J. Eisert, and C. Schütte. Multidimensional approximation of nonlinear dynamical systems. *J. Comput. Nonlinear Dyn.*, 14:061006, 2019.

M. Genkin and T. A. Engel. Moving beyond generalization to accurate interpretation of flexible models. *Nat. Mach. Intell.*, 2(11):674–683, Nov. 2020.

M. Genkin, O. Hughes, and T. A. Engel. Learning non-stationary Langevin dynamics from stochastic observations of latent trajectories. *Nature Communications*, 12(1):5986, 2021.

A. Gerhardus. Characterization of causal ancestral graphs for time series with latent confounders. *Preprint arXiv:2112.08417*, 2021.

A. Gerhardus and J. Runge. High-recall causal discovery for autocorrelated time series with latent confounders. *Advances in Neural Information Processing Systems*, 33:12615–12625, 2020.

W. Gerstner, W. M. Kistler, R. Naud, and L. Paninski. *Neuronal Dynamics: From Single Neurons to Networks and Models of Cognition.* Cambridge University Press, 2014.

J. Geweke. Measurement of linear dependence and feedback between multiple time series. *Journal of the American statistical association*, 77(378):304–313, 1982.

D. Gillies. Artificial intelligence and scientific method. *Mind*, 107(428), 1998.

T. A. Glass, S. N. Goodman, M. A. Hernán, and J. M. Samet. Causal inference in public health. *Annual review of public health*, 34:61–75, 2013.

C. Glymour, K. Zhang, and P. Spirtes. Review of causal discovery methods based on graphical models. *Frontiers in Genetics*, 10, 2019.

K. Gödel. Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für mathematik und physik*, 38(1):173–198, 1931.

E. Gokcen, A. I. Jasper, J. D. Semedo, A. Zandvakili, A. Kohn, C. K. Machens, and B. M. Yu. Disentangling the flow of signals between populations of neurons. *Nat. Comput. Sci.*, 2(8): 512–525, Aug. 2022.

W. Gong, J. Jennings, C. Zhang, and N. Pawlowski. Rhino: Deep causal temporal relationship learning with history-dependent noise. In *The Eleventh International Conference on Learning Representations*, 2023.

P. J. Gonçalves, J.-M. Lueckmann, M. Deistler, M. Nonnenmacher, K. Öcal, G. Bassetto, C. Chintaluri, W. F. Podlaski, S. A. Haddad, T. P. Vogels, D. S. Greenberg, and J. H. Macke. Training deep neural density estimators to identify mechanistic models of neural dynamics. *eLife*, 9: e56261, Sept. 2020.

A. E. Goodwell, P. Jiang, B. L. Ruddell, and P. Kumar. Debates—Does information theory provide a new paradigm for Earth science? Causality, interaction, and feedback. *Water Resources Research*, 56(2):e2019WR024940, 2020.

A. Gordon, A. Moore, and A. Carlson. Using genetic algorithms to discover good representations. *Machine Learning*, 15(1):239–263, 1994.

A. Gozolchiani, S. Havlin, and K. Yamasaki. Emergence of El Niño as an autonomous component in the climate network. *Phys. Rev. Lett.*, 107(14):148501, 2011.

P. Gradu, T. Zrnic, Y. Wang, and M. Jordan. Valid inference after causal discovery. In *NeurIPS 2022 Workshop on Causality for Real-world Impact*, 2022.

C. W. J. Granger. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: journal of the Econometric Society*, 37:424–438, 1969.

C. W. J. Granger. Testing for causality: A personal viewpoint. *Journal of Economic Dynamics and control*, 2:329–352, 1980.

J. K. Green, S. I. Seneviratne, A. M. Berg, K. L. Findell, S. Hagemann, D. M. Lawrence, and P. Gentine. Large influence of soil moisture on long-term terrestrial carbon uptake. *Nature*, 565 (7740):476–479, 2019.

A. Gretton, K. Fukumizu, C. Teo, L. Song, B. Schölkopf, and A. Smola. A kernel statistical test of independence. *Advances in neural information processing systems*, 20, 2007.

N. Groun, M. Villalba-Orero, E. Lara-Pezzi, E. Valero, J. Garicano-Mena, and S. Le Clainche. Higher order dynamic mode decomposition: From fluid dynamics to heart disease analysis. *Preprint arXiv:2201.03030*, 2022.

E. Grunberg and F. Modigliani. The predictability of social events. *Journal of Political Economy*, 62(6):465–478, 1954.

Y. Guan, S. L. Brunton, and I. Novosselov. Sparse nonlinear models of chaotic electroconvection. *R. Soc. Open Sci.*, 8(8):202367, 2021.

L. Guastoni, A. Güemes, A. Ianiro, S. Discetti, P. Schlatter, H. Azizpour, and R. Vinuesa. Convolutional-network models to predict wall-bounded turbulence from wall quantities. *J. Fluid Mech.*, 928:A27, 2021.

A. Güemes, S. Discetti, A. Ianiro, B. Sirmacek, H. Azizpour, and R. Vinuesa. From coarse wall measurements to turbulent velocity fields through deep learning. *Phys. Fluids*, 33:075121, 2021.

R. Guimerà, I. Reichardt, A. Aguilar-Mogas, F. A. Massucci, M. Miranda, J. Pallarès, and M. Sales-Pardo. A Bayesian machine scientist to aid in the solution of challenging scientific problems. *Science advances*, 6(5):eaav6971, 2020.

J. Y. Halpern. *Actual causality*. MiT Press, 2016.

A. Hannart, J. Pearl, F. E. Otto, P. Naveau, and M. Ghil. Causal counterfactual theory for the attribution of weather and climate-related events. *Bull. Am. Meteorol. Soc.*, 97(1):99–110, 2016.

N. Hansen and A. Sokol. Causal interpretation of stochastic differential equations. *Electronic Journal of Probability*, 19:1 – 24, 2014.

T. Hastie, R. Tibshirani, and M. Wainwright. Statistical learning with sparsity. *Monographs on statistics and applied probability*, 143:143, 2015.

S. Haufe, F. Meinecke, K. Görgen, S. Dähne, J.-D. Haynes, B. Blankertz, and F. Bießmann. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *NeuroImage*, 87: 96–110, 2014.

C. Heinze-Deml, M. H. Maathuis, and N. Meinshausen. Causal structure learning. *Annual Review of Statistics and Its Application*, 5:371–391, 2018a.

C. Heinze-Deml, J. Peters, and N. Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2), 2018b.

W. Heisenberg. Über quantentheoretische Umdeutung kinematischer und mechanischer Beziehungen. *Zeitschrift für Physik*, 33(3):879–893, 1925.

C. G. Hempel. *The philosophy of Carl G. Hempel: Studies in science, explanation, and rationality*. Oxford University Press, 2001.

M. Hernan and J. Robins. *Causal Inference: What if*. Chapman & Hill/CRC, 2020.

M. A. Hernán. The C-word: Scientific euphemisms do not improve causal inference from observational data. *American journal of public health*, 108(5):616–619, 2018.

J. Hicks et al. *Causality in Economics*. Australian National University Press, 1980.

C. Hiemstra and J. D. Jones. Testing for linear and nonlinear Granger causality in the stock price-volume relation. *The Journal of Finance*, 49(5):1639–1664, 1994.

G. E. Hinton and R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313:504–507, 2006.

P. Hoyer, D. Janzing, J. M. Mooij, J. Peters, and B. Schölkopf. Nonlinear causal discovery with additive noise models. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems*, volume 21. Curran Associates, Inc., 2008.

B. Huang, K. Zhang, and B. Schölkopf. Identification of time-dependent causal model: A Gaussian process treatment. In *Proceedings of the 24th International Conference on Artificial Intelligence*, Ijcai'15, page 3561–3568. AAAI Press, 2015. ISBN 9781577357384.

B. Huang, K. Zhang, M. Gong, and C. Glymour. Causal Discovery and Forecasting in Nonstationary Environments with State-Space Models. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2901–2910. Pmlr, 09–15 Jun 2019.

B. Huang, K. Zhang, J. Zhang, J. Ramsey, R. Sanchez-Romero, C. Glymour, and B. Schölkopf. Causal discovery from heterogeneous/nonstationary data. *The Journal of Machine Learning Research*, 21(1):3482–3534, 2020.

Y. Huang and M. Valtorta. Pearl's calculus of intervention is complete. In R. Dechter and T. Richardson, editors, *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence*, Uai'06, page 217–224, Arlington, Virginia, USA, 2006. AUAI Press. ISBN 0974903922.

J. C. R. Hunt, A. A. Wray, and P. Moin. Eddies, streams, and convergence zones in turbulent flows. *Center for Turbulence Research (CTR) Proceedings of Summer Program*, 1998.

J. M. Hyman, L. Ma, E. Balaguer-Ballester, D. Durstewitz, and J. K. Seamans. Contextual encoding by ensembles of medial prefrontal cortex neurons. *Proc. Natl. Acad. Sci. U. S. A.*, 109(13): 5086–5091, Mar. 2012.

A. Hyttinen, F. Eberhardt, and P. O. Hoyer. Learning linear cyclic causal models with latent variables. *Journal of Machine Learning Research*, 13(109):3387–3439, 2012.

A. Hyvärinen, J. Karhunen, and E. Oja. *Independent Component Analysis*. Adaptive and Cognitive Dynamic Systems: Signal Processing, Learning, Communications and Control. Wiley, 2004. ISBN 9780471464198.

A. Hyvärinen, S. Shimizu, and P. O. Hoyer. Causal modelling combining instantaneous and lagged effects: An identifiable model based on non-Gaussianity. In *Proceedings of the 25th international conference on Machine learning*, pages 424–431, 2008.

A. Hyvärinen, K. Zhang, S. Shimizu, and P. O. Hoyer. Estimation of a structural vector autoregression model using non-Gaussianity. *Journal of Machine Learning Research*, 11(May):1709–1731, 2010.

G. W. Imbens and D. B. Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.

M. F. Jansen and I. M. Held. Parameterizing subgrid-scale eddy effects using energetically consistent backscatter. *Ocean Modelling*, 80:36–48, 2014.

D. Janzing. On causally asymmetric versions of Occam's Razor and their relation to thermodynamics. *arXiv preprint arXiv:0708.3411*, 2007.

J. Jiménez. Cascades in wall-bounded turbulence. *Annu. Rev. Fluid Mech.*, 44:27–45, 2012.

J. Jiménez. Optimal fluxes and Reynolds stresses. *J. Fluid Mech.*, 809:585–600, 2016.

J. Jiménez. Machine-aided turbulence theory. *J. Fluid Mech.*, 854:R1, 2018a.

J. Jiménez. Coherent structures in wall-bounded turbulence. *J. Fluid Mech.*, 842:P1, 2018b.

H. Johnson, G. Harris, and K. Williams. BRAINSFit: Mutual information registrations of whole-brain 3D images, using the Insight Toolkit. *The Insight Journal*, 2007.

J. E. Johnson, V. Laparra, and G. Camps-Valls. Disentangling derivatives, uncertainty and error in Gaussian process models. In *IGARSS 2018 - 2018 IEEE International Geoscience and Remote Sensing Symposium*, pages 4051–4054, July 2018.

P. N. Johnson-Laird and R. M. Byrne. *Deduction*. Lawrence Erlbaum Associates, Inc, 1991.

J. Juang and R. Pappa. An eigensystem realization algorithm for modal parameter identification and model reduction. *J. Guid. Control Dyn.*, 8:620, 1985.

M. Jung, M. Reichstein, P. Ciais, S. I. Seneviratne, J. Sheffield, M. L. Goulden, G. Bonan, A. Cescatti, J. Chen, R. De Jeu, et al. Recent decline in the global land evapotranspiration trend due to limited moisture supply. *Nature*, 467(7318):951–954, 2010.

J. Kaddour, A. Lynch, Q. Liu, M. J. Kusner, and R. Silva. Causal machine learning: A survey and open problems. *arXiv preprint arXiv:2206.15475*, 2022.

K. Kaheman, S. L. Brunton, and J. N. Kutz. Automatic differentiation to simultaneously identify nonlinear dynamics and extract noise probability distributions from data. *Machine Learning: Science and Technology*, 3(1):015031, Mar. 2022. Publisher: IOP Publishing.

E. Kaiser, J. N. Kutz, and S. L. Brunton. Data-driven approximations of dynamical systems operators for control. In *The Koopman Operator in Systems and Control*, pages 197–234. Springer, 2020.

E. Kaiser, J. N. Kutz, and S. L. Brunton. Data-driven discovery of Koopman eigenfunctions for control. *Machine Learning: Science and Technology*, 2(3):035023, June 2021a.

E. Kaiser, J. N. Kutz, and S. L. Brunton. Data-driven discovery of Koopman eigenfunctions for control. *Machine Learning: Science and Technology*, 2(3):035023, 2021b.

M. Kalisch and P. Bühlman. Estimating high-dimensional directed acyclic graphs with the PC-algorithm. *Journal of Machine Learning Research*, 8(3), 2007.

O. Kallenberg. *Probabilistic symmetries and invariance principles*, volume 9. Springer, 2005.

E. Kalnay, M. Kanamitsu, R. Kistler, W. Collins, D. Deaven, L. Gandin, M. Iredell, S. Saha, G. White, J. Woollen, et al. The NCEP/NCAR 40-year reanalysis project. *Bulletin of the American meteorological Society*, 77(3):437–472, 1996.

F. Kemeth, T. Bertalan, T. Thiem, F. Dietrich, S. Moon, C. Laing, and Y. Kevrekidis. Learning emergent partial differential equations in a learned emergent space. *Nature Communications*, 13, 06 2022.

J. M. Keynes. *A treatise on probability*. Courier Corporation, 2013.

M. Khodkar, P. Hassanzadeh, and A. Antoulas. A Koopman-based framework for forecasting the spatiotemporal evolution of chaotic dynamics with nonlinearities modeled as exogenous forcings. *Preprint arXiv:1909.00076*, 2019.

H. T. Kim, S. J. Kline, and W. C. Reynolds. The production of turbulence near a smooth wall in a turbulent boundary layer. *J. Fluid Mech.*, 50:133–160, 1971.

S. Kim, P. Y. Lu, S. Mukherjee, M. Gilbert, L. Jing, V. Čeperić, and M. Soljačić. Integration of Neural Network-Based Symbolic Regression in Deep Learning for Scientific Discovery. *IEEE Transactions on Neural Networks and Learning Systems*, 32(9):4166–4177, 2021.

R. King, P. Langley, and H. Simon. Automated Discovery in the Biological Sciences. *AI Magazine*, 25(3):21–36, 2004.

R. D. King, J. Rowland, W. Aubrey, M. Liakata, M. Markham, L. N. Soldatova, K. E. Whelan, A. Clare, M. Young, A. Sparkes, et al. The robot scientist Adam. *Computer*, 42(8):46–54, 2009.

D. Klahr and H. A. Simon. Studies of scientific discovery: Complementary approaches and convergent findings. *Psychological Bulletin*, 125(5):524, 1999.

S. J. Kline, W. C. Reynolds, F. A. Schraub, and P. W. Runstadler. The structure of turbulent boundary layers. *J. Fluid Mech.*, 30:741–773, 1967.

S. Klus, P. Koltai, and C. Schütte. On the numerical approximation of the Perron–Frobenius and Koopman operator. *Journal of Computational Dynamics*, 3:51–79, 2016.

S. Klus, F. Nüske, P. Koltai, H. Wu, I. Kevrekidis, C. Schütte, and F. Noé. Data-driven model reduction and transfer operator approximation. *Journal of Nonlinear Science*, 1010:9437–7, 2018.

S. Klus, I. Schuster, and K. Muandet. Eigendecompositions of transfer operators in reproducing kernel Hilbert spaces. *Journal of Nonlinear Science*, 30(1):283–315, 2020.

A. Kocabas. A Genetic Programming System for Automated Discovery in the Physical Sciences. *Machine Learning*, 7(3-4):295–314, 1991.

C. Koch. *Biophysics of Computation: Information Processing in Single Neurons (Computational Neuroscience Series)*. Oxford University Press, Inc., Usa, 2004. ISBN 0195181999.

M. Kokar. Knowledge acquisition: A realization of new artificial intelligence. *Artificial Intelligence*, 32:251–290, 1986.

A. N. Kolmogorov. On tables of random numbers. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 369–376, 1963.

M. Kommenda, B. Burlacu, G. Kronberger, and M. Affenzeller. Parameter identification for symbolic regression using nonlinear least squares. *Genetic Programming and Evolvable Machines*, 21(3):471–501, 2020.

G. Koppe, H. Toutounji, P. Kirsch, S. Lis, and D. Durstewitz. Identifying nonlinear dynamical systems via generative recurrent neural networks with applications to fMRI. *PLOS Computational Biology*, 15(8):1–35, 08 2019.

R. D. Koster, P. A. Dirmeyer, Z. Guo, G. Bonan, E. Chan, P. Cox, C. Gordon, S. Kanae, E. Kowalczyk, D. Lawrence, et al. Regions of strong coupling between soil moisture and precipitation. *Science*, 305(5687):1138–1140, 2004.

R. D. Koster, Y. Sud, Z. Guo, P. A. Dirmeyer, G. Bonan, K. W. Oleson, E. Chan, D. Verseghy, P. Cox, H. Davies, et al. GLACE: the global land–atmosphere coupling experiment. Part I: overview. *Journal of Hydrometeorology*, 7(4):590–610, 2006.

V. Kostic, P. Novelli, A. Maurer, C. Ciliberto, L. Rosasco, and M. Pontil. Learning dynamical systems via Koopman operator regression in reproducing kernel Hilbert spaces. In *NeurIPS 2022*, pages 1–9, 2022.

S. Kotz and D. Drouet. *Correlation and dependence*. World Scientific, 2001.

J. Koza, F. Bennett, D. Andre, and M. Keane. Nonlinear Genetic Programming: Automatic Discovery of Reusable Programs. *Machine Learning*, 42(1):185–223, 2001.

J. R. Koza. Genetic Programming: A Paradigm for Genetically Breeding Populations of Computer Programs to Solve Problems. Technical report, Dept. of Computer Science, Stanford University, Stanford, CA, USA, 1990.

J. R. Koza. Genetic programming: On the programming of computers by means of natural selection. *MIT Press*, 1992.

J. R. Koza. Genetic programming as a means for programming computers by natural selection. *Statistics and computing*, 4(2):87–112, 1994.

M. Kretschmer, D. Coumou, J. F. Donges, and J. Runge. Using causal effect networks to analyze different Arctic drivers of midlatitude winter circulation. *Journal of climate*, 29(11):4069–4081, 2016.

M. Kretschmer, J. Runge, and D. Coumou. Early prediction of extreme stratospheric polar vortex states based on causal precursors. *Geophysical research letters*, 44(16):8592–8600, 2017.

C. Krich, J. Runge, D. G. Miralles, M. Migliavacca, O. Perez-Priego, T. El-Madany, A. Carrara, and M. D. Mahecha. Estimating causal networks in biosphere–atmosphere interaction with the PCMCI approach. *Biogeosciences*, 17(4):1033–1061, 2020.

C. Krich, M. Migliavacca, D. G. Miralles, G. Kraemer, T. S. El-Madany, M. Reichstein, J. Runge, and M. D. Mahecha. Functional convergence of biosphere–atmosphere interactions in response to meteorological conditions. *Biogeosciences*, 18(7):2379–2404, 2021.

G. Kronberger, F. O. de Franca, B. Burlacu, C. Haider, and M. Kommenda. Shape-constrained symbolic regression—Improving extrapolation with prior knowledge. *Evolutionary Computation*, 30(1):75–98, 03 2022.

T. S. Kuhn. *The Structure of Scientific Revolutions*. University of Chicago Press, 1962.

J. Kutz, S. Brunton, B. Brunton, and J. Proctor. Dynamic mode decomposition: Data-driven modeling of complex systems. *Siam*, 2016.

J. Kwapień and S. Drożdż. Physical approach to complex systems. *Physics Reports*, 515(3-4): 115–226, 2012.

S. Lachapelle, P. Brouillard, T. Deleu, and S. Lacoste-Julien. Gradient-Based Neural DAG Learning. In *International Conference on Learning Representations*, 2020.

M. T. Landahl and M. T. Landahlt. Wave breakdown and turbulence. *SIAM J. Appl. Maths*, 28: 735–756, 1975.

P. Langley. Scientific discovery, causal explanation, and process model induction. *Mind & Society*, 18(1):43–56, 2019.

P. Langley, H. Simon, and G. Bradshaw. Scientific discovery: Computational explorations of the creative process. *AI Magazine*, 8(3):30–44, 1987a.

P. Langley, H. A. Simon, G. L. Bradshaw, and J. M. Zytkow. *Scientific discovery: Computational explorations of the creative processes*. MIT press, 1987b.

P. Langley, H. Simon, and G. Bradshaw. Automated discovery in the Physical Sciences. *AI Magazine*, 23(3):11– 28, 2002a.

P. Langley, H. Simon, and G. Bradshaw. Scientific discovery and the future of AI. *AI Magazine*, 23 (3):29–39, 2002b.

C. C. Lapish, E. Balaguer-Ballester, J. K. Seamans, A. G. Phillips, and D. Durstewitz. Amphetamine Exerts Dose-Dependent Changes in Prefrontal Cortex Attractor Dynamics during Working Memory. *Journal of Neuroscience*, 35(28):10172–10187. EB–B and CCL contributed equally., 2015.

E. Lazpita, A. Martínez-Sánchez, A. Corrochano, S. Hoyas, S. Le Clainche, and R. Vinuesa. On the generation and destruction mechanisms of arch vortices in urban fluid flows. *Phys. Fluids*, 34:051702, 2022.

S. Le Clainche and J. M. Vega. Higher Order Dynamic Mode Decomposition. *SIAM J. Appl. Dyn. Syst.*, 16:882–925, 2017.

J. Lee and T. A. Zaki. Detection algorithm for turbulent interfaces and large-scale structures in intermittent flows. *Comput. Fluids*, 175(1):142–158, 2018.

K. Lee and K. T. Carlberg. Model reduction of dynamical systems on nonlinear manifolds using deep convolutional autoencoders. *J. Comput. Phys.*, 404:108973, 2020.

S. LeRoy. *Causality in economics*. London School of Economics, Centre for Philosophy of Natural and Social Sciences, 2004.

S.-C. Lin and M. Oettel. A classical density functional from machine learning and a convolutional neural network. *SciPost Phys.*, 6:025, 2019.

S.-C. Lin, G. Martius, and M. Oettel. Analytical classical density functionals from an equation learning network. *The Journal of Chemical Physics*, 152(2):021102, 2020.

D. V. Lindley. *Understanding uncertainty*. John Wiley & Sons, 2013.

L. Ljung and T. Glad. On global identifiability for arbitrary model parametrizations. *Automatica*, 30(2):265–276, 1994.

G. Lohmann, K. Erfurth, K. Müller, and R. Turner. Critical comments on dynamic causal modelling. *NeuroImage*, 59(3):2322–2329, 2012.

J.-C. Loiseau. Data-driven modeling of the chaotic thermal convection in an annular thermosyphon. *Theor. Comput. Fluid Dyn.*, 34(4):339–365, 2020.

Z. Long, Y. Lu, and B. Dong. PDE-Net 2.0: Learning PDEs from data with a numeric-symbolic hybrid deep network. *Journal of Computational Physics*, 399:108925, Dec. 2019.

C. Louizos, M. Welling, and D. P. Kingma. Learning Sparse Neural Networks through $L_0$ Regularization. In *International Conference on Learning Representations*, 2018.

A. Lozano-Durán and J. Jiménez. Time-resolved evolution of coherent structures in turbulent channels: characterization of eddies and cascades. *J. Fluid Mech.*, 759:432–471, 2014.

A. Lozano-Durán, O. Flores, and J. Jiménez. The three-dimensional structure of momentum transfer in turbulent channels. *J. Fluid Mech.*, 694:100–130, 2012.

A. Lozano-Durán, H. J. Bae, and M. P. Encinar. Causality of energy-containing eddies in wall turbulence. *J. Fluid Mech.*, 882:A2, 2020.

S. S. Lu and W. W. Willmarth. Measurements of the structure of the Reynolds stress in a turbulent boundary layer. *J. Fluid Mech.*, 60:481–511, 1973.

J. Ludescher, M. Martin, N. Boers, A. Bunde, C. Ciemer, J. Fan, S. Havlin, M. Kretschmer, J. Kurths, J. Runge, et al. Network-based forecasting of climate phenomena. *Proceedings of the National Academy of Sciences*, 118(47):e1922872118, 2021.

J.-M. Lueckmann, P. J. Gonçalves, G. Bassetto, K. Öcal, M. Nonnenmacher, and J. H. Macke. Flexible Statistical Inference for Mechanistic Models of Neural Dynamics. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 1289–1299, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.

J. L. Lumley. The structure of inhomogeneous turbulence. *Atmospheric turbulence and wave propagation, A. M. Yaglom and V. I. Tatarski (eds). Nauka, Moscow*, pages 166–178, 1967.

S. Lun-Chau, R. Hu, J. Gonzalez, and D. Sejdinovic. RKHS-SHAP: Shapley values for kernel methods. *Preprint arXiv:2110.09167v2*, 2022.

M. Lungarella, A. Pitti, and Y. Kuniyoshi. Information transfer at multiple scales. *Physical Review E*, 76(5):056117, 2007.

B. Lusch, J. N. Kutz, and S. L. Brunton. Deep learning for universal linear embeddings of nonlinear dynamics. *Nature Communications*, 9(1), 2018.

J. M. Mooij and T. Claassen. Constraint-Based Causal Discovery using Partial Ancestral Graphs in the presence of Cycles. In J. Peters and D. Sontag, editors, *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 124 of *Proceedings of Machine Learning Research*, pages 1159–1168. Pmlr, 03–06 Aug 2020.

N. Madani, N. C. Parazoo, J. S. Kimball, A. P. Ballantyne, R. H. Reichle, M. Maneta, S. Saatchi, P. I. Palmer, Z. Liu, and T. Tagesson. Recent amplified global gross primary productivity due to temperature increase is offset by reduced productivity due to water constraints. *AGU Advances*, 1(4):e2020AV000180, 2020.

D. Malinsky and P. Spirtes. Causal structure learning from multivariate time series in settings with unmeasured confounding. In T. D. Le, K. Zhang, E. Kıcıman, A. Hyvärinen, and L. Liu, editors, *Proceedings of 2018 ACM SIGKDD Workshop on Causal Disocvery*, volume 92 of *Proceedings of Machine Learning Research*, pages 23–47, London, UK, 20 Aug 2018. PMLR.

D. Maraun, T. G. Shepherd, M. Widmann, G. Zappa, D. Walton, J. M. Gutiérrez, S. Hagemann, I. Richter, P. M. Soares, A. Hall, et al. Towards process-informed bias correction of climate change simulations. *Nature Climate Change*, 7(11):764–773, 2017.

D. Marinazzo, M. Pellicoro, and S. Stramaglia. Kernel-Granger causality and the analysis of dynamical networks. *Physical review E*, 77(5):056215, 2008a.

D. Marinazzo, M. Pellicoro, and S. Stramaglia. Kernel method for nonlinear Granger causality. *Phys. Rev. Lett.*, 100:144103, Apr. 2008b.

M. M. Marini and B. Singer. Causality in the social sciences. *Sociological methodology*, 18: 347–409, 1988.

A. C. Marreiros, K. E. Stephan, and K. J. Friston. Dynamic causal modeling. *Scholarpedia*, 5(7): 9568, 2010.

D. P. Marshall and A. J. Adcroft. Parameterization of ocean eddies: Potential vorticity mixing, energetics and arnold's first stability theorem. *Ocean Modelling*, 32(3-4):188–204, 2010.

A. Martínez-Sánchez, E. López, S. Le Clainche, A. Lozano-Durán, A. Srivastava, and R. Vinuesa. Causality analysis of large-scale structures in the flow around a wall-mounted square cylinder. *Preprint arXiv:2209.15356*, 2022.

A. Martínez-Sánchez, E. Lazpita, A. Corrochano, S. Le Clainche, S. Hoyas, and R. Vinuesa. Data-driven assessment of arch vortices in urban flows. *Int. J. Heat Fluid Flow, To Appear. Preprint arXiv:2202.01667v1*, 2023.

G. Martius and C. H. Lampert. Extrapolation and learning equations, 2016. https://arxiv.org/abs/1610.02995.

M. Mattia and P. Del Giudice. Population dynamics of interacting spiking neurons. *Phys. Rev. E*, 66:051917, Nov. 2002.

M. Mattia, M. Biggio, A. Galluzzi, and M. Storace. Dimensional reduction in networks of non-Markovian spiking neurons: Equivalence of synaptic filtering and heterogeneous propagation delays. *PLOS Computational Biology*, 15(10):1–35, Oct. 2019.

R. M. May. Will a large complex system be stable? *Nature*, 238(5364):413–414, 1972.

T. McConaghy. FFX: Fast, scalable, deterministic symbolic regression technology. In *Genetic Programming Theory and Practice IX*, Genetic and Evolutionary Computation, chapter 13, pages 235–260. Springer, Ann Arbor, USA, 2011.

A. McDavid, R. Gottardo, N. Simon, and M. Drton. Graphical models for zero-inflated single cell gene expression. *The annals of applied statistics*, 13(2):848, 2019.

R. T. McGibbon and V. S. Pande. Variational cross-validation of slow dynamical modes in molecular kinetics. *The Journal of Chemical Physics*, 142, 2015.

P. E. McKnight, K. M. McKnight, S. Sidani, and A. J. Figueredo. *Missing data: A gentle introduction*. Guilford Press, 2007.

J. D. Medaglia, M.-E. Lynall, and D. S. Bassett. Cognitive network neuroscience. *Journal of cognitive neuroscience*, 27(8):1471–1491, 2015.

C. Meek. Causal inference and causal explanation with background knowledge. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, Uai'95, page 403–410, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1558603859.

C. Meek. *Graphical models: Selecting causal and statistical models*. PhD thesis, Carnegie Mellon University., 1997.

D. M. A. Mehler and K. P. Kording. The lure of causal statements: Rampant mis-inference of causality in estimated connectivity, 2018.

C. Meneveau and J. Katz. Scale-invariance and turbulence models for large-eddy simulation. *Annual Review of Fluid Mechanics*, 32(1):1–32, 2000.

L. Menzly, T. Santos, and P. Veronesi. Understanding predictability. *Journal of Political Economy*, 112(1):1–47, 2004.

A. S. Meyer-Lindenberg, R. K. Olsen, P. D. Kohn, T. Brown, M. F. Egan, D. R. Weinberger, and K. F. Berman. Regionally Specific Disturbance of Dorsolateral Prefrontal–Hippocampal Functional Connectivity in Schizophrenia. *Archives of General Psychiatry*, 62(4):379–386, 04 2005.

I. Mezić. Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dynamics*, 41(1):309–325, 2005.

M. Milano and P. Koumoutsakos. Neural network modeling for near wall turbulent flow. *J. Comput. Phys.*, 182:1–26, 2002.

P. Milly. Potential evaporation and soil moisture in general circulation models. *Journal of climate*, 5(3):209–226, 1992.

J. Moehlis, H. Faisst, and B. Eckhardt. A low-dimensional model for turbulent shear flows. *New J. Phys.*, 6:56, 2004.

S. W. Mogensen. Equality Constraints in Linear Hawkes Processes. In B. Schölkopf, C. Uhler, and K. Zhang, editors, *Proceedings of the First Conference on Causal Learning and Reasoning*, volume 177 of *Proceedings of Machine Learning Research*, pages 576–593. Pmlr, 11–13 Apr 2022.

S. W. Mogensen and N. R. Hansen. Markov equivalence of marginalized local independence graphs. *The Annals of Statistics*, 48(1):539–559, 2020.

S. W. Mogensen, D. Malinsky, and N. R. Hansen. Causal Learning for Partially Observed Stochastic Dynamical Systems. In *Thirty-Fourth Conference on Uncertainty in Artifical Intelligence*. AUAI Press Corvallis, Oregon, 2018.

B. Monnier, B. Neiswander, and C. Wark. Stereoscopic particle image velocimetry measurements in an urban type boundary layer: Insight into flow regimes and incidence angle effect. *Boundary-Layer Meteorol.*, 135:243—-268, 2010.

D. Mønster, R. Fusaroli, K. Tylén, A. Roepstorff, and J. F. Sherson. Causal inference from noisy time-series data—Testing the convergent cross-mapping algorithm in the presence of noise and external influence. *Future Generation Computer Systems*, 73:52–62, 2017.

E. Montbrió, D. Pazó, and A. Roxin. Macroscopic description for networks of spiking neurons. *Phys. Rev. X*, 5:021028, June 2015.

R. P. Monti, l. Khemakhem, and A. Hyvarinen. Autoregressive flow-based causal discovery and inference. *arXiv preprint: 2007.09390*, 07 2020.

J. Mooij, D. Janzing, J. Peters, and B. Schölkopf. Regression by dependence minimization and its application to causal inference in additive noise models. In *Proceedings of the 26th annual international conference on machine learning*, pages 745–752, 2009.

J. M. Mooij, D. Janzing, and B. Schölkopf. From ordinary differential equations to structural causal models: The deterministic case. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*, Uai'13, page 440–448, Arlington, Virginia, USA, 2013a. AUAI Press.

J. M. Mooij, D. Janzing, and B. Schölkopf. From Ordinary Differential Equations to Structural Causal Models: the deterministic case. *arXiv preprint arXiv:1304.7920*, 2013b.

J. M. Mooij, J. Peters, D. Janzing, J. Zscheischler, and B. Schölkopf. Distinguishing cause from effect using observational data: methods and benchmarks. *The Journal of Machine Learning Research*, 17(1):1103–1204, 2016.

J. M. Mooij, S. Magliacane, and T. Claassen. Joint Causal Inference from Multiple Contexts. *Journal of Machine Learning Research*, 21(99):1–108, 2020.

R. Moraffah, P. Sheth, M. Karami, A. Bhattacharya, Q. Wang, A. Tahir, A. Raglin, and H. Liu. Causal inference for time series analysis: Problems, methods and evaluation. *Knowledge and Information Systems*, 63:3041–3085, 2021.

S. L. Morgan and C. Winship. *Counterfactuals and causal inference*. Cambridge University Press, 2015.

P. Moulet. Learning Rules from Structured Data. *Machine Learning*, 8(1):47–75, 1992.

P. Munz. *Our knowledge of the growth of knowledge: Popper or Wittgenstein?* Routledge, 2014.

T. Murata and K. Tanaka. A Constructive Induction Algorithm Incorporating Prior Knowledge. *Machine Learning*, 14(1):71–96, 1994.

T. Murata, K. Fukami, and K. Fukagata. Nonlinear mode decomposition with convolutional neural networks for fluid dynamics. *J. Fluid Mech.*, 882:A13, 2020.

R. O. Ness, K. Sachs, P. Mallick, and O. Vitek. A Bayesian active learning experimental design for inferring signaling networks. In *Research in Computational Molecular Biology: 21st Annual International Conference, RECOMB 2017, Hong Kong, China, May 3-7, 2017, Proceedings 21*, pages 134–156. Springer, 2017.

W.-J. Neumann, A. Horn, and A. A. Kühn. Insights and opportunities for deep brain stimulation as a brain circuit intervention, 2023.

I. Newton. *Philosophiae naturalis principia mathematica*, volume 1. G. Brookman, 1833.

I. Ng, A. Ghassami, and K. Zhang. On the Role of Sparsity and DAG Constraints for Learning Linear DAGs. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 17943–17954. Curran Associates, Inc., 2020.

I. Ng, S. Lachapelle, N. Rosemary Ke, S. Lacoste-Julien, and K. Zhang. On the Convergence of Continuous Constrained Optimization for Structure Learning. In G. Camps-Valls, F. J. R. Ruiz, and I. Valera, editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 8176–8198. Pmlr, 28–30 Mar 2022.

E. Nieh, M. Schottdorf, N. Freeman, R. Low, S. Lewallen, S. Koay, L. Pinto, J. Gauthier, C. Brody, and D. Tank. Geometry of abstract learned knowledge in the hippocampus. *Nature*, 595(7865): 80–84, July 2021. doi: 10.1038/s41586-021-03652-7.

D. Niemeijer and R. S. de Groot. Framing environmental indicators: moving from causal chains to causal networks. *Environment, development and sustainability*, 10(1):89–106, 2008.

B. R. Noack, K. Afanasiev, M. Morzynski, G. Tadmor, and F. Thiele. A hierarchy of low-dimensional models for the transient and post-transient cylinder wake. *J. Fluid Mech.*, 497: 335–363, 2003.

F. Noé and F. Nüske. A variational approach to modeling slow processes in stochastic dynamical systems. *Multiscale Modeling & Simulation*, 11:635–655, 2013.

K. Nordhausen and P. Langley. Inverse Entailment and Progol. *Machine Learning*, 5(1):25–38, 1990.

P. Nowack, J. Runge, V. Eyring, and J. D. Haigh. Causal networks for climate model evaluation and constrained projections. *Nature communications*, 11(1):1–11, 2020.

F. Nüske, B. G. Keller, G. Pérez-Hernández, A. S. J. S. Mey, and F. Noé. Variational approach to molecular kinetics. *Journal of Chemical Theory and Computation*, 10:1739–1752, 2014.

M. L. Observatory. https://www.climate.gov/teaching/resources/atmospheric-co2-mauna-loa-observatory, 2020.

P. J. Olver. *Classical Invariant Theory*. London Mathematical Society Student Texts. Cambridge University Press, 1999.

W. M. Orr. The stability or instability of the steady motions of a perfect liquid and of a viscous liquid. Part II. A viscous liquid. *Math. Proc. R. Irish Acad.*, 27:69–138, 1907.

R. Pamfil, N. Sriwattanaworachai, S. Desai, P. Pilgerstorfer, K. Georgatzis, P. Beaumont, and B. Aragam. DYNOTEARS: Structure Learning from Time-Series Data. In S. Chiappa and R. Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 1595–1605. Pmlr, 26–28 Aug 2020.

C. Pandarinath, D. O'Shea, J. Collins, R. Jozefowicz, S. Stavisky, J. Kao, E. Trautmann, M. Kaufman, S. Ryu, L. Hochberg, J. Henderson, K. Shenoy, L. Abbott, and D. Sussillo. Inferring single-trial neural population dynamics using sequential auto-encoders. *Nature Methods*, 15, 10 2018. doi: 10.1038/s41592-018-0109-9.

J. Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1988. ISBN 0934613737.

J. Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688, 1995.

J. Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, New York, NY, USA, 2nd edition, 2009a. ISBN 052189560x, 9780521895606.

J. Pearl. Causal inference in statistics: An overview. *Stat. Surv.*, 3:96–146, 2009b.

J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, Cambridge, UK, 2nd edition, 2009c.

J. Pearl. Statistics and Causality: Separated to Reunite—Commentary on Bryan Dowd's "Separated at Birth". *Health Services Research*, 46(2):421–429, 2011.

J. Pearl and D. Mackenzie. *The Book of Why: The New Science of Cause and Effect*. Basic books, New York, 2018.

J. Pearl, M. Glymour, and N. P. Jewell. *Causal inference in statistics: A Primer*. John Wiley & Sons, 2016. ISBN 1935-7516.

W. D. Penny, K. E. Stephan, A. Mechelli, and K. J. Friston. Comparing dynamic causal models. *Neuroimage*, 22(3):1157–1172, 2004.

W. D. Penny, K. E. Stephan, J. Daunizeau, M. J. Rosa, K. J. Friston, T. M. Schofield, and A. P. Leff. Comparing families of dynamic causal models. *PLoS computational biology*, 6(3):e1000709, 2010.

A. Pérez-Suay and G. Camps-Valls. Causal Inference in Geoscience and Remote Sensing from Observational Data. *IEEE Transactions on Geoscience and Remote Sensing*, 57(3):1502–1513, 2019.

J. Peters, J. M. Mooij, D. Janzing, and B. Schölkopf. Identifiability of Causal Graphs Using Functional Models. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, Uai'11, page 589–598, Arlington, Virginia, USA, 2011. AUAI Press. ISBN 9780974903972.

J. Peters, D. Janzing, and B. Schölkopf. Causal Inference on Time Series using Restricted Structural Equation Models. In C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.

J. Peters, P. Bühlmann, and N. Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):947–1012, 2016.

J. Peters, D. Janzing, and B. Schölkopf. *Elements of Causal Inference: Foundations and Learning Algorithms*. MIT Press, Cambridge, MA, USA, 2017a.

J. Peters, D. Janzing, and B. Schölkopf. *Elements of causal inference: foundations and learning algorithms*. The MIT Press, 2017b.

J. Peters, S. Bauer, and N. Pfister. *Causal Models for Dynamical Systems*, page 671–690. Association for Computing Machinery, New York, NY, USA, 1 edition, 2022. ISBN 9781450395861.

J. M. Peters. *Restricted structural equation models for causal inference*. PhD thesis, ETH Zurich and MPI for Intelligent Systems, 2012.

B. K. Petersen, M. L. Larma, T. N. Mundhenk, C. P. Santiago, S. K. Kim, and J. T. Kim. Deep symbolic regression: Recovering mathematical expressions from data via risk-seeking policy gradients. In *International Conference on Learning Representations*, 2021.

L. Petersen and N. R. Hansen. Testing Conditional Independence via Quantile Regression Based Partial Copulas. *Journal of Machine Learning Research*, 22(70):1–47, 2021.

N. Pfister, P. Bühlmann, and J. Peters. Invariant causal prediction for sequential data. *Journal of the American Statistical Association*, 114(527):1264–1276, 2019.

S. B. Pope. Turbulent Flows. *Cambridge University Press*, 2000.

K. Popper. *The logic of scientific discovery*. Routledge, 2005.

R. Potthast. *Amari Model*, pages 1–6. Springer New York, New York, NY, 2013. ISBN 978-1-4614-7320-6.

R. Praksova. Eureqa: Software review. *Genet. Program. Evol. M.*, 12(1):173–178, 2011.

S. J. Press. A compound events model for security prices. *Journal of business*, pages 317–335, 1967.

A. Pukrittayakamee, M. Malshe, M. Hagan, L. Raff, R. Narulkar, S. Bukkapatnum, and R. Komanduri. Simultaneous fitting of a potential-energy surface and its corresponding force fields using feedforward neural networks. *The Journal of chemical physics*, 130(13):134101, 2009.

M. Quade, J. Gout, and M. Abel. Glyph: Symbolic regression tools. *Journal of Open Research Software*, June 2019.

M. I. Rabinovich and P. Varona. Discrete Sequential Information Coding: Heteroclinic Cognitive Dynamics. *Frontiers in Computational Neuroscience*, 12, 2018.

M. I. Rabinovich, R. Huerta, P. Varona, and V. S. Afraimovich. Transient Cognitive Dynamics, Metastability, and Decision Making. *PLOS Computational Biology*, 4(5):1–9, 05 2008.

F. Raia. Causality in complex dynamic systems: A challenge in earth systems science education. *Journal of Geoscience Education*, 56(1):81–94, 2008.

J. Ramsey, P. Spirtes, and J. Zhang. Adjacency-Faithfulness and Conservative Causal Inference. In *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence*, Uai'06, page 401–408, Arlington, Virginia, USA, 2006. AUAI Press. ISBN 0974903922.

J. Ramsey, M. Glymour, R. Sanchez-Romero, and C. Glymour. A million variables and more: the fast greedy equivalence search algorithm for learning high-dimensional graphical causal models, with an application to functional magnetic resonance images. *International journal of data science and analytics*, 3(2):121–129, 2017.

C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 2006.

A. Razi and K. J. Friston. The connected brain: causality, models, and intrinsic dynamics. *IEEE Signal Processing Magazine*, 33(3):14–35, 2016.

H. Reichenbach. *The direction of time*, volume 65. Univ of California Press, 1991.

A. Reid, D. Headley, R. Mill, R. sanchez romero, L. Uddin, D. Marinazzo, D. Lurie, P. Valdés-Sosa, S. Hanson, B. Biswal, V. Calhoun, R. Poldrack, and M. Cole. Advancing functional connectivity research from association to causation. *Nature Neuroscience*, 22:1–10, Oct. 2019.

A. G. Reisach, C. Seiler, and S. Weichwald. Beware of the Simulated DAG! Causal Discovery Benchmarks May Be Easy To Game. In *Advances in Neural Information Processing Systems 34 (NeurIPS)*, 2021.

F. Reitsma. Geoscience explanations: Identifying what is needed for generating scientific narratives from data models. *Environmental Modelling & Software*, 25(1):93–99, 2010.

O. Reynolds. On the dynamical theory of incompressible viscous fluids and the determination of the criterion. *Phil. Trans. R. Soc. A*, 186:123–164, 1895.

T. Richardson. A Discovery Algorithm for Directed Cyclic Graphs. In *Proceedings of the Twelfth International Conference on Uncertainty in Artificial Intelligence*, Uai'96, page 454–461, San Francisco, CA, USA, 1996. Morgan Kaufmann Publishers Inc. ISBN 155860412x.

T. Richardson and P. Spirtes. Ancestral Graph Markov Models. *The Annals of Statistics*, 30(4): 962–1030, 2002.

J. Richiardi, S. Achard, H. Bunke, and D. Van De Ville. Machine learning with brain graphs: predictive modeling approaches for functional imaging in systems neuroscience. *IEEE Signal processing magazine*, 30(3):58–70, 2013.

A. Rinaldo, L. Wasserman, and M. G'Sell. Bootstrapping and sample splitting for high-dimensional, assumption-lean inference. *The Annals of Statistics*, 47(6):3438–3469, 2019.

J. M. Robins, R. Scheines, P. Spirtes, and L. Wasserman. Uniform consistency in causal inference. *Biometrika*, 90(3):491–515, 2003.

M. Robins and M. Hernan. Causal inference: what if. *Found Agnostic Stat*, pages 235–281, 2020.

M. Rolinek, D. Zietlow, and G. Martius. Variational Autoencoders Pursue PCA Directions (by Accident). In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 12406–12415, June 2019.

A. S. Ross, Z. Li, P. Perezhogin, C. Fernandez-Granda, and L. Zanna. Benchmarking of machine learning ocean subgrid parameterizations in an idealized model. *Journal of Advances in Modeling Earth Systems*, 15:e2022MS003258, 2023.

L. N. Ross. Dynamical Models and Explanation in Neuroscience. *Philosophy of Science*, 82(1): 32–54, 2015.

C. Rowley, I. Mezic, S. Bagheri, P. Schlatter, and D. Henningson. Spectral analysis of nonlinear flows. *J. Fluid Mech.*, 641:115–127, 2009a.

C. W. Rowley and S. T. Dawson. Model reduction for flow analysis and control. *Annu. Rev. Fluid Mech.*, 49:387–417, 2017.

C. W. Rowley, I. Mezić, S. Bagheri, P. Schlatter, and D. S. Henningson. Spectral analysis of nonlinear flows. *Journal of Fluid Mechanics*, 641:115–127, 2009b.

P. K. Rubenstein, B. Bongers, S. Bernhard, and J. M. Mooij. From Deterministic ODEs to Dynamic Structural Causal Models. In *Thirty-Fourth Conference on Uncertainty in Artifical Intelligence*. AUAI Press Corvallis, Oregon, 2018.

J. Runge. Quantifying information transfer and mediation along causal pathways in complex systems. *Physical Review E*, 92(6):062829, 2015.

J. Runge. Conditional independence testing based on a nearest-neighbor estimator of conditional mutual information. In A. Storkey and F. Perez-Cruz, editors, *International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 938–947. Pmlr, 2018.

J. Runge. Discovering contemporaneous and lagged causal relations in autocorrelated nonlinear time series datasets. In J. Peters and D. Sontag, editors, *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 124 of *Proceedings of Machine Learning Research*, pages 1388–1397. Pmlr, 03–06 Aug 2020.

J. Runge. Necessary and sufficient graphical conditions for optimal adjustment sets in causal graphical models with hidden variables. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, 2021.

J. Runge, J. Heitzig, V. Petoukhov, and J. Kurths. Escaping the Curse of Dimensionality in Estimating Multivariate Transfer Entropy. *Physical Review Letters*, 108:258701, June 2012.

J. Runge, V. Petoukhov, and J. Kurths. Quantifying the strength and delay of climatic interactions: The ambiguities of cross correlation and a novel measure based on graphical models. *Journal of climate*, 27(2):720–739, 2014.

J. Runge, R. V. Donner, and J. Kurths. Optimal model-free prediction from multivariate time series. *Physical Review E*, 91(5):052909, 2015a.

J. Runge, V. Petoukhov, J. F. Donges, J. Hlinka, N. Jajcay, M. Vejmelka, D. Hartman, N. Marwan, M. Paluš, and J. Kurths. Identifying causal gateways and mediators in complex spatio-temporal systems. *Nature communications*, 6(1):1–10, 2015b.

J. Runge, S. Bathiany, E. Bollt, G. Camps-Valls, D. Coumou, E. Deyle, C. Glymour, M. Kretschmer, M. D. Mahecha, J. Muñoz-Marí, et al. Inferring causation from time series in Earth system sciences. *Nature communications*, 10(1):1–13, 2019a.

J. Runge, P. Nowack, M. Kretschmer, S. Flaxman, and D. Sejdinovic. Detecting and quantifying causal associations in large nonlinear time series datasets. *Science advances*, 5(11):eaau4996, 2019b.

J. Runge, X.-A. Tibau, M. Bruhns, J. Muñoz-Marí, and G. Camps-Valls. The causality for climate competition. In *NeurIPS 2019 Competition and Demonstration Track*, pages 110–120. Pmlr, 2020.

J. Runge, A. Gerhardus, G. Varando, V. Eyring, and G. Camps-Valls. Causal inference for time series. *Nature Reviews Earth & Environment*, 10:2553, 2023.

S. Russell. Human-compatible artificial intelligence. *Human-Like Machine Intelligence*, pages 3–23, 2021.

F. Russo. *Causality and causal modelling in the social sciences*. Springer, 2010.

V. Rutten, A. Bernacchia, M. Sahani, and G. Hennequin. Non-reversible Gaussian processes for identifying latent dynamical structure in neural data. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 9622–9632. Curran Associates, Inc., 2020.

E. Saggioro, J. de Wiljes, M. Kretschmer, and J. Runge. Reconstructing regime-dependent causal relationships from observational time series. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 30(11):113115, 2020.

S. S. Sahoo, C. H. Lampert, and G. Martius. Learning equations for extrapolation and control. In J. Dy and A. Krause, editors, *Proc. 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden*, volume 80, pages 4442–4450. Pmlr, 2018.

S. Salcedo-Sanz, D. Casillas-Pérez, J. Del Ser, C. Casanova-Mateo, L. Cuadra, M. Piles, and G. Camps-Valls. Persistence in complex systems. *Physics Reports*, 957:1–73, 2022.

O. Sani, H. Abbaspourazad, Y. Wong, B. Pesaran, and M. Shanechi. Modeling behaviorally relevant neural dynamics enabled by preferential subspace identification. *Nature Neuroscience*, 24: 140–149, 2021. doi: 10.1038/s41593-020-00733-0.

C. Schaffer. Constructing Explanations for Propositional Knowledge Bases. *Machine Learning*, 4 (4):321–353, 1990.

M. Schirner, L. Domide, D. Perdikis, P. Triebkorn, L. Stefanovski, R. Pai, P. Prodan, B. Valean, J. Palmer, C. Langford, A. Blickensdörfer, M. van der Vlag, S. Diaz-Pier, A. Peyser, W. Klijn, D. Pleiter, A. Nahm, O. Schmid, M. Woodman, L. Zehl, J. Fousek, S. Petkoski, L. Kusch, M. Hashemi, D. Marinazzo, J.-F. Mangin, A. Flöel, S. Akintoye, B. C. Stahl, M. Cepic, E. Johnson, G. Deco, A. R. McIntosh, C. C. Hilgetag, M. Morgan, B. Schuller, A. Upton, C. McMurtrie, T. Dickscheid, J. G. Bjaalie, K. Amunts, J. Mersmann, V. Jirsa, and P. Ritter. Brain simulation as a cloud service: The Virtual Brain on EBRAINS. *NeuroImage*, 251:118973, 2022.

D. Schmekel, F. Alcántara-Ávila, S. Hoyas, and R. Vinuesa. Predicting coherent turbulent structures via deep learning. *Front. Phys.*, 10:888832, 2022.

P. J. Schmid. Dynamic mode decomposition of numerical and experimental data. *Journal of Fluid Mechanics*, 656:5–28, 2010.

M. Schmidt and H. Lipson. Distilling Free-Form Natural Laws from Experimental Data. *Science*, 324(5923):81–85, 2009a.

M. Schmidt and H. Lipson. Distilling Free-Form Natural Laws from Experimental Data. *Science*, 324(5923):81–85, 2009b.

M. D. Schmidt, R. R. Vallabhajosyula, J. W. Jenkins, J. E. Hood, A. S. Soni, J. P. Wikswo, and H. Lipson. Automated refinement and inference of analytical models for metabolic networks. *Physical biology*, 8(5):055011, 2011.

S. Schneider, J. H. Lee, and M. W. Mathis. Learnable latent embeddings for joint behavioural and neural analysis. *Nature*, May 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06031-6. URL https://doi.org/10.1038/s41586-023-06031-6.

B. Schölkopf and A. Smola. *Learning with Kernels*. MIT Press, Cambridge, MA, USA, 2008.

B. Schölkopf, F. Locatello, S. Bauer, N. R. Ke, N. Kalchbrenner, A. Goyal, and Y. Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.

T. Schreiber. Measuring information transfer. *Physical review letters*, 85(2):461, 2000.

K. T. Schütt, F. Arbabzadah, S. Chmiela, K. R. Müller, and A. Tkatchenko. Quantum-chemical insights from deep tensor neural networks. *Nature communications*, 8(1):13890, 2017.

C. R. Schwantes and V. S. Pande. Modeling molecular kinetics with tICA and the kernel trick. *Journal of Chemical Theory and Computation*, 11:600–608, 2015.

J. Scott Armstrong and F. Collopy. Causal forces: Structuring knowledge for time-series extrapolation. *Journal of Forecasting*, 12(2):103–115, 1993.

W. Sellars et al. Empiricism and the Philosophy of Mind. *Minnesota studies in the philosophy of science*, 1(19):253–329, 1956.

R. D. Shah and J. Peters. The hardness of conditional independence testing and the generalised covariance measure. *The Annals of Statistics*, 48(3):1514 – 1538, 2020.

C. E. Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27 (3):379–423, 1948.

D. Shea, S. Brunton, and J. Kutz. SINDy-BVP: Sparse identification of nonlinear dynamics for boundary value problems. *Phys. Rev. Res.*, 3:023255, 2021.

A. Sheikhattar, S. Miran, J. Liu, J. B. Fritz, S. A. Shamma, P. O. Kanold, and B. Babadi. Extracting neuronal functional network dynamics via adaptive Granger causality analysis. *Proceedings of the National Academy of Sciences*, 115(17):E3869–e3878, 2018.

T. G. Shepherd. Storyline approach to the construction of regional climate change information. *Proceedings of the Royal Society A*, 475(2225):20190013, 2019.

S. Shimizu, P. O. Hoyer, A. Hyväinen, and A. Kerminen. A Linear Non-Gaussian Acyclic Model for Causal Discovery. *Journal of Machine Learning Research*, 7(72):2003–2030, 2006.

S. Shimizu, T. Inazumi, Y. Sogawa, A. Hyvärinen, Y. Kawahara, T. Washio, P. O. Hoyer, and K. Bollen. DirectLiNGAM: A direct method for learning a linear non-Gaussian structural equation model. *The Journal of Machine Learning Research*, 12:1225–1248, 2011.

I. Shpitser and J. Pearl. Identification of Conditional Interventional Distributions. In R. Dechter and T. Richardson, editors, *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence*, Uai'06, page 437–444, Arlington, Virginia, USA, 2006. AUAI Press. ISBN 0974903922.

I. Shpitser and J. Pearl. Complete identification methods for the causal hierarchy. *Journal of Machine Learning Research*, 9:1941–1979, 2008.

J. Shrager and P. Langley. *Computational Models of Scientific Discovery and Theory Formation*. Morgan Kaufmann, 1990.

S. H. Siddiqi, K. P. Kording, J. Parvizi, and M. D. Fox. Causal mapping of human brain function. *Nature reviews neuroscience*, 23(6):361–375, 2022.

N. Simidjievski, L. Todorovski, J. Kocijan, and S. Džeroski. Equation discovery for nonlinear system identification. *IEEE Access*, 8:29930–29943, 2020.

H. A. Simon et al. The scientist as problem solver. *Complex information processing: The impact of Herbert A. Simon*, pages 375–398, 1989.

P. R. Spalart. Strategies for turbulence modelling and simulations. *Int. J. Heat Fluid Flow*, 21: 252–263, 2000.

P. Spirtes and C. Glymour. An algorithm for fast recovery of sparse causal graphs. *Social science computer review*, 9(1):62–72, 1991.

P. Spirtes, C. Meek, and T. Richardson. Causal Inference in the Presence of Latent Variables and Selection Bias. In P. Besnard and S. Hanks, editors, *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, Uai'95, page 499–506, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1558603859.

P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search*. MIT Press, Boston, 2000.

O. Stegle, D. Janzing, K. Zhang, J. M. Mooij, and B. Schölkopf. Probabilistic latent variable models for distinguishing between cause and effect. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, volume 23. Curran Associates, Inc., 2010.

K. E. Stephan, W. D. Penny, R. J. Moran, H. E. den Ouden, J. Daunizeau, and K. J. Friston. Ten simple rules for dynamic causal modeling. *Neuroimage*, 49(4):3099–3109, 2010.

T. Stephens. gplearn: Genetic programming in Python with a scikit-learn inspired API. https://gplearn.readthedocs.io/en/stable, 2022.

J. D. Sterman. Learning in and about complex systems. *System dynamics review*, 10(2-3):291–330, 1994.

T. F. Stocker, D. Qin, G. Plattner, M. Tignor, S. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex, and P. Midgley. Climate change 2013: the physical science basis. Intergovernmental panel on climate change, working group I contribution to the IPCC fifth assessment report (AR5). *New York*, 2013.

P. A. Stokes and P. L. Purdon. A study of problems encountered in Granger causality analysis from a neuroscience perspective. *Proceedings of the national academy of sciences*, 114(34): E7063–e7072, 2017.

E. V. Strobl. A constraint-based algorithm for causal discovery with cycles, latent variables and selection bias. *International Journal of Data Science and Analytics*, 8(1):33–56, 2019.

E. V. Strobl, S. Visweswaran, and P. L. Spirtes. Fast causal inference with non-random missingness by test-wise deletion. *International journal of data science and analytics*, 6:47–62, 2018.

G. Sugihara, R. May, H. Ye, C.-h. Hsieh, E. Deyle, M. Fogarty, and S. Munch. Detecting causality in complex ecosystems. *science*, 338(6106):496–500, 2012.

M. Sugiyama and M. Kawanabe. *Machine learning in non-stationary environments: Introduction to covariate shift adaptation*. MIT press, 2012.

X. Sun, O. Schulte, G. Liu, and P. Poupart. NTS-NOTEARS: Learning nonparametric DBNs with prior knowledge. In *The 26th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2023.

S. S. Suseela, Y. Feng, and K. Mao. A comparative study on machine learning algorithms for knowledge discovery. In *2022 17th International Conference on Control, Automation, Robotics and Vision (ICARCV)*, pages 131–136. Ieee, 2022.

J. D. Swearingen and R. F. Blackwelder. The growth and breakdown of streamwise vortices in the presence of a wall. *J. Fluid Mech.*, 182:255–290, 1987.

A. Tabas, M. Andermann, V. Schuberth, H. Riedel, E. Balaguer-Ballester, and A. Rupp. Modeling and MEG evidence of early consonance processing in auditory cortex. *PLOS Computational Biology*, 15(2):1–28, 02 2019.

K. Taira, S. L. Brunton, S. Dawson, C. W. Rowley, T. Colonius, B. J. McKeon, O. T. Schmidt, S. Gordeyev, V. Theofilis, and L. S. Ukeiley. Modal analysis of fluid flows: An overview. *Aiaa J.*, 55(12):4013–4041, 2017.

K. Takagi. Principles of mutual information maximization and energy minimization affect the activation patterns of large scale networks in the brain. *Frontiers in Computational Neuroscience*, 13, 2020.

N. Takeishi, Y. Kawahara, and T. Yairi. Learning koopman invariant subspaces for dynamic mode decomposition. *Advances in neural information processing systems*, 30, 2017.

F. Takens. Detecting strange attractors in turbulence. In D. A. Rand and L.-S. Young, editors, *Dynamical Systems and Turbulence, Warwick 1980*, volume 898 of *Lecture Notes in Mathematics*, pages 366–381. Springer, Berlin, 1981. ISBN 978-3-540-11171-9.

A. Tarski and J. Tarski. *Introduction to Logic and to the Methodology of the Deductive Sciences*. Number 24. Oxford University Press on Demand, 1994.

H. Tennekes and J. L. Lumley. A first course in turbulence. *MIT press*, 1972.

T. N. Thiem, M. Kooshkbaghi, T. Bertalan, C. R. Laing, and I. G. Kevrekidis. Emergent spaces for coupled oscillators. *Front. Comput. Neurosci.*, 14:36, May 2020.

X.-A. Tibau, C. Reimers, A. Gerhardus, J. Denzler, V. Eyring, and J. Runge. A spatiotemporal stochastic climate model for benchmarking causal discovery methods for teleconnections. *Environmental Data Science*, 1:e12, 2022.

R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society (Series B)*, 58:267–288, 1996.

M. E. Tipping. Sparse Bayesian learning and the relevance vector machine. *Journal of machine learning research*, 1(Jun):211–244, 2001.

M. T. Todd, L. E. Nystrom, and J. D. Cohen. Confounds in multivariate pattern analysis: Theory and rule representation case study. *NeuroImage*, 77:157–165, 2013.

E. Tognoli and J. S. Kelso. The metastable brain. *Neuron*, 81(1):35–48, 2014.

J. H. Tu, C. W. Rowley, D. M. Luchtenburg, S. L. Brunton, and J. N. Kutz. On dynamic mode decomposition: Theory and applications. *Journal of Computational Dynamics*, 1, 2014.

R. Tu, C. Zhang, P. Ackermann, K. Mohan, H. Kjellström, and K. Zhang. Causal discovery in the presence of missing data. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1762–1770. Pmlr, 2019.

S.-M. Udrescu and M. Tegmark. AI Feynman: A physics-inspired method for symbolic regression. *Science Advances*, 6(16):eaay2631, 2020.

S.-M. Udrescu, A. Tan, J. Feng, O. Neto, T. Wu, and M. Tegmark. AI Feynman 2.0: Pareto-optimal symbolic regression exploiting graph modularity. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 4860–4871. Curran Associates, Inc., 2020.

C. Uhler, G. Raskutti, P. Bühlmann, and B. Yu. Geometry of the faithfulness assumption in causal inference. *The Annals of Statistics*, pages 436–463, 2013.

S. M. Ulam. *A Collection of Mathematical Problems*. Interscience Publisher NY, 1960.

L. Van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11), 2008.

T. VanderWeele. *Explanation in causal inference: Methods for mediation and interaction*. Oxford University Press, 2015.

V. N. Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10 (5):988–999, 1999.

G. Varando and N. R. Hansen. Graphical continuous Lyapunov models. In *Conference on Uncertainty in Artificial Intelligence*, pages 989–998. PMLR, 2020.

A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

T. Verma and J. Pearl. Causal networks: Semantics and expressiveness. In R. D. Shachter, T. S. Levitt, L. N. Kanal, and J. F. Lemmer, editors, *Uncertainty in Artificial Intelligence*, volume 9 of *Machine Intelligence and Pattern Recognition*, pages 69–76. North-Holland, 1990a.

T. Verma and J. Pearl. Equivalence and synthesis of causal models. In P. P. Bonissone, M. Henrion, L. N. Kanal, and J. F. Lemmer, editors, *Proceedings of the Sixth Annual Conference on Uncertainty in Artificial Intelligence*, Uai '90, page 255–270, New York, NY, USA, 1990b. Elsevier Science Inc. ISBN 0444892648.

P. Versteeg, J. Mooij, and C. Zhang. Local constraint-based causal discovery under selection bias. In B. Schölkopf, C. Uhler, and K. Zhang, editors, *Proceedings of the First Conference on Causal Learning and Reasoning*, volume 177 of *Proceedings of Machine Learning Research*, pages 840–860. Pmlr, 11–13 Apr 2022.

R. Vinuesa and B. Sirmacek. Interpretable deep-learning models to help achieve the Sustainable Development Goals. *Nat. Mach. Intell.*, 3:926, 2021.

R. Vinuesa, P. Schlatter, and H. M. Nagib. Secondary flow in turbulent ducts with increasing aspect ratio. *Phys. Rev. Fluids*, 3:054606, 2018.

H. Von Storch and F. W. Zwiers. *Statistical analysis in climate research*. Cambridge university press, 2002.

A. Wagner. Causality in complex systems. *Biology and Philosophy*, 14:83–101, 1999.

F. Waleffe. Hydrodynamic stability and turbulence: Beyond transients to a self-sustaining process. *Stud. Appl. Maths*, 95:319–343, 1995.

G. T. Walker. Correlation in seasonal variations of weather, VIII: A preliminary study of world weather. *Mem. Indian Meteorol. Dep.*, 24(4):75–131, 1923.

J. M. Wallace, H. Eckelman, and R. S. Brodkey. The wall region in turbulent shear flow. *J. Fluid Mech.*, 54:39–48, 1972.

D. Waltz and B. G. Buchanan. Automating science. *Science*, 324(5923):43–44, 2009.

Z. Wang, I. Akhtar, J. Borggaard, and T. Iliescu. Proper orthogonal decomposition closure models for turbulent flows: A numerical comparison. *Comput. Methods Appl. Mech. Eng.*, 237:10–26, 2012.

A. Warne. Causality and regime inference in a Markov switching VAR. Technical report, Sveriges Riksbank Working Paper Series, 2000.

T. Washio and H. Motoda. Inductive inference of first-order rules with non-linear structures. *Machine Learning*, 27(2):153–172, 1997.

J. D. Watson and F. H. Crick. Molecular structure of nucleic acids: A structure for deoxyribose nucleic acid. *Nature*, 171(4356), 1953.

J. Weatheritt and R. D. Sandberg. A novel evolutionary algorithm applied to algebraic modifications of the RANS stress-strain relationship. *J. Comput. Phys.*, 325:22–37, 2016.

J. Weatheritt and R. D. Sandberg. The development of algebraic stress models using a novel evolutionary algorithm. *Int. J. Heat Fluid Flow*, 68:298–318, 2017.

S. Weichwald and J. Peters. Causality in cognitive neuroscience: concepts, challenges, and distributional robustness. *Journal of Cognitive Neuroscience*, 33(2):226–247, 2021.

S. Weichwald, T. Meyer, O. Özdenizci, B. Schölkopf, T. Ball, and M. Grosse-Wentrup. Causal interpretation rules for encoding and decoding models in neuroimaging. *NeuroImage*, 110:48–59, 2015.

M. Werner, A. Junginger, P. Hennig, and G. Martius. Informed equation learning, 2021.

M. Werner, A. Junginger, P. Hennig, and G. Martius. Uncertainty in equation learning. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion (GECCO)*, pages 2298–2305. Association for Computing Machinery, 2022.

M. O. Williams, I. G. Kevrekidis, and C. W. Rowley. A data-driven approximation of the Koopman operator: Extending dynamic mode decomposition. *Journal of Nonlinear Science*, 25:1307–1346, 2015a.

M. O. Williams, C. W. Rowley, and I. G. Kevrekidis. A kernel-based method for data-driven Koopman spectral analysis. *Journal of Computational Dynamics*, 2:247–265, 2015b.

H. R. Wilson and J. D. Cowan. Excitatory and inhibitory interactions in localized populations of model neurons. *Biophysical Journal*, 12:1–24, 1972.

H. R. Wilson and J. D. Cowan. Evolution of the Wilson-Cowan equations. *Biol. Cybern.*, 115(6): 643–653, Dec. 2021.

W. C. Wimsatt and W. K. Wimsatt. *Re-engineering philosophy for limited beings: Piecewise approximations to reality*. Harvard University Press, 2007.

C. Winship and S. L. Morgan. The estimation of causal effects from observational data. *Annual review of sociology*, pages 659–706, 1999.

E. Winter. The shapley value. *Handbook of Game Theory with Economic Applications*, 3:2025–2054, 2002.

A. Woolgar, P. Golland, and S. Bode. Coping with confounds in multivoxel pattern analysis: What should we do about reaction time differences? A comment on Todd, Nystrom & Cohen 2013. *NeuroImage*, 98:506–512, 2014.

D. Wootton. *The invention of science: A new history of the scientific revolution*. Penguin UK, 2015.

Z. Wu, S. L. Brunton, and S. Revzen. Challenges in dynamic mode decomposition. *Journal of the Royal Society Interface*, 18(185):20210686, 2021.

H. Ye, E. R. Deyle, L. J. Gilarranz, and G. Sugihara. Distinguishing time-delayed causal interactions using convergent cross mapping. *Scientific reports*, 5(1):1–9, 2015.

M. Z. Yousif, M. Zhang, L. Yu, R. Vinuesa, and H. Lim. A transformer-based synthetic-inflow generator for spatially-developing turbulent boundary layers. *J. Fluid Mech., To Appear. Preprint arXiv:2206.01618*, 2022.

B. M. Yu, J. P. Cunningham, G. Santhanam, S. I. Ryu, K. V. Shenoy, and M. Sahani. Gaussian-process factor analysis for low-dimensional single-trial analysis of neural population activity. *Journal of Neurophysiology*, 102(1):614–635, 2009. Pmid: 19357332.

S. Yu, M. Drton, and A. Shojaie. Directed graphical models and causal discovery for zero-inflated data. *arXiv preprint arXiv:2004.04150*, 2020.

Y. Yu, T. Gao, N. Yin, and Q. Ji. DAGs with No Curl: An Efficient DAG Structure Learning Approach. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 12156–12166. PMLR, 18–24 Jul 2021.

L. Zanna and T. Bolton. Data-driven equation discovery of ocean mesoscale closures. *Geophysical Research Letters*, 47(17):e2020GL088376, 2020.

J. Zhang. Causal reasoning with ancestral graphs. *Journal of Machine Learning Research*, 9(47):1437–1474, 2008a.

J. Zhang. On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artificial Intelligence*, 172(16):1873–1896, 2008b.

J. Zhang and P. Spirtes. Strong faithfulness and uniform consistency in causal inference. In *Proceedings of the Nineteenth Conference on Uncertainty in Artificial Intelligence*, Uai'03, page 632–639, San Francisco, CA, USA, 2002. Morgan Kaufmann Publishers Inc. ISBN 0127056645.

J. Zhang and P. Spirtes. Detection of unfaithfulness and robust causal inference. *Minds and Machines*, 18(2):239–271, 2008.

X. Zheng, B. Aragam, P. K. Ravikumar, and E. P. Xing. DAGs with NO TEARS: Continuous optimization for structure learning. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

X. Zheng, C. Dan, B. Aragam, P. Ravikumar, and E. Xing. Learning sparse nonparametric DAGs. In S. Chiappa and R. Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 3414–3425. Pmlr, 26–28 Aug 2020.

H. Zou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *J. Comput. Graph. Stat.*, 15(1):265–286, 2012.

J. Zscheischler, S. Westra, B. J. Van Den Hurk, S. I. Seneviratne, P. J. Ward, A. Pitman, A. AghaKouchak, D. N. Bresch, M. Leonard, T. Wahl, et al. Future climate risk from compound events. *Nature Climate Change*, 8(6):469–477, 2018.

J. Zscheischler, O. Martius, S. Westra, E. Bevacqua, C. Raymond, R. M. Horton, B. van den Hurk, A. AghaKouchak, A. Jézéquel, M. D. Mahecha, et al. A typology of compound weather and climate events. *Nature reviews earth & environment*, 1(7):333–347, 2020.

J. Żytkow, R. Michalski, and R. Stepp. Representation and learning of categorical structures. *Machine Learning*, 5(1):7–48, 1990.