# FOUNDATIONS OF STRUCTURAL CAUSAL MODELS WITH CYCLES AND LATENT VARIABLES

BY STEPHAN BONGERS[1,*], PATRICK FORRÉ[1,†], JONAS PETERS[2,‡], BERNHARD SCHÖLKOPF[3] AND JORIS M. MOOIJ[4,*,†]

[1]*Informatics Institute, University of Amsterdam, s.r.bongers@uva.nl; p.d.forre@uva.nl*

[2]*Department of Mathematical Sciences, University of Copenhagen, jonas.peters@math.ku.dk*

[3]*Empirical Inference, Max Planck Institute for Intelligent Systems, bernhard.schoelkopf@tuebingen.mpg.de*

[4]*Korteweg-De Vries Institute, University of Amsterdam, j.m.mooij@uva.nl*

Structural causal models (SCMs), also known as (non-parametric) structural equation models (SEMs), are widely used for causal modeling purposes. In particular, acyclic SCMs, also known as recursive SEMs, form a well-studied subclass of SCMs that generalize causal Bayesian networks to allow for latent confounders. In this paper, we investigate SCMs in a more general setting, allowing for the presence of both latent confounders and cycles. We show that in the presence of cycles, many of the convenient properties of acyclic SCMs do not hold in general: they do not always have a solution; they do not always induce unique observational, interventional and counterfactual distributions; a marginalization does not always exist, and if it exists the marginal model does not always respect the latent projection; they do not always satisfy a Markov property; and their graphs are not always consistent with their causal semantics. We prove that for SCMs in general each of these properties does hold under certain solvability conditions. Our work generalizes results for SCMs with cycles that were only known for certain special cases so far. We introduce the class of simple SCMs that extends the class of acyclic SCMs to the cyclic setting, while preserving many of the convenient properties of acyclic SCMs. With this paper we aim to provide the foundations for a general theory of statistical causal modeling with SCMs.

**1. Introduction** Structural causal models (SCMs), also known as (non-parametric) structural equation models (SEMs), are widely used for causal modeling purposes [5, 73, 51, 55]. They form the basis for many statistical methods that aim at inferring knowledge of the underlying causal structure from data [see e.g., 37, 45, 56, 7, 48]. In these models, the causal relationships between the variables are expressed in the form of deterministic, functional relationships, and probabilities are introduced through the assumption that certain variables are exogenous latent random variables. SCMs arose out of certain causal models that were first introduced in genetics [79], econometrics [25], electrical engineering [39, 40], and the social sciences [23, 12].

Acyclic SCMs, also known as recursive SEMs, form a special well-studied subclass of SCMs that generalize causal Bayesian networks [51]. They have many convenient properties [see e.g., 50, 35, 78, 34, 60, 15, 16]: (i) they induce a unique distribution over the variables; (ii) they are closed under perfect interventions; (iii) they are closed under marginalizations; (iv) their marginalization respects the latent projection; (v) they obey (various equiva-

lent versions of) the Markov property; and (vi) their graphs express the causal relationships encoded by the SCM in an intuitive manner.

One important limitation of acyclic SCMs is that they cannot model systems that involve causal cycles. In many systems occurring in the real world, there are feedback loops between observed variables. For example, in economics the price of a product may be a function of the demanded or supplied quantities, and vice versa, the demanded and supplied quantities may be functions of the price. The underlying dynamic processes describing such systems have an acyclic causal structure over time. However, causal cycles may arise when one approximates such systems over time [17, 43, 42] or when one describes the equilibrium states of these systems [29, 33, 27, 46, 6, 3, 57]. In particular, in [6] it was shown that the equilibrium states of a system governed by (random) differential equations can be described by an SCM that represents their causal semantics, which gives rise to a plethora of SCMs that include cycles (we provide some examples of such feedback systems in Appendix D.1 in the Supplementary Material). In contrast to their acyclic counterparts, SCMs with cycles have enjoyed less attention in the literature and are not as well understood. In general, none of the above properties (i)–(vi) hold in the class of SCMs. However, some progress has been made in the case of discrete [52, 49] and linear models [70, 71, 72, 63, 31, 27], and more recently, for more general cyclic models the Markov properties have been elucidated [18].

*Contributions*   The purpose of this paper is to provide the foundations for a general theory of statistical causal modeling with SCMs. We study properties of SCMs and allow for cycles, latent variables and non-linear functional relationships between the variables. We investigate to which extent and under which sufficient conditions each of the properties (i)–(vi) holds, in particular, in the presence of cycles. In the next paragraphs, we describe our contributions in more detail.

When there are cyclic functional relationships between variables, one encounters various technical complications, which even arise in the linear setting. The structural equations of an acyclic SCM trivially have a unique solution. This unique solvability property ensures that the SCM gives rise to a unique, well-defined probability distribution on the variables. In the case of cycles, however, this property may be violated, and consequently, the SCM may not have a solution at all, or may allow for multiple different probability distributions [26]. Even if one starts with a cyclic SCM that is uniquely solvable, performing an intervention on the SCM may lead to an intervened SCM that is not uniquely solvable. Hence, a cyclic SCM may not give rise to a unique, well-defined probability distribution corresponding to that intervention, and whether or not this happens may depend on the intervention. We provide sufficient conditions for the existence and uniqueness of these probability distributions after intervention. In general, it is not clear whether the solutions of the structural equations of an SCM are measurable if cycles are present. In addition, we provide sufficient and necessary conditions for the measurability of solution functions of cyclic SCMs.

SCMs provide a detailed modeling description of a system. Not all information may be necessary for a certain modeling task, which motivates to consider certain classes of SCMs to be equivalent. In this paper, we formally introduce several of such equivalence relations. For example, we consider two SCMs observationally equivalent if they cannot be distinguished based on observations alone. Observationally equivalent SCMs can often still be distinguished by interventions. We consider two SCMs interventionally equivalent if they cannot be distinguished based on observations and interventions. While these concepts have been around in implicit form for acyclic SCMs, we formulate them in such a way that they also apply to cyclic SCMs that have either no solution at all or have multiple different induced probability distributions on the variables. Finally, we consider two SCMs counterfactually equivalent if they cannot be distinguished based on observations and interventions and in addition encode the same counterfactual distributions, which are the distributions induced

by the so-called twin SCM via the twin network method [1]. These different equivalence relations formalize the different levels of abstraction in the so-called causal hierarchy [69, 53]. In addition, we add another, strong version of equivalence, such that equivalent SCMs have the same solutions. This notion clarifies ambiguities when a function is constant in one of its arguments, for example.

Marginalization becomes useful if not all variables are observed: given a joint probability distribution on some variables, we obtain a marginal distribution on a subset of the variables by integrating out the remaining variables. Analogously, we can marginalize an acyclic SCM by substituting the solutions of the structural equations of a subset of the endogenous variables into the structural equations of the remaining endogenous variables. For acyclic SCMs, the induced observational and interventional distributions of the marginalized SCM coincide with the marginals of the distributions induced by the original SCM [see 78, 75, 15, 16, a.o.]. In other words, for acyclic SCMs the operation of marginalization preserves the probabilistic and causal semantics (restricted to the remaining variables). We show that for cyclic SCMs a marginalization does not always exist without further assumptions. In [18] it is shown that for modular SCMs, which can be seen as an SCM together with an additional structure of a compatible system of solution functions, a marginalization can be defined that preserves the probabilistic and causal semantics. We prove that this additional structure is not necessary and use a local unique solvability condition instead. Under this condition, we show that an SCM and its marginalization are observationally, interventionally and counterfactually equivalent on the remaining endogenous variables. Analogously, we define a marginalization operation on the associated graph of an SCM, which generalizes the latent projection [78, 76, 15]. In general, the marginalization of an SCM does not respect the latent projection of its associated graph, but we show that it does so under an additional local ancestral unique solvability condition.

In graphical models, Markov properties allow one to read off conditional independencies in a distribution directly from a graph. Various equivalent formulations of Markov properties exist for acyclic SCMs [34], one prominent example being the $d$-separation criterion, also known as the directed global Markov property, which was originally derived for Bayesian networks [50]. Markov properties have been of key importance to derive various central results regarding causal reasoning and causal discovery. For cyclic SCMs, however, the usual Markov properties do not hold in general, as was already pointed out by Spirtes [71]. His solution in terms of collapsed graphs was recently generalized and reformulated for a general class of causal graphical models [18] by adapting the notion of $d$-separation into what has been termed $\sigma$-separation. This resulted in a general directed global Markov property expressed in terms of $\sigma$-separation instead of $d$-separation. Here, we formulate these general Markov properties specifically within the framework of SCMs. Again, they only hold under certain unique solvability conditions.

In addition to its interpretation in terms of conditional independencies, the graph of an acyclic SCM also has a direct causal interpretation [51]. As was already observed in [49], the causal interpretation of SCMs with cycles can be counterintuitive, as the causal semantics under interventions no longer needs to be compatible with the structure imposed by the functional relations between the variables. We resolve this issue by showing that under certain ancestral unique solvability conditions the causal interpretation of SCMs is consistent with its graph.

Cycles lead to several technical complications related to solvability issues. We introduce a special subclass of (possibly cyclic) SCMs, the class of simple SCMs, for which most of these technical complications are absent and which preserves much of the simplicity of the theory for acyclic SCMs. A simple SCM is an SCM that is uniquely solvable with respect to every subset of the variables. Because of this strong solvability assumption, simple SCMs

Fig 1: *Overview of the objects constructed from an SCM and the mappings between them. The numbers correspond to the definition, proposition or theorem of the corresponding object, mapping, or result. When an arrow is dashed, the relation only holds under non-trivial assumptions that can be found in the corresponding definition or theorem. The symbol "⊆" stands for the subgraph of a directed mixed graph (see Definition A.1 in the Supplementary Material) and the symbol "↻" denotes that the surrounding diagram commutes. Table 1 gives an overview of the commutativity results for each pair of mappings between the objects with the names in bold.*

have all the convenient properties (i)–(vi): they always have uniquely defined observational, interventional and counterfactual distributions; we can perform every perfect intervention and marginalization on them and the result is again a simple SCM; marginalization does respect the latent projection; they obey the general directed global Markov property, and for special cases (including the acyclic, linear and discrete case) they obey the (stronger) directed global Markov property; their graphs have a direct and intuitive causal interpretation.

The scope of this paper is limited to establishing the foundations for statistical causal modeling with cyclic SCMs (Figure 9 in Appendix A.5 in the Supplementary Material shows an overview of how SCMs relate to other causal graphical models). For a detailed discussion of causal reasoning, causal discovery and causal prediction with cyclic SCMs we refer the reader to other literature [e.g., 58, 59, 61, 14, 27, 28, 21]. Several recent results (generalizations of the do-calculus, adjustment criteria and an identification algorithm) for modular SCMs [19, 20] directly apply to the subclass of simple SCMs, as well. Finally, many causal discovery algorithms that have been designed for the acyclic case also apply to simple SCMs with no or only minor changes [47, 44].

*Overview*   Figure 1 gives an overview of the different objects that can be constructed from an SCM and the different mappings between them. For pairs of mappings between the objects with the names in bold we prove commutativity results which are summarized in Table 1.

*Outline*   This paper is structured as follows: In Section 2, we provide a formal definition of SCMs and a natural notion of equivalence between SCMs, define the (augmented) graph corresponding to an SCM, and describe perfect interventions and counterfactuals. In Section 3, we discuss the concept of (unique) solvability, its properties and how it relates to self-cycles. In Section 4, we define and relate various equivalence relations between SCMs. In Section 5, we define a marginalization operation that is applicable to cyclic SCMs under certain conditions. We discuss several properties of this marginalization operation and discuss the relation with a marginalization operation defined on directed mixed graphs. In Section 6, we discuss Markov properties of SCMs. In Section 7, we discuss the causal interpretation of the graphs of SCMs. Section 8 introduces and discusses the class of simple SCMs.

| **SCMs** | do | twin | marg |
|---|---|---|---|
| $\mathcal{G}, \mathcal{G}^a$ | 2.15 | 2.20 | (5.12) |
| do | 2.16.(1) | 2.22.(1) | 5.5.(1) |
| twin | $\cdots$ | - | 5.5.(2) |
| marg | $\cdots$ | $\cdots$ | 5.4 |

| **Graphs** | do | twin | marg |
|---|---|---|---|
| do | 2.16.(1) | 2.22.(2) | 5.10.(1) |
| twin | $\cdots$ | - | 5.10.(2) |
| marg | $\cdots$ | $\cdots$ | 5.9 |

Table 1: *Overview of the commutativity results of different pairs of mappings, defined on SCMs (left table) and on graphs (right table). All results apply under the assumptions stated in the corresponding proposition. The entries denoted by dots are omitted due to symmetry. We do not consider the commutativity of the twin operation with itself in this paper. Proposition 5.12 (in parentheses) is not a commutativity result but a weaker relation. The graphical twin operator is only defined for directed graphs.*

The Supplementary Material introduces causal graphical models in Appendix A. This section also contains details on Markov properties and modular SCMs. Appendix B provides additional (unique) solvability properties, some results for linear SCMs are discussed in Appendix C, other examples in Appendix D, and the proofs of all the theoretical results are in Appendix E. Appendix F contains some lemmas and measurable selection theorems that are used in several proofs.

**2. Structural causal models**   In this section, we provide the definition and properties of structural causal models (SCMs). Our definition of SCMs slightly deviates from existing definitions [5, 51, 73], because we make the definition of the SCM independent of the random variables that solve it. This enables us to deal with the various technical complications that arise in the presence of cycles.

2.1. *Structural causal models and their solutions*

DEFINITION 2.1 (Structural causal model).    *A structural causal model (SCM) is a tuple*[1]

$$\mathcal{M} := \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle,$$

*where*

1. $\mathcal{I}$ *is a finite index set of* endogenous variables,
2. $\mathcal{J}$ *is a disjoint finite index set of* exogenous variables,
3. $\boldsymbol{\mathcal{X}} = \prod_{i \in \mathcal{I}} \mathcal{X}_i$ *is the product of the domains of the endogenous variables, where each domain $\mathcal{X}_i$ is a standard measurable space (see Definition F.1),*
4. $\boldsymbol{\mathcal{E}} = \prod_{j \in \mathcal{J}} \mathcal{E}_j$ *is the product of the domains of the exogenous variables, where each domain $\mathcal{E}_j$ is a standard measurable space,*
5. $\boldsymbol{f} : \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{E}} \to \boldsymbol{\mathcal{X}}$ *is a measurable function that specifies the* causal mechanism,
6. $\mathbb{P}_{\boldsymbol{\mathcal{E}}} = \prod_{j \in \mathcal{J}} \mathbb{P}_{\mathcal{E}_j}$ *is a product measure, the* exogenous distribution, *where $\mathbb{P}_{\mathcal{E}_j}$ is a probability measure on $\mathcal{E}_j$ for each $j \in \mathcal{J}$.*[2]

In SCMs, the functional relationships between variables are expressed in terms of deterministic equations, where each equation expresses an endogenous variable (on the left-hand side) in terms of a causal mechanism depending on endogenous and exogenous variables (on the right-hand side). This allows us to model interventions in an unambiguous way by changing the causal mechanisms that target specific endogenous variables (see Section 2.4).

---

[1]We often use boldface for variables that have multiple components, e.g., vectors in a Cartesian product.

[2]For the case $\mathcal{J} = \emptyset$ we have that $\boldsymbol{\mathcal{E}}$ is the singleton $\mathbf{1}$ and $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$ is the degenerate probability measure $\mathbb{P}_{\mathbf{1}}$.

DEFINITION 2.2 (Structural equations).   *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM. We call the set of equations*

$$x_i = f_i(\boldsymbol{x}, \boldsymbol{e}) \qquad \boldsymbol{x} \in \boldsymbol{\mathcal{X}}, \boldsymbol{e} \in \boldsymbol{\mathcal{E}}$$

*for $i \in \mathcal{I}$ the* structural equations *of the structural causal model $\mathcal{M}$.*

Although it is common to assume the absence of cyclic functional relations (see Definition 2.10), we make no such assumption here. In particular, we allow for self-cycles, which we will discuss in more detail in Section 2.2 and 3.3.

The solutions of an SCM in terms of random variables are defined up to almost sure equality. Random variables that are almost surely equal are generally considered to be equivalent to each other for all practical purposes.

DEFINITION 2.3 (Solution).   *A pair $(\boldsymbol{X}, \boldsymbol{E})$ of random variables $\boldsymbol{X} : \Omega \to \boldsymbol{\mathcal{X}}, \boldsymbol{E} : \Omega \to \boldsymbol{\mathcal{E}}$, where $\Omega$ is a probability space, is a* solution *of the SCM $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ if*

1. *$\mathbb{P}^{\boldsymbol{E}} = \mathbb{P}_{\boldsymbol{\mathcal{E}}}$, i.e., the distribution of $\boldsymbol{E}$ is equal to $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$,[3] and*
2. *the* structural equations *are satisfied, i.e.,*

$$\boldsymbol{X} = \boldsymbol{f}(\boldsymbol{X}, \boldsymbol{E}) \text{ a.s..}$$

*For convenience, we call a random variable $\boldsymbol{X}$ a* solution *of $\mathcal{M}$ if there exists a random variable $\boldsymbol{E}$ such that $(\boldsymbol{X}, \boldsymbol{E})$ forms a solution of $\mathcal{M}$.*

Often, the endogenous random variables $\boldsymbol{X}$ can be observed, while the exogenous random variables $\boldsymbol{E}$ are treated as latent. Latent exogenous variables are often referred to as "disturbance terms" or "noise variables". For a solution $\boldsymbol{X}$, we call the distribution $\mathbb{P}^{\boldsymbol{X}}$ the *observational distribution of $\mathcal{M}$ associated to $\boldsymbol{X}$.* In general there may be multiple different observational distributions associated to an SCM due to the existence of different solutions of the structural equations. This is a consequence of the allowance of cycles in SCMs, as the following simple example illustrates.

EXAMPLE 2.4 (Cyclic SCMs).   *For brevity, we use throughout this paper the notation $\boldsymbol{n} := \{1, 2, \ldots, n\}$ for $n \in \mathbb{N}$. Let $\mathcal{M} = \langle \boldsymbol{2}, \boldsymbol{1}, \mathbb{R}^2, \mathbb{R}, \boldsymbol{f}, \mathbb{P}_{\mathbb{R}} \rangle$ be an SCM[4] with $f_1(\boldsymbol{x}, e) = x_2$ and $f_2(\boldsymbol{x}, e) = x_1$, and $\mathbb{P}_{\mathbb{R}}$ an arbitrary probability measure on $\mathbb{R}$. Then $(X, X)$ is a solution of $\mathcal{M}$ for any arbitrary random variable $X$ with values in $\mathbb{R}$. Hence, any probability distribution on $\{(x, x) : x \in \mathbb{R}\}$ is an observational distribution associated to $\mathcal{M}$. Now consider instead the same SCM but with $f_1(\boldsymbol{x}, e) = x_2 + 1$. This SCM has no solutions at all, and hence induces no observational distribution.*

Due to the fact that the structural equations only need to be satisfied almost surely, there may exist many different SCMs representing the same set of solutions.

---

[3]This implies that the components $E_j$ of $\boldsymbol{E}$ are mutually independent, since $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$ is a product measure on $\prod_{j \in \mathcal{J}} \mathcal{E}_j$.

[4]In our examples, we will abuse notation by using non-disjoint subsets of the natural numbers to index both endogenous and exogenous variables; these should be understood to be disjoint copies of the natural numbers: if we write $\mathcal{I} = \boldsymbol{n}$ and $\mathcal{J} = \boldsymbol{m}$, we mean instead $\mathcal{I} = \{1, 2, \ldots, n\}$ and $\mathcal{J} = \{1', 2', \ldots, m'\}$ where $k'$ is a copy of $k$.

EXAMPLE 2.5 (Structural equations up to almost sure equality).   *Consider the SCM $\mathcal{M} = \langle \mathbf{1}, \mathbf{1}, \mathcal{X}, \mathcal{E}, f, \mathbb{P}_{\mathcal{E}} \rangle$ with $\mathcal{X} = \mathcal{E} = \{-1, 0, 1\}$, $\mathbb{P}_{\mathcal{E}}(\{-1\}) = \mathbb{P}_{\mathcal{E}}(\{1\}) = \frac{1}{2}$ and $f(x, e) = e^2 + e - 1$. Let $\tilde{\mathcal{M}}$ be the SCM $\mathcal{M}$ but with a different causal mechanism $\tilde{f}(x, e) = e$. Then the set of solutions of the structural equations agree for both SCMs for $e \in \{-1, +1\}$, while they differ only for $e = 0$, which occurs with probability zero. Hence, a pair of random variables $(X, E)$ is a solution of $\mathcal{M}$ if and only if it is a solution of $\tilde{\mathcal{M}}$.*

It therefore seems natural not to differentiate between structural equations that have different solutions on at most a $\mathbb{P}_{\mathcal{E}}$-null set of exogenous variables. This leads to the following equivalence relation between SCMs. To be able to state the equivalence relation concisely, we introduce the following notation: For subsets $\mathcal{U} \subseteq \mathcal{I}$ and $\mathcal{V} \subseteq \mathcal{J}$ we write $\boldsymbol{\mathcal{X}}_{\mathcal{U}} := \prod_{i \in \mathcal{U}} \mathcal{X}_i$ and $\boldsymbol{\mathcal{E}}_{\mathcal{V}} := \prod_{j \in \mathcal{V}} \mathcal{E}_j$. In particular, $\boldsymbol{\mathcal{X}}_{\emptyset}$ and $\boldsymbol{\mathcal{E}}_{\emptyset}$ are defined by the singleton $\mathbf{1}$. Moreover, for a subset $\mathcal{W} \subseteq \mathcal{I} \cup \mathcal{J}$, we use the convention that we write $\boldsymbol{\mathcal{X}}_{\mathcal{W}}$ and $\boldsymbol{\mathcal{E}}_{\mathcal{W}}$ instead of $\boldsymbol{\mathcal{X}}_{\mathcal{W} \cap \mathcal{I}}$ and $\boldsymbol{\mathcal{E}}_{\mathcal{W} \cap \mathcal{J}}$ respectively and we adopt a similar notation for the (random) variables in those spaces, that is, we write $\boldsymbol{x}_{\mathcal{W}}$ and $\boldsymbol{e}_{\mathcal{W}}$ instead of $\boldsymbol{x}_{\mathcal{W} \cap \mathcal{I}}$ and $\boldsymbol{e}_{\mathcal{W} \cap \mathcal{J}}$ respectively. This allows us to define the following natural equivalence relation for SCMs.[5,6]

DEFINITION 2.6 (Equivalence).   *The two SCMs $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ and $\tilde{\mathcal{M}} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \tilde{\boldsymbol{f}}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ are* equivalent, *denoted by $\mathcal{M} \equiv \tilde{\mathcal{M}}$, if for all $i \in I$, for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$*

$$x_i = f_i(\boldsymbol{x}, \boldsymbol{e}) \quad \Longleftrightarrow \quad x_i = \tilde{f}_i(\boldsymbol{x}, \boldsymbol{e}).$$

Thus, two equivalent SCMs can only differ in terms of their causal mechanism. Importantly, equivalent SCMs have the same solutions and, as we will see in Section 2.4 and 2.5, they have the same causal and counterfactual semantics (see Definition 2.13 and 2.18 respectively). This equivalence relation on the set of all SCMs gives rise to the quotient set of equivalence classes of SCMs. In this paper we prove properties and define operations on the equivalence classes of SCMs, by first proving the property and defining the operation for an SCM and then showing that this property and operation preserves the equivalence relation.

2.2. *The (augmented) graph*   We will now define two types of graphs that can be used for representing structural properties of the SCM. These graphical representations are related to Wright's path diagrams [79]. The structural properties of the functional relations between variables modeled by an SCM are specified by the causal mechanism of the SCM and can be encoded in an (augmented) graph. For the graphical notation and standard terminology on directed (mixed) graphs that is used throughout this paper, we refer the reader to Appendix A.1.
We first define the parents of an endogenous variable.

---

[5]An attempt at coarsening this notion of equivalence by replacing the quantifier "for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$" by "for almost every $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$ under the observational distribution $\mathbb{P}^{\boldsymbol{X}}$" will not lead to a well-defined equivalence relation, since in general the observational distribution $\mathbb{P}^{\boldsymbol{X}}$ may be non-unique or even non-existent. Refining it by replacing the quantifier "for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$" by "for all $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$" would make it too fine for our purposes, since we assume the exogenous distribution to be fixed and we assume as usual that random variables that are almost surely identical are indistinguishable in practice.

[6]We may extend this definition to allow $\tilde{\mathcal{J}} \neq \mathcal{J}$ and for a larger class of SCMs such that the exogenous distribution does not factorize. Then, for any $\mathcal{M}$ that satisfies Definition 2.1, except for that it may have a non-factorizing exogenous distribution, there exists an equivalent SCM with a factorizing exogenous distribution (and a different $\mathcal{J}$); the latter can be obtained by partitioning the exogenous components into independent tuples. This motivates why we can restrict ourselves in Definition 2.1 to factorizing exogenous distributions only.

DEFINITION 2.7 (Parent). *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM. We call $k \in \mathcal{I} \cup \mathcal{J}$ a parent of $i \in \mathcal{I}$ if and only if there does not exist a measurable function[7] $\tilde{f}_i : \boldsymbol{\mathcal{X}}_{\backslash k} \times \boldsymbol{\mathcal{E}}_{\backslash k} \to \mathcal{X}_i$ such that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$*

$$x_i = f_i(\boldsymbol{x}, \boldsymbol{e}) \quad \Longleftrightarrow \quad x_i = \tilde{f}_i(\boldsymbol{x}_{\backslash k}, \boldsymbol{e}_{\backslash k}).$$

Exogenous variables have no parents by definition. These parental relations are preserved under the equivalence relation $\equiv$ on SCMs. They can be represented by a directed graph or a directed mixed graph.[8]

DEFINITION 2.8 (Graph and augmented graph). *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM. We define:*

1. *the* augmented graph $\mathcal{G}^a(\mathcal{M})$ *as the directed graph with nodes $\mathcal{I} \cup \mathcal{J}$ and directed edges $u \to v$ if and only if $u \in \mathcal{I} \cup \mathcal{J}$ is a parent of $v \in \mathcal{I}$;*
2. *the* graph $\mathcal{G}(\mathcal{M})$ *as the directed mixed graph with nodes $\mathcal{I}$, directed edges $u \to v$ if and only if $u \in \mathcal{I}$ is a parent of $v \in \mathcal{I}$ and bidirected edges $u \leftrightarrow v$ if and only if there exists a $j \in \mathcal{J}$ that is a parent of both $u \in \mathcal{I}$ and $v \in \mathcal{I}$.*

*We call the mappings $\mathcal{G}^a$ and $\mathcal{G}$, that map $\mathcal{M}$ to $\mathcal{G}^a(\mathcal{M})$ and $\mathcal{G}(\mathcal{M})$, the* augmented graph mapping *and the* graph mapping *respectively.*

In particular, the augmented graph contains no directed edges pointing towards an exogenous variable, i.e., $u \in \mathcal{I} \cup \mathcal{J}$ cannot be a parent of $v \in \mathcal{J}$, because they are not functionally related through the causal mechanism. We call a directed edge $i \to i$ in $\mathcal{G}^a(\mathcal{M})$ and $\mathcal{G}(\mathcal{M})$ (here, $i$ is a parent of itself) a *self-cycle* at $i$. By definition, the mappings $\mathcal{G}^a$ and $\mathcal{G}$ are invariant under the equivalence relation $\equiv$ on SCMs and hence the equivalence class of an SCM $\mathcal{M}$ is mapped to a unique augmented graph $\mathcal{G}^a(\mathcal{M})$ and a unique graph $\mathcal{G}(\mathcal{M})$.

EXAMPLE 2.9 (Graphs of an SCM). *Let $\mathcal{M} = \langle \boldsymbol{5}, \boldsymbol{3}, \mathbb{R}^5, \mathbb{R}^3, \boldsymbol{f}, \mathbb{P}_{\mathbb{R}^3} \rangle$ be an SCM with causal mechanism given by*

$$f_1(\boldsymbol{x}, \boldsymbol{e}) = x_1 - x_1^2 + \alpha e_1^2, \qquad f_3(\boldsymbol{x}, \boldsymbol{e}) = -x_4 + e_2, \quad f_5(\boldsymbol{x}, \boldsymbol{e}) = x_4 \cdot e_3,$$

$$f_2(\boldsymbol{x}, \boldsymbol{e}) = x_1 + x_3 + x_4 + e_1, \quad f_4(\boldsymbol{x}, \boldsymbol{e}) = x_2 + e_2,$$

*where $\alpha \neq 0$ and $\mathbb{P}_{\mathbb{R}^3}$ is a product of three probability measures $\mathbb{P}_{\mathbb{R}}$ over $\mathbb{R}$ that are non-degenerate. The augmented graph $\mathcal{G}^a(\mathcal{M})$ and the graph $\mathcal{G}(\mathcal{M})$ of $\mathcal{M}$ are depicted[9] in Figure 2 (left). Observe that if $\alpha$ had been equal to zero, then the endogenous variable 1 would not have any parents in $\mathcal{G}^a(\mathcal{M})$, i.e., it would not have a self-cycle and directed edge from any exogenous variables in $\mathcal{G}^a(\mathcal{M})$, and it would not have a self-cycle and bidirected edge from any other variable in $\mathcal{G}(\mathcal{M})$. Moreover, if one of the probability measures $\mathbb{P}_{\mathbb{R}}$ over $\mathbb{R}$ were degenerate, then some of the directed edges from the exogenous variables to the endogenous variables in the augmented graph $\mathcal{G}^a(\mathcal{M})$ and bidirected edges in the graph $\mathcal{G}(\mathcal{M})$ would be missing.*

---

[7]For $\boldsymbol{\mathcal{X}} = \prod_{i \in \mathcal{I}} \mathcal{X}_i$, $\mathcal{I}$ some index set, $I \subseteq \mathcal{I}$ and $k \in \mathcal{I}$, we denote $\boldsymbol{\mathcal{X}}_{\backslash I} = \prod_{i \in \mathcal{I} \backslash I} \mathcal{X}_i$ and $\boldsymbol{\mathcal{X}}_{\backslash k} = \prod_{i \in \mathcal{I} \backslash \{k\}} \mathcal{X}_i$, and similarly for their elements.

[8]A *directed mixed graph* $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ consists of a set of nodes $\mathcal{V}$, a set of directed edges $\mathcal{E}$ and a set of bidirected edges $\mathcal{B}$ (see Definition A.1 for a more precise definition).

[9]For visualizing an (augmented) graph, we adapt the common convention of using random variables, with the index set as a subscript, instead of using the index set itself. With a slight abuse of notation, we still use the random variables notation in the (augmented) graph in the case that the SCM has no solution at all.
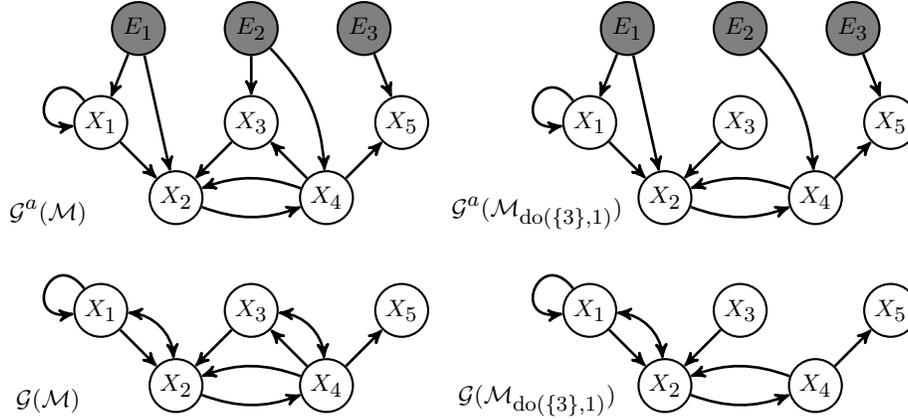
Fig 2: *The augmented graph (top) and the graph (bottom) of the SCM $\mathcal{M}$ of Example 2.9 (left) and of the intervened SCM $\mathcal{M}_{\mathrm{do}(\{3\},1)}$ of Example 2.17 (right).*

As is illustrated in this example, the augmented graph provides a more detailed representation than the graph. Therefore, we use the augmented graph as the standard graphical representation for SCMs, unless stated otherwise. For an SCM $\mathcal{M}$, we denote the sets $\mathrm{pa}_{\mathcal{G}^a(\mathcal{M})}(\mathcal{U})$, $\mathrm{ch}_{\mathcal{G}^a(\mathcal{M})}(\mathcal{U})$, $\mathrm{an}_{\mathcal{G}^a(\mathcal{M})}(\mathcal{U})$, etc., for some subset $\mathcal{U} \subseteq \mathcal{I} \cup \mathcal{J}$, by respectively $\mathrm{pa}(\mathcal{U})$, $\mathrm{ch}(\mathcal{U})$, $\mathrm{an}(\mathcal{U})$, etc., when the notation is clear from the context.

DEFINITION 2.10. *We call an SCM $\mathcal{M}$ acyclic if $\mathcal{G}^a(\mathcal{M})$ is a directed acyclic graph (DAG). Otherwise, we call $\mathcal{M}$ cyclic.*

Equivalently, an SCM $\mathcal{M}$ is acyclic if $\mathcal{G}(\mathcal{M})$ is an acyclic directed mixed graph (ADMG) [60]. Acyclic SCMs are also known as semi-Markovian SCMs [51, 76]. A commonly considered class of acyclic SCMs are the Markovian SCMs, which are acyclic SCMs for which each exogenous variable has at most one child. Several Markov properties were first shown for these models [51, 35, 76].

2.3. *Structurally minimal representations*　We have discussed an equivalence relation between SCMs in Section 2.1. In this subsection we show that for each SCM there exists a representative of the equivalence class of that SCM for which each component of the causal mechanism does not depend on its non-parents [see also 55].

DEFINITION 2.11 (Structurally minimal SCM). *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \mathcal{X}, \mathcal{E}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM. We call $\mathcal{M}$ structurally minimal if for all $i \in \mathcal{I}$ there exists a mapping $\tilde{f}_i : \mathcal{X}_{\mathrm{pa}(i)} \times \mathcal{E}_{\mathrm{pa}(i)} \to \mathcal{X}_i$ such that $f_i(\boldsymbol{x}, \boldsymbol{e}) = \tilde{f}_i(\boldsymbol{x}_{\mathrm{pa}(i)}, \boldsymbol{e}_{\mathrm{pa}(i)})$ for all $\boldsymbol{e} \in \mathcal{E}$ and all $\boldsymbol{x} \in \mathcal{X}$.*

We already encountered a structurally minimal SCM $\mathcal{M}$ in Example 2.9. Taking instead $\alpha = 0$ in that example gives an SCM $\mathcal{M}$ that is not structurally minimal, since the endogenous variable 1 is then not a parent of itself, while $f_1(\boldsymbol{x}, \boldsymbol{e})$ depends on $x_1$. However, the equivalent SCM where we have replaced the causal mechanism of 1 by $f_1(\boldsymbol{x}, \boldsymbol{e}) = 0$ yields a structurally minimal SCM. In general, there always exists an equivalent structurally minimal SCM.

PROPOSITION 2.12 (Existence of a structurally minimal SCM). *For an SCM $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \mathcal{X}, \mathcal{E}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$, there exists an equivalent SCM $\tilde{\mathcal{M}} = \langle \mathcal{I}, \mathcal{J}, \mathcal{X}, \mathcal{E}, \hat{\boldsymbol{f}}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ that is structurally minimal.*

For a causal mechanism $\boldsymbol{f} : \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{E}} \to \boldsymbol{\mathcal{X}}$ and a subset $\mathcal{U} \subseteq \mathcal{I}$, we write $\boldsymbol{f}_{\mathcal{U}} : \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{E}} \to \boldsymbol{\mathcal{X}}_{\mathcal{U}}$ for the $\mathcal{U}$ components[10] of $\boldsymbol{f}$. A structurally minimal representation is compatible with the (augmented) graph, in the sense that for every $\mathcal{U} \subseteq \mathcal{I}$ there exists a unique measurable mapping $\tilde{\boldsymbol{f}}_{\mathcal{U}} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{U})} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{U})} \to \boldsymbol{\mathcal{X}}_{\mathcal{U}}$ such that $\boldsymbol{f}_{\mathcal{U}}(\boldsymbol{x}, \boldsymbol{e}) = \tilde{\boldsymbol{f}}_{\mathcal{U}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{U})}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{U})})$ for all $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$. Moreover, for any $\mathcal{U} \subseteq \mathcal{I}$ there exists a unique measurable mapping $\tilde{\boldsymbol{f}}_{\mathrm{an}(\mathcal{U})} : \boldsymbol{\mathcal{X}}_{\mathrm{an}(\mathcal{U})} \times \boldsymbol{\mathcal{E}}_{\mathrm{an}(\mathcal{U})} \to \boldsymbol{\mathcal{X}}_{\mathrm{an}(\mathcal{U})}$ with $\boldsymbol{f}_{\mathrm{an}(\mathcal{U})}(\boldsymbol{x}, \boldsymbol{e}) = \tilde{\boldsymbol{f}}_{\mathcal{U}}(\boldsymbol{x}_{\mathrm{an}(\mathcal{U})}, \boldsymbol{e}_{\mathrm{an}(\mathcal{U})})$ for all $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$.

2.4. *Interventions*  To define the causal semantics of SCMs, we consider here an idealized class of interventions introduced by Pearl [51] that we refer to as perfect interventions. Other types of interventions, like mechanism changes [77], fat-hand interventions [13], activity interventions [45], and stochastic versions of all these are at least as relevant, but we do not consider them here.

DEFINITION 2.13 (Perfect intervention on an SCM).  *Let* $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ *be an SCM,* $I \subseteq \mathcal{I}$ *a subset of endogenous variables and* $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$ *a value. The* perfect intervention $\mathrm{do}(I, \boldsymbol{\xi}_I)$ *maps* $\mathcal{M}$ *to the SCM* $\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)} := \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \tilde{\boldsymbol{f}}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$, *where the* intervened causal mechanism $\tilde{\boldsymbol{f}}$ *is given by*

$$\tilde{f}_i(\boldsymbol{x}, \boldsymbol{e}) = \begin{cases} \xi_i & i \in I \\ f_i(\boldsymbol{x}, \boldsymbol{e}) & i \in \mathcal{I} \setminus I. \end{cases}$$

This operation $\mathrm{do}(I, \boldsymbol{\xi}_I)$ preserves the equivalence relation (see Definition 2.6) on the set of all SCMs and hence this mapping induces a well-defined mapping on the set of equivalence classes of SCMs. Previous work has considered interventions only on a specific subset of endogenous variables [67, 2, 3]. Instead, we assume that we can intervene on any subset of endogenous variables in the model.

We define an analogous operation $\mathrm{do}(I)$ on directed mixed graphs.

DEFINITION 2.14 (Perfect intervention on a directed mixed graph).  *Let* $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ *be a directed mixed graph and* $I \subseteq \mathcal{V}$ *a subset. The perfect intervention* $\mathrm{do}(I)$ *maps* $\mathcal{G}$ *to the directed mixed graph* $\mathrm{do}(I)(\mathcal{G}) := (\mathcal{V}, \tilde{\mathcal{E}}, \tilde{\mathcal{B}})$, *where* $\tilde{\mathcal{E}} = \mathcal{E} \setminus \{ v \to i : v \in \mathcal{V}, i \in I \}$ *and* $\tilde{\mathcal{B}} = \mathcal{B} \setminus \{ v \leftrightarrow i : v \in \mathcal{V}, i \in \mathcal{I} \}$.

This operation simply removes all incoming edges on the nodes in $I$. The two notions of intervention are compatible with the (augmented) graph mapping.

PROPOSITION 2.15.  *Let* $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ *be an SCM,* $I \subseteq \mathcal{I}$ *a subset of endogenous variables and* $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$ *a value. Then* $\big(\mathcal{G}^a \circ \mathrm{do}(I, \boldsymbol{\xi}_I)\big)(\mathcal{M}) = \big(\mathrm{do}(I) \circ \mathcal{G}^a\big)(\mathcal{M})$ *and* $\big(\mathcal{G} \circ \mathrm{do}(I, \boldsymbol{\xi}_I)\big)(\mathcal{M}) = \big(\mathrm{do}(I) \circ \mathcal{G}\big)(\mathcal{M})$.

The two notions of perfect intervention satisfy the following elementary properties.

PROPOSITION 2.16.  *For an SCM and a directed mixed graph we have the following properties:*

1. *perfect interventions on disjoint subsets of variables commute;*
2. *acyclicity is preserved under perfect intervention.*

---

[10]For $\mathcal{U} = \emptyset$ we always consider the trivial mapping $\boldsymbol{f}_{\emptyset} : \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{E}} \to \boldsymbol{\mathcal{X}}_{\emptyset}$ where $\boldsymbol{\mathcal{X}}_{\emptyset}$ is the singleton $\boldsymbol{1}$.

The following example shows that an SCM with a solution may not have a solution anymore after performing a perfect intervention on the SCM, and vice versa, that an SCM without a solution may yield an SCM with a solution after intervention.

EXAMPLE 2.17 (Intervened SCM and its graphs). *Consider the SCM $\mathcal{M}$ of Example 2.9 which has a solution if and only if $\alpha \geq 0$. Applying the perfect intervention $\mathrm{do}(\{3\}, 1)$ to $\mathcal{M}$ gives the intervened model $\mathcal{M}_{\mathrm{do}(\{3\},1)}$ with the intervened causal mechanism*

$$\tilde{f}_1(\boldsymbol{x}, \boldsymbol{e}) = x_1 - x_1^2 + \alpha e_1^2, \qquad \tilde{f}_3(\boldsymbol{x}, \boldsymbol{e}) = 1, \qquad \tilde{f}_5(\boldsymbol{x}, \boldsymbol{e}) = x_4 \cdot e_3,$$

$$\tilde{f}_2(\boldsymbol{x}, \boldsymbol{e}) = x_1 + x_3 + x_4 + e_1, \qquad \tilde{f}_4(\boldsymbol{x}, \boldsymbol{e}) = x_2 + e_2,$$

*for which the augmented graph $\mathcal{G}^a(\mathcal{M}_{\mathrm{do}(\{3\},1)})$ and the graph $\mathcal{G}(\mathcal{M}_{\mathrm{do}(\{3\},1)})$ are depicted in Figure 2 (right). This is an example where a perfect intervention leads to an intervened SCM $\mathcal{M}_{\mathrm{do}(\{3\},1)}$ that does not have a solution anymore. In addition, performing a perfect intervention $\mathrm{do}(\{4\}, 1)$ on $\mathcal{M}_{\mathrm{do}(\{3\},1)}$ yields again an SCM with a solution for $\alpha \geq 0$.*

Remember that for each solution $\boldsymbol{X}$ of an SCM $\mathcal{M}$ we call the distribution $\mathbb{P}^{\boldsymbol{X}}$ the observational distribution of $\mathcal{M}$ associated to $\boldsymbol{X}$. For cyclic SCMs the observational distribution is in general not unique.[11] For example, the SCM $\mathcal{M}$ of Example 2.9 has two different observational distributions if $\alpha > 0$. Similarly, an intervened SCM may induce a distribution that is not unique. Whenever the intervened SCM $\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}$ has a solution $\boldsymbol{X}$ we therefore call the distribution $\mathbb{P}^{\boldsymbol{X}}$ the *interventional distribution of $\mathcal{M}$ under the perfect intervention* $\mathrm{do}(I,\boldsymbol{\xi}_I)$ *associated to* $\boldsymbol{X}$.[12]

2.5. *Counterfactuals* The causal semantics of an SCM are described by the interventions on the SCM. Adding another layer of complexity, one can describe the counterfactual semantics of an SCM by the interventions on the so-called twin SCM, an idea introduced in [1].

DEFINITION 2.18 (Twin SCM). *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM. The* twin operation *maps $\mathcal{M}$ to the* twin structural causal model (twin SCM)

$$\mathcal{M}^{\mathrm{twin}} := \langle \mathcal{I} \cup \mathcal{I}', \mathcal{J}, \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \tilde{\boldsymbol{f}}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle,$$

*where $\mathcal{I}' = \{i' : i \in \mathcal{I}\}$ is a copy of $\mathcal{I}$ and the causal mechanism $\tilde{\boldsymbol{f}} : \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{E}} \to \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{X}}$ is the measurable function given by $\tilde{\boldsymbol{f}}(\boldsymbol{x}, \boldsymbol{x}', \boldsymbol{e}) = \big(\boldsymbol{f}(\boldsymbol{x}, \boldsymbol{e}), \boldsymbol{f}(\boldsymbol{x}', \boldsymbol{e})\big)$.*

The twin operation on SCMs preserves the equivalence relation $\equiv$ on the set of all SCMs. We define an analogous twin operation $\mathrm{twin}(\mathcal{I})$ on directed graphs.

DEFINITION 2.19 (Twin graph). *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a directed graph and $\mathcal{I} \subseteq \mathcal{V}$ a subset such that $\mathcal{J} := \mathcal{V} \setminus \mathcal{I}$ is exogenous, i.e., $\mathrm{pa}_{\mathcal{G}}(\mathcal{J}) = \emptyset$. The $\mathrm{twin}(\mathcal{I})$ operation maps $\mathcal{G}$ to the twin graph w.r.t. $\mathcal{I}$ defined by $\mathrm{twin}(\mathcal{I})(\mathcal{G}) := (\tilde{\mathcal{V}}, \tilde{\mathcal{E}})$, where*

1. $\tilde{\mathcal{V}} = \mathcal{V} \cup \mathcal{I}'$, where $\mathcal{I}'$ is a copy of $\mathcal{I}$,

---

[11]In order to assure the existence of a unique observational distribution it is common to consider only SCMs for which the structural equations have a unique solution (see for example Definition 7.1.1 in [51]). Although these SCMs induce a unique observational distribution, they generally do not induce a unique distribution after a perfect intervention.

[12]In the literature, one often finds the notation $p(\boldsymbol{x})$ and $p(\boldsymbol{x} \mid \mathrm{do}(\boldsymbol{X}_I = \boldsymbol{x}_I))$ for the densities of the observational and interventional distribution, respectively, in case these are uniquely defined by the SCM [e.g. 51].

2. $\tilde{\mathcal{E}} = \mathcal{E} \cup \mathcal{E}'$, *where $\mathcal{E}'$ is given by*

$$\mathcal{E}' = \{j \to i' : j \in \mathcal{J}, i \in \mathcal{I}, j \to i \in \mathcal{E}\} \cup \{\tilde{i}' \to i' : \tilde{i}, i \in \mathcal{I}, \tilde{i} \to i \in \mathcal{E}\}$$

*with $i', \tilde{i}' \in \mathcal{I}'$ the respective copies of $i, \tilde{i} \in \mathcal{I}$.*

These two twin operations are compatible with the augmented graph mapping and preserve acyclicity.

PROPOSITION 2.20. *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM. Then $(\mathcal{G}^a \circ \mathrm{twin})(\mathcal{M}) = (\mathrm{twin}(\mathcal{I}) \circ \mathcal{G}^a)(\mathcal{M})$.*

PROPOSITION 2.21. *For SCMs and directed graphs we have that acyclicity is preserved under the twin operation.*

The perfect intervention and the twin operation for SCMs and directed graphs commute with each other in the following way.

PROPOSITION 2.22. *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM and $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ a directed graph. Then we have that perfect intervention commutes with the twin operation on both*

1. *the SCM $\mathcal{M}$: for a subset $I \subseteq \mathcal{I}$ and value $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$, $(\mathrm{do}(I \cup I', \boldsymbol{\xi}_{I \cup I'})) \circ \mathrm{twin})(\mathcal{M}) = (\mathrm{twin} \circ \mathrm{do}(I, \boldsymbol{\xi}_I))(\mathcal{M})$, and*
2. *the directed graph $\mathcal{G}$: for subsets $I \subseteq \mathcal{I} \subseteq \mathcal{V}$ such that $\mathcal{J} := \mathcal{V} \setminus \mathcal{I}$ is exogenous, $(\mathrm{do}(I \cup I') \circ \mathrm{twin}(\mathcal{I}))(\mathcal{G}) = (\mathrm{twin}(\mathcal{I}) \circ \mathrm{do}(I))(\mathcal{G})$,*

*where $I'$ is the copy of $I$ in $\mathcal{I}'$ and $\boldsymbol{\xi}_{I'} = \boldsymbol{\xi}_I$.*

Whenever the intervened twin SCM $(\mathcal{M}^{\mathrm{twin}})_{\mathrm{do}(\tilde{I}, \boldsymbol{\xi}_{\tilde{I}})}$, where $\tilde{I} \subseteq \mathcal{I} \cup \mathcal{I}'$ and $\boldsymbol{\xi}_{\tilde{I}} \in \boldsymbol{\mathcal{X}}_{\tilde{I}}$, has a solution $(\boldsymbol{X}, \boldsymbol{X}')$, we call the distribution $\mathbb{P}^{(\boldsymbol{X}, \boldsymbol{X}')}$ the *counterfactual distribution of $\mathcal{M}$ under the perfect intervention* $\mathrm{do}(\tilde{I}, \boldsymbol{\xi}_{\tilde{I}})$ *associated to* $(\boldsymbol{X}, \boldsymbol{X}')$. In Example D.3 we provide an example of how counterfactuals can be sensibly formulated for a well-known market equilibrium model described in terms of a cyclic SCM.

The interpretation of counterfactual statements has received a lot of attention in the literature [36, 66, 8, 1, 51]. For acyclic graphs, an alternative graphical approach to counterfactuals is the framework of Single World Intervention Graphs (SWIGs) [64]. One topic of discussion is that there exist SCMs that induce the same observational and interventional distributions, but differ in their counterfactual statements [11] (see also Example D.5). This raises the question how one can estimate such SCMs from data.

**3. Solvability** In this section we introduce the notions of solvability and unique solvability with respect to a subset of the endogenous variables of an SCM. They describe the existence and uniqueness of measurable solution functions for the subsystem of structural equations that correspond with a certain subset of the endogenous variables. These notions play a central role in formulating sufficient conditions under which several properties of acyclic SCMs may be extended to the cyclic setting. For example, we show that solvability of an SCM is a sufficient and necessary condition for the existence of a solution of an SCM. Further, unique solvability of an SCM implies the uniqueness of the induced observational distribution.
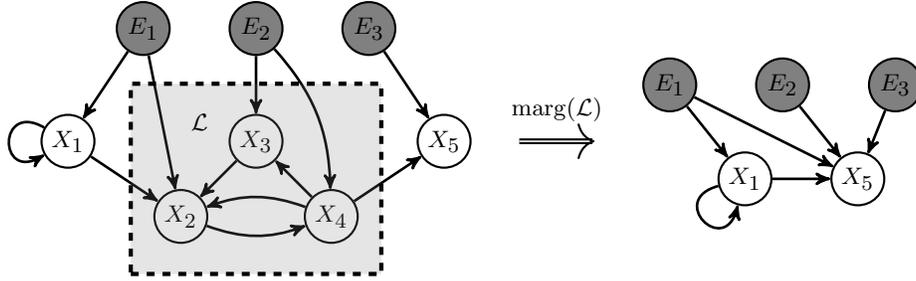
Fig 3: *The augmented graphs of the SCM $\mathcal{M}$ (left) and $\tilde{\mathcal{M}}$ (right) of Example 2.9, 3.2 and 5.2, where the SCM $\tilde{\mathcal{M}}$ is a marginalization of $\mathcal{M}$ w.r.t. $\mathcal{L}$.*

3.1. *Definition of solvability*   Intuitively, one can think of the structural equations corresponding to a subset of endogenous variables $\mathcal{O} \subseteq \mathcal{I}$ as a description of how the subsystem formed by the variables $\mathcal{O}$ interacts with the rest of the system $\mathcal{I} \setminus \mathcal{O}$ through the variables $\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}$. A solution function w.r.t. $\mathcal{O}$ assigns each input value $(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})})$ of this subsystem to a specific output value $\boldsymbol{x}_{\mathcal{O}}$ of the subsystem. This is formalized as follows.

DEFINITION 3.1 (Solvability).   *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM. We call $\mathcal{M}$ solvable w.r.t. $\mathcal{O} \subseteq \mathcal{I}$ if there exists a measurable mapping $\boldsymbol{g}_{\mathcal{O}} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{O})} \to \boldsymbol{\mathcal{X}}_{\mathcal{O}}$ such that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$*

$$\boldsymbol{x}_{\mathcal{O}} = \boldsymbol{g}_{\mathcal{O}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}) \quad \Longrightarrow \quad \boldsymbol{x}_{\mathcal{O}} = \boldsymbol{f}_{\mathcal{O}}(\boldsymbol{x}, \boldsymbol{e}).$$

*We then call $\boldsymbol{g}_{\mathcal{O}}$ a measurable solution function w.r.t. $\mathcal{O}$ for $\mathcal{M}$. We call $\mathcal{M}$ solvable if it is solvable w.r.t. $\mathcal{I}$.*

By definition, solvability w.r.t. a subset respects the equivalence relation $\equiv$ on SCMs.

EXAMPLE 3.2 (Different cases of solvability).   *Consider the SCM $\mathcal{M}$ of Example 2.9 and the subset of endogenous variables $\{2, 3, 4\}$ which is depicted by the box around the nodes in the augmented graph in Figure 3 (left). For each input value $x_1 \in \mathcal{X}_1$ and $(e_1, e_2) \in \boldsymbol{\mathcal{E}}_{\{1,2\}}$ of the box, the structural equations for the variables $\{2, 3, 4\}$ have a unique output for $x_2, x_3$ and $x_4$, which is given by the mapping $\boldsymbol{g}_{\{2,3,4\}} : \mathbb{R}^3 \to \mathbb{R}^3$ defined by $\boldsymbol{g}_{\{2,3,4\}}(x_1, e_1, e_2) := (x_1 + e_1 + e_2, -x_1 - e_1 - e_2, x_1 + e_1 + 2e_2)$. The existence of such a mapping means that $\mathcal{M}$ is solvable w.r.t. $\{2, 3, 4\}$. Solvability does not require the uniqueness of the function $\boldsymbol{g}_{\{2,3,4\}}$. For example, if we consider the subset $\{1\}$ and take $\alpha > 0$, then there exist two measurable solution functions $g_1^+, g_1^- : \mathbb{R}^1 \to \mathbb{R}^1$ of $\mathcal{M}$ w.r.t. $\{1\}$ defined by $g_1^{\pm}(e_1) := \pm\sqrt{\alpha e_1^2}$. In general, solvability does not hold w.r.t. every subset. For example, $\mathcal{M}$ is not solvable w.r.t. the subset $\{2, 4\}$, because the equality $x_1 + x_3 + e_1 + e_2 = 0$ does not hold for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$.*

The following theorem states that various possible notions of "solvability" are actually equivalent.

THEOREM 3.3 (Sufficient and necessary conditions for solvability).   *For an SCM $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ the following are equivalent:*

1. *$\mathcal{M}$ has a solution (see Definition 2.3);*
2. *for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ the structural equations*

$$\boldsymbol{x} = \boldsymbol{f}(\boldsymbol{x}, \boldsymbol{e})$$

   *have a solution $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$;*

3. $\mathcal{M}$ *is solvable (see Definition 3.1).*

While in the acyclic case, the above theorem is almost trivial, in the cyclic case the measure-theoretic aspects are not that obvious. In particular, to prove the existence of a *measurable* solution function $g : \mathcal{E}_{\mathrm{pa}(\mathcal{I})} \to \mathcal{X}$ in case the structural equations have a solution for almost every $e \in \mathcal{E}$, we make use of a strong measurable selection theorem (see Theorem F.8 or [30]). This theorem implies that if there exists a solution $X : \Omega \to \mathcal{X}$, then there necessarily exists a random variable $E : \Omega \to \mathcal{E}$ and a mapping $g : \mathcal{E}_{\mathrm{pa}(\mathcal{I})} \to \mathcal{X}$ such that $g(E_{\mathrm{pa}(\mathcal{I})})$ is a solution. However, it does not imply that there necessarily exists a random variable $E : \Omega \to \mathcal{E}$ and a mapping $g : \mathcal{E}_{\mathrm{pa}(\mathcal{I})} \to \mathcal{X}$ such that $X = g(E_{\mathrm{pa}(\mathcal{I})})$ holds a.s., e.g. if $X$ is a non-trivial mixture of such solutions, as in the following example.

EXAMPLE 3.4 (Mixtures of solutions are solutions). *Let* $\mathcal{M} = \langle \mathbf{1}, \emptyset, \mathbb{R}, \mathbf{1}, f, \mathbb{P}_{\mathbf{1}} \rangle$ *be an SCM with causal mechanism* $f : \mathcal{X} \times \mathcal{E} \to \mathcal{X}$ *defined by* $f(x, e) = x - x^2 + 1$. *There exist only two measurable solution functions* $g_{\pm} : \mathcal{E} \to \mathcal{X}$ *for* $\mathcal{M}$, *defined by* $g_{\pm}(e) = \pm 1$. *Let* $X : \Omega \to \mathbb{R}$ *be a random variable that is a non-trivial mixture of point masses on* $\{-1, +1\}$. *Then* $X$ *is a solution of* $\mathcal{M}$, *however neither* $g_{+}(E) = X$ *a.s., nor* $g_{-}(E) = X$ *a.s., for any random variable* $E$ *such that* $\mathbb{P}^E = \mathbb{P}_{\mathcal{E}}$.

Solvability w.r.t. a strict subset of $\mathcal{I}$ is in general neither sufficient nor necessary for the existence of a (global) solution of the SCM. Consider for example the SCM $\mathcal{M}$ in Example 2.9 with $\alpha < 0$. Even though this SCM is solvable w.r.t. $\{2, 3, 4\}$, it is not (globally) solvable, and hence does not have any solution. In Proposition B.1 we provide a sufficient condition for solvability w.r.t. a strict subset of $\mathcal{I}$ that is similar to condition (2) in Theorem 3.3 in the sense that it is formulated in terms of the solutions of (a subset of) the structural equations without requiring measurability of the solutions. For the class of linear SCMs we provide in Proposition C.2 a sufficient and necessary condition for solvability w.r.t. a subset of $\mathcal{I}$.

3.2. *Unique solvability*    The notion of unique solvability w.r.t. a subset $\mathcal{O} \subseteq \mathcal{I}$ is similar to the notion of solvability, but with the additional requirement that the measurable solution function $g_{\mathcal{O}} : \mathcal{X}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}_{\mathcal{O}}$ is unique up to a $\mathbb{P}_{\mathcal{E}}$-null set.

DEFINITION 3.5 (Unique solvability). *Let* $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \mathcal{X}, \mathcal{E}, f, \mathbb{P}_{\mathcal{E}} \rangle$ *be an SCM. We call* $\mathcal{M}$ *uniquely solvable w.r.t.* $\mathcal{O} \subseteq \mathcal{I}$ *if there exists a measurable mapping* $g_{\mathcal{O}} : \mathcal{X}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}_{\mathcal{O}}$ *such that for* $\mathbb{P}_{\mathcal{E}}$-*almost every* $e \in \mathcal{E}$ *and for all* $x \in \mathcal{X}$

$$x_{\mathcal{O}} = g_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}) \iff x_{\mathcal{O}} = f_{\mathcal{O}}(x, e).$$

*We call* $\mathcal{M}$ *uniquely solvable if it is uniquely solvable w.r.t.* $\mathcal{I}$.

If $\mathcal{M} \equiv \tilde{\mathcal{M}}$ and $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{O}$, then $\tilde{\mathcal{M}}$ is uniquely solvable w.r.t. $\mathcal{O}$ as well, and the same mapping $g_{\mathcal{O}}$ is a measurable solution function w.r.t. $\mathcal{O}$ for both $\mathcal{M}$ and $\tilde{\mathcal{M}}$.

The following result explains why the notions of (unique) solvability do not play an important role in the theory of acyclic SCMs.

PROPOSITION 3.6.    *An acyclic SCM* $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \mathcal{X}, \mathcal{E}, f, \mathbb{P}_{\mathcal{E}} \rangle$ *is uniquely solvable w.r.t. every subset* $\mathcal{O} \subseteq \mathcal{I}$.

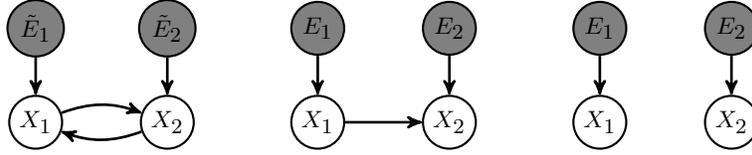The next example illustrates that also cyclic SCMs can be uniquely solvable w.r.t. every subset.

Fig 4: *The augmented graph of the SCM $\tilde{\mathcal{M}}$ (left) of Example 3.7, $\mathcal{M}$ (center) of Example 4.2 and 4.4, and $\bar{\mathcal{M}}$ (right) of Example 4.4. The SCMs $\tilde{\mathcal{M}}$ and $\mathcal{M}$ are observationally equivalent in Example 4.2, but not interventionally equivalent.*

EXAMPLE 3.7 (Cyclic SCM, uniquely solvable w.r.t. each subset).    *Consider the linear SCM $\tilde{\mathcal{M}} = \langle \mathbf{2}, \mathbf{2}, \mathbb{R}^2, \mathbb{R}^2, \tilde{\boldsymbol{f}}, \mathbb{P}_{\tilde{\boldsymbol{\mathcal{E}}}} \rangle$ with causal mechanism given by*

$$\tilde{f}_1(\boldsymbol{x}, \tilde{\boldsymbol{e}}) = \alpha x_2 + \tilde{e}_1, \quad \tilde{f}_2(\boldsymbol{x}, \tilde{\boldsymbol{e}}) = \beta x_1 + \tilde{e}_2,$$

*with $\alpha, \beta \neq 0$, $\alpha\beta \neq 1$, and $\mathbb{P}_{\tilde{\boldsymbol{\mathcal{E}}}} = \mathbb{P}^{\tilde{\boldsymbol{E}}}$ with normal distributions $\tilde{E}_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$, $\tilde{E}_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$ and $\tilde{E}_1 \perp\!\!\!\perp \tilde{E}_2$. This SCM $\tilde{\mathcal{M}}$ is uniquely solvable w.r.t. every subset and its (augmented) graph includes a cycle (see Figure 12 on the left).*

Theorem 3.3 provides sufficient and necessary conditions for (global) solvability. The next theorem states that under the additional uniqueness requirement there exists a sufficient and necessary condition for unique solvability w.r.t. any subset (for solvability w.r.t. a subset we only have the sufficient condition provided in Proposition B.1), and moreover, that all solutions of a uniquely solvable SCM induce the same observational distribution.

THEOREM 3.8 (Sufficient and necessary conditions for unique solvability).    *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM and $\mathcal{O} \subseteq \mathcal{I}$ a subset of endogenous variables. The following are equivalent:*

1. *for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x}_{\setminus \mathcal{O}} \in \boldsymbol{\mathcal{X}}_{\setminus \mathcal{O}}$ the structural equations*

$$\boldsymbol{x}_{\mathcal{O}} = \boldsymbol{f}_{\mathcal{O}}(\boldsymbol{x}, \boldsymbol{e})$$

   *have a unique solution $\boldsymbol{x}_{\mathcal{O}} \in \boldsymbol{\mathcal{X}}_{\mathcal{O}}$;*
2. *$\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{O}$.*

*Furthermore, if $\mathcal{M}$ is uniquely solvable, then there exists a solution, and all solutions have the same observational distribution.*

It is well-known that under acyclicity the observational distribution is unique. Theorem 3.8 generalizes this result to settings with cycles. For linear SCMs the unique solvability condition w.r.t. a subset of endogenous variables is equivalent to a matrix invertibility condition (see Proposition C.3).

In general, (unique) solvability w.r.t. $\mathcal{O} \subseteq \mathcal{I}$ does not imply (unique) solvability w.r.t. a strict superset $\mathcal{O} \subsetneq \mathcal{V} \subseteq \mathcal{I}$ nor w.r.t. a strict subset $\mathcal{W} \subsetneq \mathcal{O}$ (see Example B.2). Moreover, (unique) solvability is in general not preserved under unions and intersections (see Appendix B.3).

3.3. *Self-cycles*    One can think of a structural equation of a single endogenous variable $i \in \mathcal{I}$ as describing a small subsystem that interacts with the rest of the system. If the output $x_i$ of this subsystem is uniquely determined by the input $(\boldsymbol{x}_{\setminus i}, \boldsymbol{e})$ from the rest of the system (up to a $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-null set), then $i$ is not a parent of itself (see Definition 2.7).

PROPOSITION 3.9 (Self-cycles).    *The SCM $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ is uniquely solvable w.r.t. $\{i\}$ for $i \in \mathcal{I}$ if and only if $\mathcal{G}^a(\mathcal{M})$ (or $\mathcal{G}(\mathcal{M})$) has no self-cycle $i \to i$ at $i \in \mathcal{I}$.*

A self-cycle at an endogenous variable denotes that that variable is not uniquely determined by its parents, up to a $\mathbb{P}_{\mathcal{E}}$-null set. This implies that an SCM with a self-cycle at an endogenous variable in its graph can be either solvable, or not solvable, w.r.t. that variable. For the SCM $\mathcal{M}$ of Example 2.9 we have indeed that it is solvable w.r.t. $\{1\}$ for $\alpha > 0$, while for $\alpha < 0$ it is not. For linear SCMs with structural equations $X_i = \sum_{j \in \mathcal{I}} B_{ij} X_j + \sum_{k \in \mathcal{J}} \Gamma_{ik} E_k$ the endogenous variable $i \in \mathcal{I}$ has a self-cycle if and only if $B_{ii} = 1$ (see Appendix C).

3.4. *Interventions* The property of (unique) solvability is in general not preserved under perfect intervention. For example, a (uniquely) solvable SCM can lead to a non-uniquely solvable SCM after intervention, which either has no solution or has solutions with multiple induced distributions.

EXAMPLE 3.10 (Solvability is not preserved under perfect intervention). *Consider the SCM $\mathcal{M} = \langle \mathbf{2}, \emptyset, \mathbb{R}^2, \mathbf{1}, \boldsymbol{f}, \mathbb{P_1} \rangle$ with the following causal mechanism*

$$f_1(\boldsymbol{x}) = x_1 + x_1^2 - x_2 + 1\,, \quad f_2(\boldsymbol{x}) = x_2(1 - \mathbf{1}_{\{0\}}(x_1)) + 1\,.$$

*This SCM is (uniquely) solvable. Doing a perfect intervention $\mathrm{do}(\{1\}, \xi_1)$ for some $\xi_1 \neq 0$, however, leads to an intervened model $\mathcal{M}_{\mathrm{do}(\{1\}, \xi_1)}$ that is not solvable. Performing instead the perfect intervention $\mathrm{do}(\{2\}, \xi_2)$ for some $\xi_2 > 1$ leads also to a non-uniquely solvable SCM $\mathcal{M}_{\mathrm{do}(\{2\}, \xi_2)}$ which has solutions with multiple induced distributions, e.g., $(X_1, X_2) = (\phi(\xi_2)\sqrt{\xi_2 - 1}, \xi_2)$ with some measurable $\phi : \mathbb{R} \to \{-1, +1\}$, but also mixtures of those.*

A sufficient condition for the intervened SCM to be (uniquely) solvable is that the original SCM has to be (uniquely) solvable w.r.t. the subset of non-intervened endogenous variables.

PROPOSITION 3.11. *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P_{\mathcal{E}}} \rangle$ be an SCM that is (uniquely) solvable w.r.t. $\mathcal{O} \subseteq \mathcal{I}$. Then, for any set $I$ such that $\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O} \subseteq I \subseteq \mathcal{I} \setminus \mathcal{O}$ and value $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$ the intervened SCM $\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$ is (uniquely) solvable w.r.t. $\mathcal{O} \cup I$.*

Proposition 3.6 shows that acyclic SCMs are uniquely solvable w.r.t. every subset and hence are uniquely solvable after every perfect intervention. This also directly follows from the fact that acyclicity is preserved under perfect intervention (see Proposition 2.16). Moreover, since acyclicity is preserved under the twin operation (see Proposition 2.21), an acyclic SCM induces unique observational, interventional and counterfactual distributions.

3.5. *Ancestral (unique) solvability* We saw that, in general, solvability w.r.t. $\mathcal{O} \subseteq \mathcal{I}$ does not imply solvability w.r.t. a strict subset of $\mathcal{O}$. Here we show that it does imply solvability w.r.t. the ancestral subsets in $\mathcal{G}(\mathcal{M})_{\mathcal{O}}$, that is, in the induced subgraph of the graph $\mathcal{G}(\mathcal{M})$ on $\mathcal{O}$. A subset $\mathcal{A} \subseteq \mathcal{O}$ is called an *ancestral subset* in $\mathcal{G}(\mathcal{M})_{\mathcal{O}}$ if $\mathcal{A} = \mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(\mathcal{A})$, where $\mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(\mathcal{A})$ are the ancestors of $\mathcal{A}$ according to the induced subgraph[13] $\mathcal{G}(\mathcal{M})_{\mathcal{O}}$.

DEFINITION 3.12 (Ancestral (unique) solvability). *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P_{\mathcal{E}}} \rangle$ be an SCM. We call $\mathcal{M}$ ancestrally (uniquely) solvable w.r.t. $\mathcal{O} \subseteq \mathcal{I}$ if $\mathcal{M}$ is (uniquely) solvable w.r.t. every ancestral subset in $\mathcal{G}(\mathcal{M})_{\mathcal{O}}$. We call $\mathcal{M}$ ancestrally (uniquely) solvable if it is ancestrally (uniquely) solvable w.r.t. $\mathcal{I}$.*

---

[13]Here, one can also use the augmented graph $\mathcal{G}^a(\mathcal{M})$ on $\mathcal{O}$ since $\mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(\mathcal{A}) = \mathrm{an}_{\mathcal{G}^a(\mathcal{M})_{\mathcal{O}}}(\mathcal{A})$ for every subset $\mathcal{A} \subseteq \mathcal{O} \subseteq \mathcal{I}$.

Fig 5: *The graphs of the SCM $\mathcal{M}$ (left) of Example 3.14 and the marginal SCM $\mathcal{M}_{\mathrm{marg}(\{2,3\})}$ (right) of Example 5.11.*

PROPOSITION 3.13 (Solvability is equivalent to ancestral solvability). *The SCM $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ is solvable w.r.t. the subset $\mathcal{O} \subseteq \mathcal{I}$ if and only if $\mathcal{M}$ is ancestrally solvable w.r.t. $\mathcal{O}$.*

A similar result does not hold for unique solvability. Although ancestral unique solvability w.r.t. $\mathcal{O} \subseteq \mathcal{I}$ implies unique solvability w.r.t. $\mathcal{O}$, the converse does not hold in general, as the following example illustrates.

EXAMPLE 3.14 (Unique solvability w.r.t. $\mathcal{O}$ does not imply ancestral unique solvability w.r.t. $\mathcal{O}$). *Consider the SCM $\mathcal{M} = \langle \boldsymbol{4}, \boldsymbol{1}, \mathbb{R}^4, \mathbb{R}, \boldsymbol{f}, \mathbb{P}_{\mathbb{R}} \rangle$ with causal mechanism given by*

$$f_1(\boldsymbol{x}, e) = e\,, \ f_2(\boldsymbol{x}, e) = x_2 \cdot (1 - \boldsymbol{1}_{\{0\}}(x_1 - x_3)) + 1\,, \ f_3(\boldsymbol{x}, e) = x_3\,, \ f_4(\boldsymbol{x}, e) = x_3$$

*and $\mathbb{P}_{\mathbb{R}}$ the standard-normal measure on $\mathbb{R}$. This SCM is uniquely solvable w.r.t. the set $\{2, 3\}$, and thus solvable w.r.t. this set. Although it is solvable w.r.t. the ancestral subset $\{3\}$ in $\mathcal{G}(\mathcal{M})_{\{2,3\}}$, depicted in Figure 5 (left), it is not uniquely solvable w.r.t. this subset. Hence, it is not ancestrally uniquely solvable w.r.t. $\{2, 3\}$.*

However, for the class of linear SCMs we have that unique solvability w.r.t. $\mathcal{O}$ always implies ancestral unique solvability w.r.t. $\mathcal{O}$ (see Proposition C.4).

Although in general unique solvability is not preserved under unions, in Proposition B.4 we show that if an SCM is uniquely solvable w.r.t. two ancestral subsets and w.r.t. their intersection, then it is uniquely solvable w.r.t. their union.

In general, the property of ancestral unique solvability is not preserved under perfect intervention, as can be seen in Example 3.10. The notion of ancestral unique solvability will appear in various results in Sections 5 and 6.

**4. Equivalences** In Section 2 we already encountered an equivalence relation on the class of SCMs (see Definition 2.6). The (augmented) graph of an SCM, its solutions and its induced observational, interventional and counterfactual distributions are preserved under this equivalence relation. In this section we give several coarser equivalence relations on the class of SCMs: observational, interventional and counterfactual equivalence.

4.1. *Observational equivalence* Observational equivalence is the property that two SCMs are indistinguishable on the basis of their observational distributions.

DEFINITION 4.1 (Observational equivalence). *Two SCMs $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ and $\tilde{\mathcal{M}} = \langle \tilde{\mathcal{I}}, \tilde{\mathcal{J}}, \tilde{\boldsymbol{\mathcal{X}}}, \tilde{\boldsymbol{\mathcal{E}}}, \tilde{\boldsymbol{f}}, \mathbb{P}_{\tilde{\boldsymbol{\mathcal{E}}}} \rangle$ are observationally equivalent w.r.t. $\mathcal{O} \subseteq \mathcal{I} \cap \tilde{\mathcal{I}}$, denoted by $\mathcal{M} \equiv_{obs(\mathcal{O})} \tilde{\mathcal{M}}$, if $\boldsymbol{\mathcal{X}}_{\mathcal{O}} = \tilde{\boldsymbol{\mathcal{X}}}_{\mathcal{O}}$ and for all solutions $\boldsymbol{X}$ of $\mathcal{M}$ there exists a solution $\tilde{\boldsymbol{X}}$ of $\tilde{\mathcal{M}}$ such that $\mathbb{P}^{\boldsymbol{X}_{\mathcal{O}}} = \mathbb{P}^{\tilde{\boldsymbol{X}}_{\mathcal{O}}}$ and for all solutions $\tilde{\boldsymbol{X}}$ of $\tilde{\mathcal{M}}$ there exists a solution $\boldsymbol{X}$ of $\mathcal{M}$ such that $\mathbb{P}^{\boldsymbol{X}_{\mathcal{O}}} = \mathbb{P}^{\tilde{\boldsymbol{X}}_{\mathcal{O}}}$. $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are called* observationally equivalent *if they are observationally equivalent w.r.t. $\mathcal{I} = \tilde{\mathcal{I}}$.*

Equivalent SCMs have the same solutions, and hence they are observationally equivalent w.r.t. every subset $\mathcal{O} \subseteq \mathcal{I}$. However, observational equivalence does not imply equivalence, as the following example illustrates.

EXAMPLE 4.2 (Observational equivalence does not imply equivalence). *Let $\tilde{\mathcal{M}}$ be the SCM of Example 3.7. Then, one can always construct an SCM $\mathcal{M} = \langle \mathbf{2}, \mathbf{2}, \mathbb{R}^2, \mathbb{R}^2, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ with a causal mechanism of the form $f_1(\boldsymbol{x}, \boldsymbol{e}) = e_1$ and $f_2(\boldsymbol{x}, \boldsymbol{e}) = \gamma x_1 + e_2$ with $\gamma \neq 0$ such that $\tilde{\mathcal{M}}$ and $\mathcal{M}$ are observationally equivalent (see Example D.4 for more details). Because both SCMs have a different (augmented) graph they are not equivalent to each other (see Figure 12).*

This example shows that if two SCMs $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are observationally equivalent, then their associated augmented graphs $\mathcal{G}^a(\mathcal{M})$ and $\mathcal{G}^a(\tilde{\mathcal{M}})$ are not necessarily equal to each other. Although the SCMs of this example are observationally equivalent, they are not interventionally equivalent, as we will see in the next subsection.

4.2. *Interventional equivalence*  We consider two SCMs to be interventionally equivalent if they induce the same interventional distributions under all perfect interventions.

DEFINITION 4.3 (Interventional equivalence). *Two SCMs $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ and $\tilde{\mathcal{M}} = \langle \tilde{\mathcal{I}}, \tilde{\mathcal{J}}, \tilde{\boldsymbol{\mathcal{X}}}, \tilde{\boldsymbol{\mathcal{E}}}, \tilde{\boldsymbol{f}}, \mathbb{P}_{\tilde{\boldsymbol{\mathcal{E}}}} \rangle$ are* interventionally equivalent w.r.t. $\mathcal{O} \subseteq \mathcal{I} \cap \tilde{\mathcal{I}}$, *denoted by $\mathcal{M} \equiv_{int(\mathcal{O})} \tilde{\mathcal{M}}$, if $\boldsymbol{\mathcal{X}}_{\mathcal{O}} = \tilde{\boldsymbol{\mathcal{X}}}_{\mathcal{O}}$ and for every $I \subseteq \mathcal{O}$ and every value $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$ their intervened models $\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$ and $\tilde{\mathcal{M}}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$ are observationally equivalent with respect to $\mathcal{O}$. $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are called* interventionally equivalent *if they are interventionally equivalent w.r.t. $\mathcal{I} = \tilde{\mathcal{I}}$.*

Equivalent SCMs have the same solutions under every perfect intervention, and hence they are interventionally equivalent w.r.t. every subset $\mathcal{O} \subseteq \mathcal{I}$. SCMs that are interventionally equivalent w.r.t. a subset $\mathcal{O} \subseteq \mathcal{I}$ are interventionally equivalent w.r.t. every strict subset $\mathcal{W} \subsetneq \mathcal{O}$. But, they are, in general, not interventionally equivalent w.r.t. a strict superset $\mathcal{O} \subsetneq \mathcal{V} \subseteq \mathcal{I}$, as can be seen in Example 4.2, where the SCMs $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are interventionally equivalent w.r.t. $\{1\}$ but are not interventionally equivalent.

Interventional equivalence w.r.t. $\mathcal{O} \subseteq \mathcal{I}$ implies observational equivalence w.r.t. $\mathcal{O}$, since the empty perfect intervention ($I = \emptyset$) is a special case of a perfect intervention. However, observational equivalence w.r.t. $\mathcal{O} \subseteq \mathcal{I}$ does, in general, not imply interventional equivalence w.r.t. $\mathcal{O}$, as can be seen in Example 4.2, where the SCM $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are observationally equivalent but not interventionally equivalent (see Figure 12).

Although interventional equivalence is a finer notion than observational equivalence, we have that if two SCMs $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are interventionally equivalent, then their associated augmented graphs $\mathcal{G}^a(\mathcal{M})$ and $\mathcal{G}^a(\tilde{\mathcal{M}})$ are not necessarily equal to each other, as is shown in the following example.

EXAMPLE 4.4 (Interventionally equivalent SCMs with different graphs). *Consider the SCM $\mathcal{M} = \langle \mathbf{2}, \mathbf{2}, \{-1, 1\}^2, \{-1, 1\}^2, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ and the SCM $\bar{\mathcal{M}}$ that is the same as $\mathcal{M}$ except for its causal mechanism $\bar{\boldsymbol{f}}$, where the causal mechanisms are given by*

$$f_1(\boldsymbol{x}, \boldsymbol{e}) = e_1, \quad f_2(\boldsymbol{x}, \boldsymbol{e}) = x_1 e_2, \qquad \bar{f}_1(\boldsymbol{x}, \boldsymbol{e}) = e_1, \quad \bar{f}_2(\boldsymbol{x}, \boldsymbol{e}) = e_2,$$

*and $\mathbb{P}_{\boldsymbol{\mathcal{E}}} = \mathbb{P}^{\boldsymbol{E}}$ with $E_1, E_2 \sim \mathcal{U}(\{-1, 1\})$ uniformly distributed and $E_1 \perp\!\!\!\perp E_2$. Then $\mathcal{M}$ and $\bar{\mathcal{M}}$ are interventionally equivalent although $\mathcal{G}^a(\mathcal{M})$ is not equal to $\mathcal{G}^a(\bar{\mathcal{M}})$ (see Figure 12).*

4.3. *Counterfactual equivalence* We consider two SCMs to be counterfactually equivalent if their twin SCMs induce the same counterfactual distributions under every perfect intervention.

DEFINITION 4.5 (Counterfactual equivalence). *Two SCMs $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ and $\tilde{\mathcal{M}} = \langle \tilde{\mathcal{I}}, \tilde{\mathcal{J}}, \tilde{\boldsymbol{\mathcal{X}}}, \tilde{\boldsymbol{\mathcal{E}}}, \tilde{\boldsymbol{f}}, \mathbb{P}_{\tilde{\boldsymbol{\mathcal{E}}}} \rangle$ are* counterfactually equivalent *with respect to $\mathcal{O} \subseteq \mathcal{I} \cap \tilde{\mathcal{I}}$, denoted by $\mathcal{M} \equiv_{cf(\mathcal{O})} \tilde{\mathcal{M}}$, if $\mathcal{M}^{\mathrm{twin}}$ and $\tilde{\mathcal{M}}^{\mathrm{twin}}$ are interventionally equivalent with respect to $\mathcal{O} \cup \mathcal{O}'$, where $\mathcal{O}'$ corresponds to the copy of $\mathcal{O}$ in $\mathcal{I}' \cap \tilde{\mathcal{I}}'$. $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are called* counterfactually equivalent *if they are counterfactually equivalent with respect to $\mathcal{I} = \tilde{\mathcal{I}}$.*

The notion of counterfactual equivalence is coarser than equivalence and finer than interventional equivalence.

PROPOSITION 4.6. *For SCMs we have that equivalence implies counterfactual equivalence w.r.t. $\mathcal{O}$, which in turn implies interventional equivalence w.r.t. $\mathcal{O}$, for any $\mathcal{O} \subseteq \mathcal{I}$.*

EXAMPLE 4.7 (Interventional equivalence does not imply counterfactual equivalence). *Consider the same SCMs as in Example 4.4. We have seen that they are interventionally equivalent. However, they are not counterfactually equivalent, as $\mathcal{M}^{\mathrm{twin}}_{\mathrm{do}(\{1',1\},(1,-1))}$ is not observationally equivalent to $\bar{\mathcal{M}}^{\mathrm{twin}}_{\mathrm{do}(\{1',1\},(1,-1))}$. To see this, consider the counterfactual query $p(X_{2'} = 1 \mid \mathrm{do}(X_{1'} = 1, X_1 = -1), X_2 = 1)$. Both SCMs give a different answer and hence $\mathcal{M}$ and $\bar{\mathcal{M}}$ are not counterfactually equivalent.*

Even interventionally equivalent SCMs with the same causal mechanism (that differ only in their exogenous distribution) may not be counterfactually equivalent (see Example D.5).

Although the notion of counterfactual equivalence is finer than the notion of observational and interventional equivalence, the (augmented) graphs for counterfactually equivalent SCMs are in general not equal to each other.

EXAMPLE 4.8 (Counterfactually equivalent SCMs with different graphs). *Consider the SCM $\mathcal{M} = \langle \mathbf{2}, \mathbf{2}, \{-1,1\}^2, \{-1,1\}^3, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ with the causal mechanism defined by $\boldsymbol{f}(\boldsymbol{x}, e_1, \boldsymbol{e}_2) := (e_1, e_{22})$, where $\boldsymbol{e}_2 = (e_{21}, e_{22}) \in \{-1,1\}^2$ and $\mathbb{P}_{\boldsymbol{\mathcal{E}}} = \mathbb{P}^{(E_1, \boldsymbol{E}_2)}$ such that $E_1, E_{21}, E_{22} \sim \mathcal{U}(\{-1,1\})$ independent. Let $\tilde{\mathcal{M}}$ be the same SCM as $\mathcal{M}$, but with the causal mechanism $\tilde{\boldsymbol{f}}(\boldsymbol{x}, e_1, \boldsymbol{e}_2) := (e_{21}, e_{22})$. Then $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are counterfactually equivalent (and, in particular, interventionally and observationally equivalent) although $1 \leftrightarrow 2 \in \mathcal{G}(\tilde{\mathcal{M}})$ but $1 \leftrightarrow 2 \notin \mathcal{G}(\mathcal{M})$.*

4.4. *Relations between equivalences* The definitions of observational, interventional and counterfactual equivalence provide equivalence relations on the set of all SCMs. For two SCMs to be observationally, interventionally or counterfactually equivalent w.r.t. $\mathcal{O} \subseteq \mathcal{I} \cap \tilde{\mathcal{I}}$, the domains of their endogenous variables $\mathcal{O}$ have to be equal, that is, $\boldsymbol{\mathcal{X}}_{\mathcal{O}} = \tilde{\boldsymbol{\mathcal{X}}}_{\mathcal{O}}$. Apart from that, the index sets of the endogenous and the exogenous variables, the spaces of the other endogenous and exogenous variables, the causal mechanism and the exogenous probability measure may all differ. The observational, interventional and counterfactual equivalence classes w.r.t. $\mathcal{O} \subseteq \mathcal{I} \cap \tilde{\mathcal{I}}$ are related in the following way (see Proposition 4.6):

$$\mathcal{M} \text{ and } \tilde{\mathcal{M}} \text{ are equivalent}$$

$$\implies \mathcal{M} \text{ and } \tilde{\mathcal{M}} \text{ are counterfactually equivalent w.r.t. } \mathcal{O}$$

$$\implies \mathcal{M} \text{ and } \tilde{\mathcal{M}} \text{ are interventionally equivalent w.r.t. } \mathcal{O}$$

$$\implies \mathcal{M} \text{ and } \tilde{\mathcal{M}} \text{ are observationally equivalent w.r.t. } \mathcal{O}.$$

This hierarchy allows us to compare SCMs at different levels of abstractions and formally establishes the "ladder" of causation (last two implications) [69, 53, 51].

**5. Marginalizations**   In this section we show how, and under which condition, one can marginalize an SCM over a subset $\mathcal{L} \subseteq \mathcal{I}$ of endogenous variables (thereby "hiding" the variables $\mathcal{L}$), to another SCM on the margin $\mathcal{I} \setminus \mathcal{L}$ that is observationally, interventionally and even counterfactually equivalent with respect to $\mathcal{I} \setminus \mathcal{L}$. In other words, we provide a formal notion of marginalization and show that this preserves the probabilistic, causal and counterfactual semantics on the margin.

The problem of marginalization of directed graphical models has been addressed for acyclic graph structures, e.g., ADMGs and mDAGs [see 78, 60, 62, 15, 16, a.o.], and more recently in [18] for certain graph structures ("HEDGes") that may include cycles. Although in the acyclic setting it has been shown that the marginalization for some of these graph structures preserves the probabilistic and causal semantics, in the cyclic setting this has only been shown for modular SCMs [18]. We show that without the additional structure of a compatible system of solution functions (see Appendix A.4) one can still define a marginalization for SCMs under certain local unique solvability conditions. Intuitively, the idea is that if the state of a subsystem of endogenous variables is uniquely determined by the parents outside of this subsystem, then one can ignore the internals of this subsystem by treating it as a "black box" that can be described by certain measurable solution functions (see Figure 3). One can marginalize over this subsystem by substituting these measurable solution functions into the rest of the model, thereby removing the functional dependencies on the variables of the subsystem from the rest of the system, while preserving the probabilistic, causal and the counterfactual semantics of the rest of the system. We show that in general this marginalization operation defined on SCMs does not respect the latent projection on its associated (augmented) graph, where the latent projection is a similar marginalization operation defined on directed mixed graphs [78, 76, 15]. We show that under certain stronger local ancestral unique solvability conditions the marginalization does respect the latent projection.

5.1. *Marginalization of a structural causal model*   Before we show how one can marginalize an SCM w.r.t. a subset of endogenous variables, we first point out that in general it is not always possible to find an SCM on the margin that preserves the causal semantics, as the following example illustrates.

EXAMPLE 5.1 (No SCM on the margin preserves the causal semantics).   *Consider the SCM* $\mathcal{M} = \langle \mathbf{3}, \emptyset, \mathbb{R}^3, \mathbf{1}, \boldsymbol{f}, \mathbb{P_1} \rangle$ *with causal mechanism*

$$f_1(\boldsymbol{x}) = x_1 + x_2 + x_3, \quad f_2(\boldsymbol{x}) = x_2, \quad f_3(\boldsymbol{x}) = 0.$$

*Then there exists no SCM* $\tilde{\mathcal{M}}$ *on the endogenous variables* $\{2, 3\}$ *that is interventionally equivalent to* $\mathcal{M}$ *w.r.t.* $\{2, 3\}$. *To see this, suppose there exists such an SCM* $\tilde{\mathcal{M}}$, *then for every* $(\xi_2, \xi_3) \in \boldsymbol{\mathcal{X}}_{\{2,3\}}$ *such that* $\xi_2 + \xi_3 \neq 0$ *the intervened model* $\tilde{\mathcal{M}}_{\mathrm{do}(\{2,3\},(\xi_2,\xi_3))}$ *has a solution but* $\mathcal{M}_{\mathrm{do}(\{2,3\},(\xi_2,\xi_3))}$ *does not.*

More generally, for an SCM $\mathcal{M}$ that is not solvable w.r.t. a subset $\mathcal{L} \subseteq \mathcal{I}$ there is no SCM $\tilde{\mathcal{M}}$ on the endogenous variables $\mathcal{I} \setminus \mathcal{L}$ that is interventionally equivalent w.r.t. $\mathcal{I} \setminus \mathcal{L}$.

The following example illustrates that for an SCM that is uniquely solvable w.r.t. a subset there exists an SCM on the margin that preserves the causal semantics.

EXAMPLE 5.2 (SCM on the margin that preserves the causal semantics).   *Consider the SCM* $\mathcal{M} = \langle \mathbf{5}, \mathbf{3}, \mathbb{R}^5, \mathbb{R}^3, \boldsymbol{f}, \mathbb{P}_{\mathbb{R}^3} \rangle$ *of Example 2.9. We saw that* $\mathcal{M}$ *is uniquely solvable w.r.t.*

$\mathcal{L} = \{2,3,4\}$ *with the measurable solution function* $\boldsymbol{g}_{\mathcal{L}}$ *given in Example 3.2. The system of structural equations for the variables* $\mathcal{L}$ *can be seen as a subsystem, that is, for* $\mathbb{P}_{\boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{L})}}$-*almost every* $\boldsymbol{e}_{\mathrm{pa}(\mathcal{L})} \in \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{L})}$ *and for every* $\boldsymbol{x}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}} \in \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}}$ *the input* $(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L})})$ *gives these equations a unique output* $\boldsymbol{x}_{\mathcal{L}} \in \boldsymbol{\mathcal{X}}_{\mathcal{L}}$. *This subsystem is depicted by the gray box in Figure 3. Substituting the components* $(g_{\mathcal{L}})_2, (g_{\mathcal{L}})_3$ *and* $(g_{\mathcal{L}})_4$ *into the causal mechanism components* $f_1, f_5$ *for the remaining endogenous variables* $\{1,5\}$ *gives a "marginal" causal mechanism*

$$\tilde{f}_1(\boldsymbol{x},\boldsymbol{e}) := x_1 - x_1^2 + \alpha e_1^2, \quad \tilde{f}_5(\boldsymbol{x},\boldsymbol{e}) := x_1 \cdot e_3 + e_1 \cdot e_3 + 2e_2 \cdot e_3 \,.$$

*These mappings define an SCM* $\tilde{\mathcal{M}} := \langle \mathbf{2}, \mathbf{3}, \mathbb{R}^2, \mathbb{R}^3, \tilde{\boldsymbol{f}}, \mathbb{P}_{\mathbb{R}^3} \rangle$ *on the margin* $\mathcal{I} \setminus \mathcal{L} = \{1,5\}$. *This constructed SCM* $\tilde{\mathcal{M}}$, *depicted in Figure 3, is interventionally equivalent w.r.t.* $\mathcal{L}$, *which can be checked manually or by applying Theorem 5.6 below.*

In general, for an SCM $\mathcal{M}$ and a given subset $\mathcal{L} \subseteq \mathcal{I}$ of endogenous variables and its complement $\mathcal{O} = \mathcal{I} \setminus \mathcal{L}$, we can consider the "subsystem" of structural equations $\boldsymbol{x}_{\mathcal{L}} = \boldsymbol{f}_{\mathcal{L}}(\boldsymbol{x}_{\mathcal{L}}, \boldsymbol{x}_{\mathcal{O}}, \boldsymbol{e})$. If $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$ with measurable solution function $\boldsymbol{g}_{\mathcal{L}} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{L})} \to \boldsymbol{\mathcal{X}}_{\mathcal{L}}$, then for each input $(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L})}) \in \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{L})}$ of the subsystem, there exists an output $\boldsymbol{x}_{\mathcal{L}} \in \boldsymbol{\mathcal{X}}_{\mathcal{L}}$, which is unique for $\mathbb{P}_{\boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{L})}}$-almost every $\boldsymbol{e}_{\mathrm{pa}(\mathcal{L})} \in \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{L})}$ and for all $\boldsymbol{x}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}} \in \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}}$. We can remove this subsystem of endogenous variables from the model by substitution. This leads to a marginal SCM that is observationally, interventionally and counterfactually equivalent to the original SCM w.r.t. the margin, as we prove in Theorem 5.6.

DEFINITION 5.3 (Marginalization of an SCM). *Let* $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ *be an SCM that is uniquely solvable w.r.t. a subset* $\mathcal{L} \subseteq \mathcal{I}$ *and let* $\mathcal{O} = \mathcal{I} \setminus \mathcal{L}$. *For* $\boldsymbol{g}_{\mathcal{L}} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{L})} \to \boldsymbol{\mathcal{L}}$ *any measurable solution function of* $\mathcal{M}$ *w.r.t.* $\mathcal{L}$, *we call the SCM* $\mathcal{M}_{\mathrm{marg}(\mathcal{L})} := \langle \mathcal{O}, \mathcal{J}, \boldsymbol{\mathcal{X}}_{\mathcal{O}}, \boldsymbol{\mathcal{E}}, \tilde{\boldsymbol{f}}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ *with the* marginal causal mechanism $\tilde{\boldsymbol{f}} : \boldsymbol{\mathcal{X}}_{\mathcal{O}} \times \boldsymbol{\mathcal{E}} \to \boldsymbol{\mathcal{X}}_{\mathcal{O}}$ *given by*

$$\tilde{\boldsymbol{f}}(\boldsymbol{x}_{\mathcal{O}}, \boldsymbol{e}) = \boldsymbol{f}_{\mathcal{O}}(\boldsymbol{g}_{\mathcal{L}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L})\setminus\mathcal{L}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L})}), \boldsymbol{x}_{\mathcal{O}}, \boldsymbol{e}) \,,$$

a *marginalization of* $\mathcal{M}$ *w.r.t.* $\mathcal{L}$. *We denote by* $\mathrm{marg}(\mathcal{L})(\mathcal{M})$ *the equivalence class of the marginalizations of* $\mathcal{M}$ *w.r.t.* $\mathcal{L}$.

The marginalization of $\mathcal{M}$ w.r.t. $\mathcal{L}$ is defined up to the equivalence $\equiv$ on SCMs, since the measurable solution functions $\boldsymbol{g}_{\mathcal{L}}$ are uniquely defined up to $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-null sets.

With this definition at hand, we can always construct a marginal SCM over a subset of the endogenous variables of an acyclic SCM by mere substitution (see also Proposition 3.6). Moreover, this definition extends that notion to SCMs that are uniquely solvable w.r.t. a certain subset. For linear SCMs this condition translates into a matrix invertibility condition, and since substitution preserves linearity, marginalization yields a linear marginal SCM (see Proposition C.5).

In general, marginalization is not always defined for all subsets. For instance, the SCM of Example 3.14 cannot be marginalized over the variable 3 (due to the self-cycle at 3), but can be marginalized over the variables 2 and 3 together. It follows from Proposition 3.9 that we can only marginalize over a single variable if that variable has no self-cycle. Note that we may introduce new self-cycles if we marginalize over a subset of variables, as can be seen, for example, from the SCM $\mathcal{M}$ in Example 2.9. This SCM has only one self-cycle, however marginalizing w.r.t. $\{2\}$ gives a marginal SCM with another self-cycle at variable 4.

The definition of marginalization satisfies an intuitive property: if we can marginalize over two disjoint subsets after each other, then we can also marginalize over the union of those subsets at once, and the respective results agree.

PROPOSITION 5.4. *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM that is uniquely solvable w.r.t. a subset $\mathcal{L}_1 \subseteq \mathcal{I}$ and let $\mathcal{L}_2 \subseteq \mathcal{I}$ be a subset disjoint from $\mathcal{L}_1$. Then $\mathcal{M}_{\mathrm{marg}(\mathcal{L}_1)}$ is uniquely solvable w.r.t. $\mathcal{L}_2$ if and only if $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}_1 \cup \mathcal{L}_2$, Moreover $\mathrm{marg}(\mathcal{L}_2) \circ \mathrm{marg}(\mathcal{L}_1)(\mathcal{M}) = \mathrm{marg}(\mathcal{L}_1 \cup \mathcal{L}_2)(\mathcal{M})$.*

In this proposition $\mathcal{L}_1$ and $\mathcal{L}_2$ have to be disjoint, since marginalizing first over $\mathcal{L}_1$ gives a marginal SCM $\mathcal{M}_{\mathrm{marg}(\mathcal{L}_1)}$ with endogenous variables $\mathcal{I} \setminus \mathcal{L}_1$.

Next we show that the distributions of a marginal SCM are identical to the marginal distributions induced by the original SCM. A simple proof of this result proceeds by showing that both the intervention and the twin operation commute with marginalization.

PROPOSITION 5.5. *Let $\mathcal{M}$ be an SCM that is uniquely solvable w.r.t. a subset $\mathcal{L} \subseteq \mathcal{I}$. Then, the marginalization $\mathrm{marg}(\mathcal{L})$ commutes with both*

1. *the perfect intervention $\mathrm{do}(I, \boldsymbol{\xi}_I)$ for a subset $I \subseteq \mathcal{I} \setminus \mathcal{L}$ and a value $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$, i.e., $(\mathrm{marg}(\mathcal{L}) \circ \mathrm{do}(I, \boldsymbol{\xi}_I))(\mathcal{M}) = (\mathrm{do}(I, \boldsymbol{\xi}) \circ \mathrm{marg}(\mathcal{L}))(\mathcal{M})$, and*
2. *the twin operation $\mathrm{twin}$, i.e., $(\mathrm{marg}(\mathcal{L} \cup \mathcal{L}') \circ \mathrm{twin})(\mathcal{M}) = (\mathrm{twin} \circ \mathrm{marg}(\mathcal{L}))(\mathcal{M})$,*

*where $\mathcal{L}'$ is the copy of $\mathcal{L}$ in $\mathcal{I}'$.*

With Proposition 5.5 at hand we can prove the main result of this subsection.

THEOREM 5.6 (Marginalization of an SCM preserves the observational, causal and counterfactual semantics). *Let $\mathcal{M}$ be an SCM that is uniquely solvable w.r.t. a subset $\mathcal{L} \subseteq \mathcal{I}$. Then $\mathcal{M}$ and $\mathrm{marg}(\mathcal{L})(\mathcal{M})$ are observationally, interventionally and counterfactually equivalent w.r.t. $\mathcal{I} \setminus \mathcal{L}$.*

This shows that our definition of marginalization (Definition 5.3) preserves the probabilistic, causal and counterfactual semantics, under a certain local unique solvability condition. Moreover, this allows us to marginalize SCMs w.r.t. a certain subset that do not satisfy the additional assumptions imposed by modular SCMs, e.g., the SCM $\mathcal{M}$ of Example 3.14 does not have any additional structure of a compatible system of solution functions, but $\mathcal{M}$ can be marginalized w.r.t. the subset $\{2, 3\}$ (see Appendix A.4).

As we saw in Example 4.7 it is generally not true that interventional equivalence implies counterfactual equivalence. However, for our definition of marginalization we arrive at a marginal SCM that is not only interventionally equivalent, but also counterfactually equivalent w.r.t. the margin.

For an SCM $\mathcal{M}$, unique solvability w.r.t. a certain subset $\mathcal{L} \subseteq \mathcal{I}$ is a sufficient, but not a necessary condition for the existence of an SCM $\tilde{\mathcal{M}}$ on the margin $\mathcal{I} \setminus \mathcal{L}$ such that $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are counterfactually equivalent w.r.t. $\mathcal{I} \setminus \mathcal{L}$. This is illustrated by the following example.

EXAMPLE 5.7 (Marginalization condition of an SCM is not a necessary condition). *Consider the SCM $\mathcal{M} = \langle \boldsymbol{4}, \boldsymbol{1}, \mathbb{R}^4, \mathbb{R}, \boldsymbol{f}, \mathbb{P}_{\mathbb{R}} \rangle$ with causal mechanism given by*

$$f_1(\boldsymbol{x}, e) = e, \quad f_2(\boldsymbol{x}, e) = x_1, \quad f_3(\boldsymbol{x}, e) = x_2, \quad f_4(\boldsymbol{x}, e) = x_4$$

*and $\mathbb{P}_{\mathbb{R}}$ is the standard-normal measure on $\mathbb{R}$. This SCM is solvable w.r.t. $\mathcal{L} = \{2, 4\}$, but not uniquely solvable w.r.t. $\mathcal{L}$, and hence we cannot apply Definition 5.3 to $\mathcal{L}$. However, the SCM $\tilde{\mathcal{M}}$ on the endogenous variables $\{1, 3\}$ with the causal mechanism $\tilde{\boldsymbol{f}}$ given by $\tilde{f}_1(\boldsymbol{x}, e) = e$ and $\tilde{f}_3(\boldsymbol{x}, e) = x_1$ is counterfactually equivalent to $\mathcal{M}$ w.r.t. $\{1, 3\}$, which can be checked easily.*

Hence, in certain cases it may be possible to relax the uniqueness condition.

5.2. *Marginalization of a graph*    We now turn to a marginalization operation for directed mixed graphs, which we call the latent projection. This name is inspired from a similar construction on directed mixed graphs in [78]. In [78], the authors concentrate on a mapping between directed mixed graphs and show that it preserves conditional independence properties [see also 76]. In this subsection, we provide a sufficient condition for the marginalization of an SCM to respect the latent projection, i.e., that the augmented graph of the marginal SCM is a subgraph of the latent projection of the augmented graph of the original SCM.

DEFINITION 5.8 (Marginalization of a directed mixed graph).    *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph and $\mathcal{L} \subseteq \mathcal{V}$ a subset. The* marginalization *of $\mathcal{G}$ w.r.t. $\mathcal{L}$ or the* latent projection *of $\mathcal{G}$ onto $\mathcal{V} \setminus \mathcal{L}$ maps $\mathcal{G}$ to the* marginal graph $\mathrm{marg}(\mathcal{L})(\mathcal{G}) := (\tilde{\mathcal{V}}, \tilde{\mathcal{E}}, \tilde{\mathcal{B}})$, *where*

1. $\tilde{\mathcal{V}} = \mathcal{V} \setminus \mathcal{L}$,
2. $i \to j \in \tilde{\mathcal{E}}$ *if and only if there exists a directed path* $i \to \ell_1 \to \cdots \to \ell_n \to j$ *in $\mathcal{G}$ with $n \geq 0$ and $\ell_1, \ldots, \ell_n \in \mathcal{L}$,*
3. $i \leftrightarrow j \in \tilde{\mathcal{B}}$ *if and only if*
   a) *there exist* $n, m \geq 0$, $\ell_1, \ldots, \ell_n \in \mathcal{L}$, $\tilde{\ell}_1, \ldots, \tilde{\ell}_m \in \mathcal{L}$ *such that* $i \leftarrow l_1 \leftarrow l_2 \leftarrow \cdots \leftarrow \ell_n \leftrightarrow \tilde{\ell}_m \to \tilde{\ell}_{m-1} \to \cdots \to \tilde{\ell}_1 \to j$ *in $\mathcal{G}$, or*
   b) *there exist* $n, m \geq 1$, $\ell_1, \ldots, \ell_n \in \mathcal{L}$, $\tilde{\ell}_1, \ldots, \tilde{\ell}_m \in \mathcal{L}$ *such that* $i \leftarrow l_1 \leftarrow l_2 \leftarrow \cdots \leftarrow \ell_n$ *and* $\tilde{\ell}_m \to \tilde{\ell}_{m-1} \to \cdots \to \tilde{\ell}_1 \to j$ *in $\mathcal{G}$ and $\ell_n = \tilde{\ell}_m$.*

Note that this gives $\mathcal{G}(\mathcal{M}) = \mathrm{marg}(\mathcal{J})(\mathcal{G}^a(\mathcal{M}))$ for any SCM $\mathcal{M}$. Further, for a subgraph $\mathcal{H} \subseteq \mathcal{G}$ we have $\mathrm{marg}(\mathcal{L})(\mathcal{H}) \subseteq \mathrm{marg}(\mathcal{L})(\mathcal{G})$ for any subset of nodes $\mathcal{L}$. It does not matter in which order we project out the nodes or if we perform several projections at once.

PROPOSITION 5.9.    *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph and $\mathcal{L}_1, \mathcal{L}_2 \subseteq \mathcal{V}$ two disjoint subsets. Then* $(\mathrm{marg}(\mathcal{L}_1) \circ \mathrm{marg}(\mathcal{L}_2))(\mathcal{G}) = (\mathrm{marg}(\mathcal{L}_2) \circ \mathrm{marg}(\mathcal{L}_1))(\mathcal{G}) = \mathrm{marg}(\mathcal{L}_1 \cup \mathcal{L}_2)(\mathcal{G})$.

Similarly to the definition of marginalization for SCMs this definition of the latent projection commutes with both the (graphical) perfect intervention and the twin operation.

PROPOSITION 5.10.    *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph and $\mathcal{L}, \mathcal{I}, I \subseteq \mathcal{V}$ subsets. Then, the marginalization $\mathrm{marg}(\mathcal{L})$ commutes with both*

1. *perfect intervention $\mathrm{do}(I)$ if $I$ is disjoint from $\mathcal{L}$, i.e.,* $(\mathrm{marg}(\mathcal{L}) \circ \mathrm{do}(I))(\mathcal{G}) = (\mathrm{do}(I) \circ \mathrm{marg}(\mathcal{L}))(\mathcal{G})$, *and*
2. *the twin operation $\mathrm{twin}(\mathcal{I})$ if $\mathcal{B} = \emptyset$, $\mathcal{J} := \mathcal{V} \setminus \mathcal{I}$ is exogenous (i.e., $\mathrm{pa}_{\mathcal{G}}(\mathcal{J}) = \emptyset$) and $\mathcal{L} \subseteq \mathcal{I}$, i.e.,* $(\mathrm{marg}(\mathcal{L} \cup \mathcal{L}') \circ \mathrm{twin}(\mathcal{I}))(\mathcal{G}) = (\mathrm{twin}(\mathcal{I} \setminus \mathcal{L}) \circ \mathrm{marg}(\mathcal{L}))(\mathcal{G})$,

*where $\mathcal{L}'$ is the copy of $\mathcal{L}$ in $\mathcal{I}'$.*

In Example 5.2 we already saw an example of a marginalization that respects the latent projection. However, not all marginalizations respect the latent projection, as is illustrated in the following example.

EXAMPLE 5.11 (Marginalization does not respect the latent projection).    *Consider the SCM $\mathcal{M}$ of Example 3.14. Although $\mathcal{M}$ and its marginalization $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ with $\mathcal{L} = \{2, 3\}$ are interventionally equivalent w.r.t. $\mathcal{I} \setminus \mathcal{L} = \{1, 4\}$, the graph $\mathcal{G}(\mathcal{M}_{\mathrm{marg}(\mathcal{L})})$ is not a subgraph of the latent projection of $\mathcal{G}(\mathcal{M})$ onto $\mathcal{I} \setminus \mathcal{L}$, as can be verified from the graphs depicted in Figure 5.*
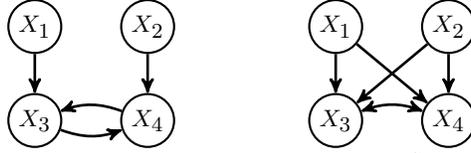
Fig 6: *The graphs of the observationally equivalent SCMs $\mathcal{M}$ (left) and $\tilde{\mathcal{M}}$ (right) of Example 6.1 and 6.2.*

Under the local ancestral unique solvability condition, which is a stronger condition than the local unique solvability condition (i.e., ancestral unique solvability w.r.t. a subset implies unique solvability w.r.t. that subset), one can prove that the marginalization of an SCM respects the latent projection.

PROPOSITION 5.12. *Let $\mathcal{M}$ be an SCM that is ancestrally uniquely solvable w.r.t. a subset $\mathcal{L} \subseteq \mathcal{I}$. Then $\big(\mathcal{G}^a \circ \mathrm{marg}(\mathcal{L})\big)(\mathcal{M}) \subseteq \big(\mathrm{marg}(\mathcal{L}) \circ \mathcal{G}^a\big)(\mathcal{M})$ and $\big(\mathcal{G} \circ \mathrm{marg}(\mathcal{L})\big)(\mathcal{M}) \subseteq \big(\mathrm{marg}(\mathcal{L}) \circ \mathcal{G}\big)(\mathcal{M})$.*

The following example illustrates why the (augmented) graph of a marginalized SCM can be a strict subgraph of the corresponding latent projection.

EXAMPLE 5.13 (Graph of the marginal SCM is a strict subgraph of the latent projection). *Consider the SCM $\mathcal{M} = \langle \mathbf{3}, \mathbf{1}, \mathbb{R}^3, \mathbb{R}, \boldsymbol{f}, \mathbb{P}_{\mathbb{R}} \rangle$ with causal mechanism given by*

$$f_1(\boldsymbol{x}, \boldsymbol{e}) = e_1, \quad f_2(\boldsymbol{x}, \boldsymbol{e}) = x_1 - x_3, \quad f_3(\boldsymbol{x}, \boldsymbol{e}) = x_1$$

*and take for $\mathbb{P}_{\mathbb{R}}$ the standard-normal measure on $\mathbb{R}$. In contrast, to the (augmented) graph of $\mathcal{M}$, there is no directed path in the (augmented) graph of the marginal SCM $\mathcal{M}_{\mathrm{marg}(\{3\})}$.*

For acyclic SCMs we recover with Proposition 5.12 the known result that this class is closed under marginalization (see Proposition 3.6) [15]. For linear SCMs we have that unique solvability w.r.t. a subset $\mathcal{L}$ holds if and only if ancestral unique solvability w.r.t. $\mathcal{L}$ holds (see Proposition C.4), and hence, a marginalization of a linear SCM always respects the latent projection.

**6. Markov properties** In this section we give a short overview of Markov properties for SCMs with cycles. We make use of the Markov properties that were recently developed by Forré and Mooij [18] for HEDGes, a graphical representation that is similar to the augmented graph of SCMs. We briefly summarize some of their main results and apply them to the class of SCMs. In Appendix A.2 we provide a more thorough introduction and give an intuitive derivation which can act as an entry point for the reader into the more extensive discussion of Markov properties provided in [18].

Markov properties associate a set of conditional independence relations to a graph. The directed global Markov property for directed acyclic graphs (see Definition A.4 and A.6), also known as the $d$-separation criterion [50], is one of the most widely used. It directly extends to a similar property for acyclic directed mixed graphs (ADMGs) [60]. It does not hold in general for cyclic SCMs, however, as was already observed earlier [71, 72].

EXAMPLE 6.1 (Directed global Markov property does not hold for cyclic SCM). *Consider the SCM $\mathcal{M} = \langle \mathbf{4}, \mathbf{4}, \mathbb{R}^4, \mathbb{R}^4, \boldsymbol{f}, \mathbb{P}_{\mathbb{R}^4} \rangle$ with causal mechanism given by*

$$f_1(\boldsymbol{x}, \boldsymbol{e}) = e_1, \quad f_2(\boldsymbol{x}, \boldsymbol{e}) = e_2, \quad f_3(\boldsymbol{x}, \boldsymbol{e}) = x_1 x_4 + e_3, \quad f_4(\boldsymbol{x}, \boldsymbol{e}) = x_2 x_3 + e_4$$

*and $\mathbb{P}_{\mathbb{R}^4}$ is the standard-normal distribution on $\mathbb{R}^4$. The graph of $\mathcal{M}$ is depicted in Figure 6 on the left. The model is uniquely solvable w.r.t. every subset. One can check that for every*

*solution $\boldsymbol{X}$ of $\mathcal{M}$, $X_1$ is not independent of $X_2$ given $\{X_3, X_4\}$. However, the variables $X_1$ and $X_2$ are d-separated given $\{X_3, X_4\}$ in $\mathcal{G}(\mathcal{M})$. Hence the global directed Markov property does not hold here.*

Although some progress has been made in the case of discrete [52, 49, 18] and linear models [70, 71, 72, 63, 31, 27, 18], only recently a general directed global Markov property has been introduced for more general cyclic models [18], that is based on $\sigma$-separation (see Definition A.16 and A.20), an extension of $d$-separation. This notion of $\sigma$-separation was derived from the notion of $d$-separation in the acyclification of the graph [18] (see Definition A.13). The acyclification of a graph generalizes the idea of the collapsed graph developed by Spirtes [71] and can, in particular, be applied to the graphs of SCMs. The main idea of the acyclification is that under the condition that the SCM is uniquely solvable w.r.t. each strongly connected component, we can replace the causal mechanisms of these strongly connected components by their measurable solution functions, which results in an acyclic SCM. This acyclified SCM (see Definition A.11) is observationally equivalent to the original SCM (see Proposition A.12).

EXAMPLE 6.2 (Construction of an observationally equivalent acyclic SCM). *Consider the SCM $\mathcal{M}$ of Example 6.1 which is uniquely solvable w.r.t. all its strongly connected components, i.e., the subsets $\{1\}$, $\{2\}$ and $\{3, 4\}$. Replacing the causal mechanisms of these strongly connected components by their measurable solution functions gives the SCM $\tilde{\mathcal{M}}$ that is the same as $\mathcal{M}$ except that its causal mechanism $\tilde{\boldsymbol{f}}$ is given by*

$$\tilde{f}_1(\boldsymbol{x}, \boldsymbol{e}) := e_1, \quad \tilde{f}_2(\boldsymbol{x}, \boldsymbol{e}) := e_2, \quad \tilde{f}_3(\boldsymbol{x}, \boldsymbol{e}) := \tfrac{x_1 e_4 + e_3}{1 - x_1 x_2}, \quad \tilde{f}_4(\boldsymbol{x}, \boldsymbol{e}) := \tfrac{x_2 e_3 + e_4}{1 - x_1 x_2}.$$

*By construction, $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are observationally equivalent. Because $\tilde{\mathcal{M}}$ is acyclic (see Figure 6 on the right) we can apply the directed global Markov property to $\tilde{\mathcal{M}}$. The fact that $X_1$ and $X_2$ are not d-separated given $\{X_3, X_4\}$ in $\mathcal{G}(\tilde{\mathcal{M}})$ is in line with $X_1$ being dependent of $X_2$ given $\{X_3, X_4\}$ for every solution $\boldsymbol{X}$ of $\tilde{\mathcal{M}}$ (and hence of $\mathcal{M}$).*

This acyclification preserves solutions, and $d$-separation in the acyclification can directly be translated into $\sigma$-separation on the original graph (see Proposition A.19). This leads to the general directed global Markov property. The following theorem summarizes the main results of [18] applied to SCMs.

THEOREM 6.3 (Global Markov properties for SCMs [18]). *Let $\mathcal{M}$ be a uniquely solvable SCM. Then its observational distribution $\mathbb{P}^{\boldsymbol{X}}$ exists, is unique and the following two statements hold.*

1. $\mathbb{P}^{\boldsymbol{X}}$ *satisfies the* directed global Markov property *("d-separation criterion") relative to $\mathcal{G}(\mathcal{M})$ (see Definition A.6) if $\mathcal{M}$ satisfies at least one of the following conditions:*
   a) *$\mathcal{M}$ is acyclic;*
   b) *all endogenous spaces $\mathcal{X}_i$ are discrete and $\mathcal{M}$ is ancestrally uniquely solvable;*
   c) *$\mathcal{M}$ is linear (see Definition C.1), each of its causal mechanisms $\{f_i\}_{i \in \mathcal{I}}$ has a non-trivial dependence on at least one exogenous variable, and $\mathbb{P}_{\boldsymbol{\varepsilon}}$ has a density w.r.t. the Lebesgue measure on $\mathbb{R}^{\mathcal{J}}$.*
2. $\mathbb{P}^{\boldsymbol{X}}$ *satisfies the* general directed global Markov property *("$\sigma$-separation criterion") relative to $\mathcal{G}(\mathcal{M})$ (see Definition A.20) if $\mathcal{M}$ is uniquely solvable w.r.t. each strongly connected component of $\mathcal{G}(\mathcal{M})$.*[14]

---

[14]Since [18] also provides results under the weaker condition that an SCM is solvable (not necessarily uniquely) w.r.t. each strongly connected component of $\mathcal{G}(\mathcal{M})$, one might believe that Theorem 6.3.(2) could

The general directed global Markov property is generally weaker than the directed global Markov property, since $\sigma$-separation implies $d$-separation. The acyclic case is well-known and was first shown in the context of linear-Gaussian structural equation models [75, 32]. The discrete case fixes the erroneous theorem by Pearl and Dechter [52], for which a counterexample was found by Neal [49], by adding the ancestral unique solvability condition, and extends it to allow for bidirected edges in the graph. The linear case is an extension of existing results for the linear-Gaussian setting without bidirected edges [71, 72, 31] to a linear (possibly non-Gaussian) setting with bidirected edges in the graph.

In constraint-based approaches to causal discovery, one usually assumes the converse of the (general) directed global Markov property to hold [73, 51], which is called $\sigma$-faithfulness respectively $d$-faithfulness (see Definition A.9 and A.23). Meek [41] showed that for multinomial and linear-Gaussian DAG (i.e., acyclic and causally sufficient SCMs) models, $d$-faithfulness holds for all parameter values up to a measure zero set. Up to our knowledge no such result has been shown for any subclass of SCMs that contains cycles, nor in more general acyclic settings.

**7. Causal interpretation of the graph of SCMs**   In Examples 4.4 and 4.8 we already saw that sometimes no information in the observational, interventional and even the counterfactual distributions suffices to decide whether a directed path or bidirected edge is present in the graph, or not. Here, we do not attempt to provide a complete characterization of all the conditions under which the presence or absence of a directed path or bidirected edge in the graph can be identified from the observational and interventional distributions. Instead, we give some sufficient conditions under which one can detect a directed path and bidirected edge in the graph.

In general, cyclic SCMs may have none, one or multiple induced observational distributions, and this may change after intervening in the system. Here, we restrict ourselves to graphs of SCMs where the induced (marginal) observational and interventional distributions are uniquely defined.

7.1. *Directed paths and edges*   For cyclic SCMs the causal interpretation of the SCM is not always consistent with its graph. This can be illustrated with the SCM $\mathcal{M}$ of Example 5.11. Here, one sees a difference in the marginal distribution $\mathbb{P}_{\mathcal{M}_{\mathrm{do}(\{1\},\xi_1)}}$ on $\mathcal{X}_4$ for different values of $\xi_1$, although variable 1 is not an ancestor of variable 4 and each marginal distribution $\mathbb{P}_{\mathcal{M}_{\mathrm{do}(\{1\},\xi_1)}}$ on $\mathcal{X}_4$ is uniquely defined. This counterintuitive behavior that an intervention on a non-ancestor of a variable can change the distribution of that variable was already observed by Neal [49]. However, under a specific unique solvability condition, we obtain a direct causal interpretation for the absence of a directed edge or directed path in the graph of an SCM.

PROPOSITION 7.1 (Sufficient condition for detecting a directed edge in the latent projection of the graph of an SCM).   *Consider an SCM $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$, a subset $\mathcal{O} \subseteq \mathcal{I}$ and $i, j \in \mathcal{O}$ such that $i \neq j$. Let $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$, where $I := \mathcal{O} \setminus \{i, j\}$, such that $\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$ is uniquely solvable w.r.t. $\mathrm{an}_{\mathcal{G}(\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}) \setminus i}(j)$. If there exist values $\xi_i \neq \tilde{\xi}_i \in \mathcal{X}_i$ such that*

---

be generalized to stating that in that case, any of its observational distributions satisfies the general directed global Markov property. However, that is not true: consider for example the SCM $\mathcal{M} = \langle \mathbf{2}, \emptyset, \mathbb{R}^2, \mathbf{1}, \boldsymbol{f}, \mathbb{P}_{\mathbf{1}} \rangle$ with $f_1(\boldsymbol{x}) = x_1$ and $f_2(\boldsymbol{x}) = x_2$. Then $\mathcal{M}$ is solvable w.r.t. each of its strongly connected components $\{1\}$ and $\{2\}$. The solution with $X_1 = X_2$ shows a dependence between $X_1$ and $X_2$ and thus $X_1 \perp\!\!\!\perp X_2$ does not hold. In general, all strongly connected components that admit multiple solutions may be dependent on any other variable(s) in the model.

*both* $(\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)})_{\mathrm{do}(\{i\},\xi_i)}$ *and* $(\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)})_{\mathrm{do}(\{i\},\tilde{\xi}_i)}$ *induce unique marginal distributions on* $\mathcal{X}_j$, *and these two induced distributions do not coincide, i.e., there exists a measurable set* $\mathcal{B}_j \subseteq \mathcal{X}_j$ *such that*

$$\mathbb{P}_{(\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)})_{\mathrm{do}(\{i\},\xi_i)}}(X_j \in \mathcal{B}_j) \neq \mathbb{P}_{(\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)})_{\mathrm{do}(\{i\},\tilde{\xi}_i)}}(X_j \in \mathcal{B}_j),$$

*then there exists a directed edge* $i \to j$ *in the latent projection* $\mathrm{marg}(\mathcal{I} \setminus \mathcal{O})(\mathcal{G}(\mathcal{M}))$ *of* $\mathcal{G}(\mathcal{M})$ *on* $\mathcal{O}$.

Two cases are of special interest: $\mathcal{O} = \mathcal{I}$, which corresponds with a directed edge $i \to j$ in $\mathcal{G}(\mathcal{M})$, and $\mathcal{O} = \{i, j\}$, which corresponds with a directed path $i \to \cdots \to j$ in $\mathcal{G}(\mathcal{M})$.

The condition in Proposition 7.1 is a sufficient condition for determining whether a directed edge or path is present in the graph. In general, not all directed edges and paths can be identified from the interventional distributions with this sufficient condition. For example, no interventional distribution satisfies the condition of Proposition 7.1 for the SCM $\mathcal{M}$ in Example 4.4, although there is a directed edge $1 \to 2$ in the graph $\mathcal{G}(\mathcal{M})$.

7.2. *Bidirected edges*  It is well-known that there exists a similar sufficient condition for detecting bidirected edges in the graph of an acyclic SCM also known as the common-cause principle [see e.g., 51]. In the two variables case, this criterion informally states that there exists a bidirected edge between the variables $i$ and $j$ in the graph of the SCM, if the marginal interventional distribution of $X_j$ under the intervention $\mathrm{do}(\{i\}, x_i)$ differs from the conditional distribution of $X_j$ given $X_i = x_i$.

EXAMPLE 7.2 (Detecting a bidirected edge in the graph of an SCM).  *Consider the acyclic SCM* $\mathcal{M}$ *of Example 4.4 and the SCM* $\tilde{\mathcal{M}}$ *that is the same as* $\mathcal{M}$ *except for its causal mechanism, which is given by* $\tilde{f}_1(\boldsymbol{x}, \boldsymbol{e}) = e_1$ *and* $\tilde{f}_2(\boldsymbol{x}, \boldsymbol{e}) = x_1 e_1$. *For the SCM* $\tilde{\mathcal{M}}$ *we observe that the marginal interventional distribution* $\mathbb{P}_{\tilde{\mathcal{M}}_{\mathrm{do}(\{1\},\xi_1)}}(X_2 = -1)$ *is not equal to the conditional distribution* $\mathbb{P}_{\tilde{\mathcal{M}}}(X_2 = -1 \mid X_1 = \xi_1)$ *for both* $\xi_1 = -1$ *and* $\xi_1 = 1$. *This observation suffices to identify the presence of the bidirected edge* $1 \leftrightarrow 2$ *in the graph* $\mathcal{G}(\tilde{\mathcal{M}})$. *For the SCM* $\mathcal{M}$, *whose graph does not contain the bidirected edge* $1 \leftrightarrow 2$, *the marginal interventional distribution and conditional distribution coincide.*

The following proposition provides a generalization of this sufficient condition for detecting bidirected edges to graphs of SCMs that may include cycles.

PROPOSITION 7.3 (Sufficient condition for detecting a bidirected edge in the latent projection of the graph of an SCM).  *Consider an SCM* $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$, *a subset* $\mathcal{O} \subseteq \mathcal{I}$ *and* $i, j \in \mathcal{O}$ *such that* $i \neq j$. *Let* $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$, *where* $I := \mathcal{O} \setminus \{i, j\}$, *such that* $\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}$ *is uniquely solvable w.r.t. both* $\mathrm{an}_{\mathcal{G}(\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)})}(i)$ *and* $\mathrm{an}_{\mathcal{G}(\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)})\setminus i}(j)$. *Assume that for every* $\xi_i \in \mathcal{X}_i$ *both* $\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}$ *and* $(\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)})_{\mathrm{do}(\{i\},\xi_i)}$ *induce a unique marginal distribution on* $\mathcal{X}_j \times \mathcal{X}_i$ *and* $\mathcal{X}_j$ *respectively. If* $j \notin \mathrm{an}_{\mathcal{G}(\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)})}(i)$ *and there exists a measurable set* $\mathcal{B}_j \subseteq \mathcal{X}_j$ *such that for every version of the regular conditional probability* $\mathbb{P}_{\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}}(X_j \in \mathcal{B}_j \mid X_i = \xi_i)$ *there exists a value* $\xi_i \in \mathcal{X}_i$ *such that*

$$\mathbb{P}_{\left(\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}\right)_{\mathrm{do}(\{i\},\xi_i)}}(X_j \in \mathcal{B}_j) \neq \mathbb{P}_{\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}}(X_j \in \mathcal{B}_j \mid X_i = \xi_i),$$

*then there exists a bidirected edge* $i \leftrightarrow j$ *in the latent projection* $\mathrm{marg}(\mathcal{I} \setminus \mathcal{O})(\mathcal{G}(\mathcal{M}))$ *of* $\mathcal{G}(\mathcal{M})$ *on* $\mathcal{O}$.

This proposition gives a sufficient condition for determining that a bidirected edge is present in the graph. In general, not all bidirected edges in the graph can be identified from the observational, interventional, and even the counterfactual distributions, as we saw in Example 4.8. There, we saw that for the SCM $\tilde{\mathcal{M}}$, there exists a bidirected edge $1 \leftrightarrow 2 \in \mathcal{G}(\tilde{\mathcal{M}})$ while the density $p(x_2 \mid \mathrm{do}(X_1 = x_1)) = p(x_2 \mid X_1 = x_1)$ for all $x_1 \in \mathcal{X}_1$. For the acyclic setting, the above criterion is generally considered as a universal way to detect a confounder (note that then one can also deal with the case $j \in \mathrm{an}_{\mathcal{G}(\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)})}(i)$ by swapping the roles of $i$ and $j$). If $i$ and $j$ are part of a cycle, the above sufficient condition cannot be applied, and in that case, to the best of our knowledge, no simple sufficient conditions for detecting the presence of a bidirected edge are known.

**8. Simple SCMs**   In this section we introduce the well-behaved class of simple SCMs. Simple SCMs satisfy all the local unique solvability conditions to ensure that this class is closed under both perfect intervention and marginalization. They extend the subclass of acyclic SCMs to the cyclic setting, while preserving many of their convenient properties.

DEFINITION 8.1 (Simple SCM).   *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM. We call $\mathcal{M}$ simple if it is uniquely solvable w.r.t. every subset $\mathcal{O} \subseteq \mathcal{I}$.*

Loosely speaking, an SCM is simple if any subset of its structural equations can be solved uniquely for its associated variables in terms of the other variables that appear in these equations. An example of a simple SCM is given in Example D.1.

On simple SCMs one can perform any number of marginalizations (see Definition 5.3) in any order (see Proposition 5.4). All these marginalizations respect the latent projection (see Proposition 5.12) and each resulting marginal SCM is again simple. Moreover, we show that this class is closed under intervention and the twin operation.

PROPOSITION 8.2.   *The class of simple SCMs is closed under marginalization, perfect intervention and the twin operation.*

The class of simple SCMs contains the acyclic SCMs as a subclass (see Proposition 3.6). In particular, a simple SCM has no self-cycles (see Proposition 3.9), since a self-cycle denotes that that variable cannot be uniquely (up to a $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-null set) determined by its parents.

From Proposition 8.2 it follows that the results summarized in Theorem 6.3 also apply to all the observational, interventional and counterfactual distributions of simple SCMs.

COROLLARY 8.3 (Global Markov properties for simple SCMs).   *Let $\mathcal{M}$ be a simple SCM. Then, the*

1. *observational distribution,*
2. *interventional distribution after perfect intervention on $I \subset \mathcal{I}$,*
3. *counterfactual distribution after perfect intervention on $\tilde{I} \subseteq \mathcal{I} \cup \mathcal{I}'$,*

*all exist, are unique and satisfy the general directed global Markov property relative to $\mathcal{G}(\mathcal{M})$, $\mathrm{do}(I)(\mathcal{G}(\mathcal{M}))$ and $\mathrm{do}(\tilde{I})(\mathrm{twin}(\mathcal{G}(\mathcal{M})))$ respectively. Moreover, if $\mathcal{M}$ satisfies at least one of the three conditions (1a), (1b), (1c) of Theorem 6.3, then they also obey the directed global Markov property relative to $\mathcal{G}(\mathcal{M})$, $\mathrm{do}(I)(\mathcal{G}(\mathcal{M}))$ and $\mathrm{do}(\tilde{I})(\mathrm{twin}(\mathcal{G}(\mathcal{M})))$ respectively.*

Many of these properties are also shown to hold for the class of *modular SCMs* [18], which contains, in particular, the class of simple SCMs (see Appendix A.4 for more details).

Moreover, simple SCMs satisfy the unique solvability conditions of Proposition 7.1 and 7.3, which allows us to define the causal relationships for simple SCMs in terms of its graph.

DEFINITION 8.4 (Causal relationships for simple SCMs).    *Let $\mathcal{M}$ be a simple SCM.*

1. *If there exists a directed edge $i \to j \in \mathcal{G}(\mathcal{M})$, i.e., $i \in \mathrm{pa}(j)$, then we call $i$* a direct cause *of $j$ according to $\mathcal{M}$;*
2. *If there exists a directed path $i \to \cdots \to j$ in $\mathcal{G}(\mathcal{M})$, i.e., $i \in \mathrm{an}(j)$, then we call $i$* a cause *of $j$ according to $\mathcal{M}$;*
3. *If there exists a bidirected edge $i \leftrightarrow j \in \mathcal{G}(\mathcal{M})$, then we call $i$ and $j$* (latently) confounded *according to $\mathcal{M}$.*

Corollary A.24 summarizes sufficient conditions for determining the different causal relationships according to a simple SCM $\mathcal{M}$. For simple SCMs it is in general not possible to identify all the causal relationships in the graph from the observational, interventional, or even the counterfactual distributions. Example 4.4 and 4.8 show that this is already impossible for acyclic SCMs without further assumptions.

Finally, there is a connection between SCMs and potential outcomes [68] that generalizes to the cyclic setting. One of the consequences of Proposition 8.2 is that all counterfactuals are defined for a simple SCM (even if it is cyclic). This allows us to define potential outcomes in terms of a simple SCM in the following way.

DEFINITION 8.5 (Potential outcome).    *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \mathcal{X}, \mathcal{E}, \boldsymbol{f}, \mathbb{P}_{\mathcal{E}} \rangle$ be a simple SCM, $I \subseteq \mathcal{I}$ a subset, $\boldsymbol{\xi}_I \in \mathcal{X}_I$ a value and $\boldsymbol{E}$ a random variable such that $\mathbb{P}^{\boldsymbol{E}} = \mathbb{P}_{\mathcal{E}}$. The* potential outcome *under the perfect intervention $\mathrm{do}(I, \boldsymbol{\xi}_I)$ is defined as $\boldsymbol{X}_{\boldsymbol{\xi}_I} := \boldsymbol{g}_{\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}}(\boldsymbol{E}_{\mathrm{pa}(\mathcal{I})})$, where $\boldsymbol{g}_{\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}} : \mathcal{E}_{\mathrm{pa}(\mathcal{I})} \to \mathcal{X}$ is a measurable solution function for $\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}$.*

**9. Discussion**    In this paper, we studied the basic properties of SCMs in the presence of cycles and latent variables without restricting to linear functional relationships between the variables. We saw that cyclic SCMs behave differently in many aspects than acyclic SCMs. Indeed, in the presence of cycles, many of the convenient properties of acyclic SCMs do not hold in general: SCMs do not always have a solution; they do not always induce unique observational, interventional and counterfactual distributions; a marginalization does not always exist, and if it exists the marginal model does not always respect the latent projection; they do not always satisfy a Markov property; and their graphs are not always consistent with their causal semantics.

We introduced various notions of (unique) solvability and showed that under appropriate (unique) solvability conditions, many of the operations and results for the acyclic setting can be extended to SCMs with cycles. For example, we introduced several equivalence relations between SCMs to compare SCMs at different levels of abstraction, we showed how to define marginal SCMs on a subset of the variables that are (in various ways) equivalent to the original SCM, we discussed under which conditions the distributions satisfy the (general) directed global Markov property relative to their graphs, and we showed under which conditions the graph of an SCM can be interpreted causally. Most of these results are shown under sufficient conditions that are not necessary (e.g., for the marginalization operation this was shown in Example 5.7). It may therefore be possible to further relax some of the conditions.

These insights led us to introduce the more well-behaved class of simple SCMs, which forms an extension of the class of acyclic SCMs to the cyclic setting that preserves many of its convenient properties: simple SCMs induce unique observational, interventional and counterfactual distributions; the class of simple SCMs is closed under both perfect intervention

and marginalization; the marginalization respects the latent projection; the induced distributions obey the general directed global Markov property and obey the directed global Markov property in the acyclic, discrete and linear case. This class does not contain SCMs that have self-cycles and graphs of simple SCMs have a direct and intuitive causal interpretation.

One key property of simple SCMs is that the solutions always satisfy the conditional independencies implied by $\sigma$-separation. By simply replacing $d$-separation with $\sigma$-separation it turns out that one can directly extend results and algorithms for acyclic SCMs to the more general class of simple SCMs. E.g., adjustment criteria (including the back-door criterion), Pearl's $\mathrm{do}$-calculus and Tian's ID algorithm for the identification of causal effects have been extended recently to the class of modular SCMs, which contains the class of simple SCMs [20]. Several causal discovery algorithms have already been proposed that work with simple SCMs, e.g., the first constraint-based causal discovery algorithm that can deal with cycles and non-linear functional relationships [19]. Also, Local Causal Discovery (LCD) [10], Y-structures [38] and the Joint Causal Inference framework (JCI) all apply to simple SCMs [47] even though they were originally developed for acyclic SCMs only. Recently it has been shown that even the well-known Fast Causal Inference (FCI) algorithm [74, 80] is directly applicable to simple SCMs [44] and provides a consistent estimate of the Markov equivalence class (under the faithfulness assumption). Moreover, a method for constructing non-linear simple SCMs using neural networks and sampling from them has been proposed [19]. This illustrates that the class of simple SCMs forms a convenient and practical extension of the class of acyclic SCMs that can be used for the purposes of causal modeling, reasoning, discovery, and prediction.

We hope that this work will provide the foundations for a general theory of statistical causal modeling with SCMs. Future work might consist of reparametrizing and reducing the space of the exogenous variables of an SCM while preserving the causal and counterfactual semantics; extending and generalizing the identifiability results for (direct) causes and confounders; extending the graphs of SCMs to represent selection bias; proving completeness results for some Markov properties for a subclass of SCMs that contains cycles.

# FOUNDATIONS OF STRUCTURAL CAUSAL MODELS WITH CYCLES AND LATENT VARIABLES

### Supplementary Material

This Supplementary Material contains a summary of the basic terminology and results for causal graphical models (Appendix A), additional (unique) solvability properties (Appendix B), some results for linear SCMs (Appendix C), other examples (Appendix D), the proofs of all the theoretical results (Appendix E) and the measurable selection theorems (Appendix F) that are used in several proofs.

### APPENDIX A: CAUSAL GRAPHICAL MODELS

In this appendix, we provide a summary of the basic terminology and results for causal graphical models. In Appendix A.1 we provide the terminology for directed (mixed) graphs. In Appendix A.2 we give an introduction and an intuitive derivation of Markov properties for SCMs with cycles. In Appendix A.3 we give a summary of useful conditions for detecting various causal relationships for simple SCMs. In Appendix A.4 we provide a definition of modular SCMs and show how they relate to SCMs. In Appendix A.5 we provide an overview of the causal graphical models related to SCMs. The proofs of the theoretical results in this appendix are given in Appendix E.

**A.1. Directed (mixed) graphs**   In this subsection we introduce the terminology for directed (mixed) graphs, where we do allow for cycles [34, 60, 51, 18].

DEFINITION A.1 (Directed (mixed) graph).

1. *A* directed graph *is a pair $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V}$ is a set of nodes and $\mathcal{E}$ is a set of directed edges, which is a subset $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ of ordered pairs of nodes. Each element $(i,j) \in \mathcal{E}$ can be represented by the directed edge $i \to j$ or equivalently $j \leftarrow i$. In particular, $(i,i) \in \mathcal{E}$ represents a* self-cycle $i \to i$.
2. *A* directed mixed graph *is a triple $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$, where the pair $(\mathcal{V}, \mathcal{E})$ forms a directed graph and $\mathcal{B}$ is a set of bidirected edges, which is a subset $\mathcal{B} \subseteq \{\{i,j\} : i,j \in \mathcal{V}, i \neq j\}$ of unordered (distinct) pairs of nodes. Each element $\{i,j\} \in \mathcal{B}$ can be represented by the bidirected edge $i \leftrightarrow j$ or equivalently $j \leftrightarrow i$. Note that a directed graph can be considered as a directed mixed graph without bidirected edges.*
3. *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph. A directed mixed graph $\tilde{\mathcal{G}} = (\tilde{\mathcal{V}}, \tilde{\mathcal{E}}, \tilde{\mathcal{B}})$ is a subgraph of $\mathcal{G}$ if $\tilde{\mathcal{V}} \subseteq \mathcal{V}$, $\tilde{\mathcal{E}} \subseteq \mathcal{E}$ and $\tilde{\mathcal{B}} \subseteq \mathcal{B}$, in which case we write $\tilde{\mathcal{G}} \subseteq \mathcal{G}$. For a subset $\mathcal{W} \subseteq \mathcal{V}$, we define the* induced subgraph *of $\mathcal{G}$ on $\mathcal{W}$ by $\mathcal{G}_{\mathcal{W}} := (\mathcal{W}, \tilde{\mathcal{E}}, \tilde{\mathcal{B}})$, where $\tilde{\mathcal{E}}$ and $\tilde{\mathcal{B}}$ are the set of directed and bidirected edges in $\mathcal{E}$ and $\mathcal{B}$ respectively that lie in $\mathcal{W} \times \mathcal{W}$ and $\{\{i,j\} : i,j \in \mathcal{W}, i \neq j\}$ respectively.*
4. *A* walk *between $i,j \in \mathcal{V}$ in a directed mixed graph $\mathcal{G}$ is a tuple $(i_0, \epsilon_1, i_1, \epsilon_2, i_2, \ldots, \epsilon_n, i_n)$ of alternating nodes and edges in $\mathcal{G}$ for some $n \geq 0$, where all $i_0, \ldots, i_n \in \mathcal{V}$, all $\epsilon_1, \ldots, \epsilon_n \in \mathcal{E} \cup \mathcal{B}$ such that $\epsilon_k \in \{i_{k-1} \to i_k, i_{k-1} \leftarrow i_k, i_{k-1} \leftrightarrow i_k\}$ for all $k = 1, \ldots, n$, and it starts with node $i_0 = i$ and ends with node $i_n = j$. Note that $n = 0$ corresponds with a trivial walk consisting of a single node. If all nodes $i_0, \ldots, i_n$ are distinct, it is called a* path. *A walk (path) of the form $i \to \cdots \to j$, i.e., $\epsilon_k$ is $i_{k-1} \to i_k$ for all $k = 1, 2, \ldots, n$, is called a* directed walk (path) *from $i$ to $j$.*

31

5. *A* cycle *through* $i \in \mathcal{V}$ *in a directed mixed graph* $\mathcal{G}$ *is a directed path from* $i$ *to some node* $j$ *extended with the edge* $j \rightarrow i \in \mathcal{E}$. *In particular, a self-cycle* $i \rightarrow i \in \mathcal{E}$ *is a cycle. Note that a path cannot contain any cycles. A directed graph and a directed mixed graph are said to be* acyclic *if they contain no cycles, and are then referred to as a* directed acyclic graph (DAG) *and an* acyclic directed mixed graph (ADMG), *respectively.*

6. *For a directed mixed graph* $\mathcal{G}$ *and a node* $i \in \mathcal{V}$ *we define the set of* parents *of* $i$ *by* $\mathrm{pa}_{\mathcal{G}}(i) := \{j \in \mathcal{V} : j \rightarrow i \in \mathcal{E}\}$, *the set of children of* $i$ *by* $\mathrm{ch}_{\mathcal{G}}(i) := \{j \in \mathcal{V} : i \rightarrow j \in \mathcal{E}\}$, *the set of* ancestors *of* $i$ *by*

$$\mathrm{an}_{\mathcal{G}}(i) := \{j \in \mathcal{V} : \text{there is a directed path from } j \text{ to } i \text{ in } \mathcal{G}\}$$

*and the set of* descendants *of* $i$ *by*

$$\mathrm{de}_{\mathcal{G}}(i) := \{j \in \mathcal{V} : \text{there is a directed path from } i \text{ to } j \text{ in } \mathcal{G}\}.$$

*Note that we have* $\{i\} \cup \mathrm{pa}_{\mathcal{G}}(i) \subseteq \mathrm{an}_{\mathcal{G}}(i)$ *and* $\{i\} \cup \mathrm{ch}_{\mathcal{G}}(i) \subseteq \mathrm{de}_{\mathcal{G}}(i)$. *We can apply all these definitions to subsets* $\mathcal{U} \subseteq \mathcal{V}$ *by taking unions, for example* $\mathrm{pa}_{\mathcal{G}}(\mathcal{U}) := \cup_{i \in \mathcal{U}} \mathrm{pa}_{\mathcal{G}}(i)$. *A subset* $\mathcal{A} \subseteq \mathcal{V}$ *is called an* ancestral subset *in* $\mathcal{G}$ *if* $\mathcal{A} = \mathrm{an}_{\mathcal{G}}(\mathcal{A})$, *i.e.,* $\mathcal{A}$ *is closed under taking ancestors of* $\mathcal{A}$ *in* $\mathcal{G}$.

7. *Let* $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ *be a directed mixed graph. We call* $\mathcal{G}$ strongly connected *if for every pair of distinct nodes* $i, j \in \mathcal{V}$, *the graph contains a cycle that passes through both* $i$ *and* $j$. *The* strongly connected component *of* $i \in \mathcal{V}$, *denoted by* $\mathrm{sc}_{\mathcal{G}}(i)$, *is the maximal subset* $\mathcal{S} \subseteq \mathcal{V}$ *such that* $i \in \mathcal{S}$ *and the induced subgraph* $\mathcal{G}_{\mathcal{S}}$ *is strongly connected. Equivalently,* $\mathrm{sc}_{\mathcal{G}}(i) = \mathrm{an}_{\mathcal{G}}(i) \cap \mathrm{de}_{\mathcal{G}}(i)$.

8. *A* loop *in a directed mixed graph* $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ *is a subset* $\mathcal{O} \subseteq \mathcal{V}$ *that is strongly connected in the induced subgraph* $\mathcal{G}_{\mathcal{O}}$ *of* $\mathcal{G}$ *on* $\mathcal{O}$.

9. *For a directed graph* $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, *we define the* graph of strongly connected components *of* $\mathcal{G}$ *as the directed graph* $\mathcal{G}^{\mathrm{sc}} := (\mathcal{V}^{\mathrm{sc}}, \mathcal{E}^{\mathrm{sc}})$, *where* $\mathcal{V}^{\mathrm{sc}}$ *are the strongly connected components of* $\mathcal{G}$, *i.e.,* $\mathcal{V}^{\mathrm{sc}}$ *are the equivalence classes in* $\mathcal{V}/\sim$ *with the equivalence relation* $i \sim j$ *if and only if* $i \in \mathrm{sc}_{\mathcal{G}}(j)$, *and* $\mathcal{E}^{\mathrm{sc}} = (\mathcal{E} \setminus \{i \rightarrow i : i \in \mathcal{V}\})/\sim$ *with the equivalence relation* $(i \rightarrow j) \sim (i' \rightarrow j')$ *if and only if* $i \sim i'$ *and* $j \sim j'$.

We omit the subscript $\mathcal{G}$ whenever it is clear which directed (mixed) graph $\mathcal{G}$ we are referring to.

LEMMA A.2 (DAG of strongly connected components). *Let* $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ *be a directed graph. Then* $\mathcal{G}^{\mathrm{sc}}$, *the graph of strongly connected components of* $\mathcal{G}$, *is a DAG.*

**A.2. Markov properties**  In this subsection we give a short overview of Markov properties for SCMs with cycles. We will make use of the Markov properties that were recently developed by Forré and Mooij [18] for HEDGes, a graphical representation that is similar to the augmented graph of SCMs. We briefly summarize some of their main results and apply them to the class of SCMs. We also provide a shorter and more intuitive derivation so that this subsection can act as an entry point for the reader into the more extensive discussion of Markov properties provided in [18].

Markov properties associate a set of conditional independence relations to a graph. The directed global Markov property for directed acyclic graphs, also known as the $d$-separation criterion [50], is one of the most widely used. It directly extends to a similar property for acyclic directed mixed graphs (ADMGs) [60]. It does not hold in general for cyclic SCMs, however, as was already observed earlier [71, 72]. Under some conditions (roughly speaking, linearity or discrete variables) the directed global Markov property can be shown to hold also in the presence of cycles [18].

Inspired by work of Spirtes [71], Forré and Mooij [18] recognized that in the general cyclic case a different extension of $d$-separation, termed $\sigma$-separation, is needed, leading to the general directed global Markov property. One key result in [18] implies that under the assumption of unique solvability w.r.t. each strongly connected component of its graph, the observational distribution of an SCM satisfies the general directed global Markov property w.r.t. its graph. The solvability assumptions are in general not preserved under interventions. Under the stronger assumption of simplicity, however, they are, and one obtains the corollary that also all interventional and counterfactual distributions of a simple SCM satisfy the general directed global Markov property w.r.t. to their corresponding graphs.

For a more extensive study of different Markov properties that can be associated to SCMs we refer the reader to [18].

A.2.1. *The directed global Markov property*   Conditional independencies in the observational distribution of an acyclic SCM can be read off from its graph by using the graphical criterion called $d$-separation [51]. The directed global Markov property associates a conditional independence relation in the observational distribution of the SCM to each $d$-separation entailed by the graph. Here, we use a formulation of $d$-separation that generalizes $d$-separation for DAGs [50] and $m$-separation for ADMGs [60] and mDAGs [15].

DEFINITION A.3 (Collider).   *Let $\pi = (i_0, \epsilon_1, i_1, \epsilon_2, i_2, \ldots, \epsilon_n, i_n)$ be a walk (path) in a directed mixed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$. A node $i_k$ on $\pi$ is called a* collider *on $\pi$ if it is a non-endpoint node ($1 \leq k < n$) and the two edges $\epsilon_k, \epsilon_{k+1}$ meet head-to-head on $i_k$ (i.e., if the subwalk $(i_{k-1}, \epsilon_k, i_k, \epsilon_{k+1}, i_{k+1})$ is of the form $i_{k-1} \to i_k \leftarrow i_{k+1}$, $i_{k-1} \leftrightarrow i_k \leftarrow i_{k+1}$, $i_{k-1} \to i_k \leftrightarrow i_{k+1}$, or $i_{k-1} \leftrightarrow i_k \leftrightarrow i_{k+1}$). The node $i_k$ is called a* non-collider *on $\pi$ otherwise, i.e., if it is an endpoint node ($k = 0$ or $k = n$) or if the subwalk $(i_{k-1}, \epsilon_k, i_k, \epsilon_{k+1}, i_{k+1})$ is of the form $i_{k-1} \to i_k \to i_{k+1}$, $i_{k-1} \leftarrow i_k \leftarrow i_{k+1}$, $i_{k-1} \leftarrow i_k \to i_{k+1}$, $i_{k-1} \leftrightarrow i_k \to i_{k+1}$, or $i_{k-1} \leftarrow i_k \leftrightarrow i_{k+1}$.*

Note in particular that the end points of a walk are non-colliders on the walk.

DEFINITION A.4 ($d$-separation).   *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph and let $C \subseteq \mathcal{V}$ be a subset of nodes. A walk (path) $\pi = (i_0, \epsilon_1, i_1, \ldots, i_n)$ in $\mathcal{G}$ is said to be $C$-$d$-blocked or $d$-blocked by $C$ if*

1. *it contains a collider $i_k \notin \mathrm{an}_{\mathcal{G}}(C)$, or*
2. *it contains a non-collider $i_k \in C$.*

*The walk (path) $\pi$ is said to be $C$-$d$-open if it is not $d$-blocked by $C$. For two subsets of nodes $A, B \subseteq \mathcal{V}$, we say that $A$ is $d$-separated from $B$ given $C$ in $\mathcal{G}$ if all paths between any node in $A$ and any node in $B$ are $d$-blocked by $C$, and write*

$$A \overset{d}{\underset{\mathcal{G}}{\perp}} B \,|\, C.$$

The next lemma is a straightforward generalization of Lemma 3.3 in [22] to the cyclic setting. It implies that it suffices to formulate $d$-separation in terms of paths rather than walks.

LEMMA A.5.   *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph, $C \subseteq \mathcal{V}$ and $i, j \in \mathcal{V}$. There exists a $C$-$d$-open walk between $i$ and $j$ in $\mathcal{G}$ if and only if there exists a $C$-$d$-open path between $i$ and $j$ in $\mathcal{G}$.*
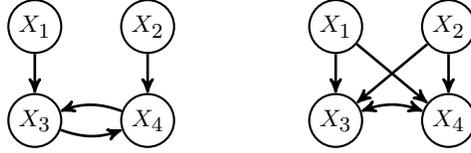
Fig 7: *The graphs of the observationally equivalent SCMs $\mathcal{M}$ (left) and $\tilde{\mathcal{M}}$ (right) of Example A.8 and A.10.*

DEFINITION A.6 (Directed global Markov property).   *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph and $\mathbb{P}_\mathcal{V}$ a probability distribution on $\boldsymbol{\mathcal{X}}_\mathcal{V} = \prod_{i \in \mathcal{V}} \mathcal{X}_i$, where each $\mathcal{X}_i$ is a standard probability space. The probability distribution $\mathbb{P}_\mathcal{V}$ satisfies the* directed global Markov property *relative to $\mathcal{G}$ if for all subsets $A, B, C \subseteq \mathcal{V}$ we have*

$$A \overset{d}{\underset{\mathcal{G}}{\perp}} B \,|\, C \quad \implies \quad \boldsymbol{X}_A \underset{\mathbb{P}_\mathcal{V}}{\perp\!\!\!\perp} \boldsymbol{X}_B \,|\, \boldsymbol{X}_C,$$

*i.e., $(X_i)_{i \in A}$ and $(X_i)_{i \in B}$ are conditionally independent given $(X_i)_{i \in C}$ under $\mathbb{P}_\mathcal{V}$, where we take the canonical projections $X_i : \boldsymbol{\mathcal{X}}_\mathcal{V} \to \mathcal{X}_i$ as random variables.*

From the results in [18] it directly follows that for the observational distribution of an SCM, the directed global Markov property w.r.t. the graph of the SCM (also known as the $d$-separation criterion), holds under one of the following assumptions.

THEOREM A.7 (Directed global Markov property for SCMs [18]).   *Let $\mathcal{M}$ be a uniquely solvable SCM that satisfies at least one of the following three conditions:*

1. *$\mathcal{M}$ is acyclic;*
2. *all endogenous spaces $\mathcal{X}_i$ are discrete and $\mathcal{M}$ is ancestrally uniquely solvable;*
3. *$\mathcal{M}$ is linear (see Definition C.1), each of its causal mechanisms $\{f_i\}_{i \in \mathcal{I}}$ has a non-trivial dependence on at least one exogenous variable, and $\mathbb{P}_\mathcal{E}$ has a density w.r.t. the Lebesgue measure on $\mathbb{R}^\mathcal{J}$.*

*Then, its observational distribution $\mathbb{P}^{\boldsymbol{X}}$ exists, is unique and satisfies the directed global Markov property relative to $\mathcal{G}(\mathcal{M})$ (see Definition A.6).*

The acyclic case is well-known and was first shown in the context of linear-Gaussian structural equation models [75, 32]. The discrete case fixes the erroneous theorem by Pearl and Dechter [52], for which a counterexample was found by Neal [49], by adding the ancestral unique solvability condition, and extends it to allow for bidirected edges in the graph. The linear case is an extension of existing results for the linear-Gaussian setting without bidirected edges [71, 72, 31] to a linear (possibly non-Gaussian) setting with bidirected edges in the graph.

The following counterexample of an SCM for which the directed global Markov property does not hold was already given in [71, 72].

EXAMPLE A.8 (Directed global Markov property does not hold for cyclic SCM).   *Consider the SCM $\mathcal{M} = \langle \boldsymbol{4}, \boldsymbol{4}, \mathbb{R}^4, \mathbb{R}^4, \boldsymbol{f}, \mathbb{P}_{\mathbb{R}^4} \rangle$ with causal mechanism given by*

$$f_1(\boldsymbol{x}, \boldsymbol{e}) = e_1, \quad f_2(\boldsymbol{x}, \boldsymbol{e}) = e_2, \quad f_3(\boldsymbol{x}, \boldsymbol{e}) = x_1 x_4 + e_3, \quad f_4(\boldsymbol{x}, \boldsymbol{e}) = x_2 x_3 + e_4$$

*and $\mathbb{P}_{\mathbb{R}^4}$ is the standard-normal distribution on $\mathbb{R}^4$. The graph of $\mathcal{M}$ is depicted in Figure 7 on the left. The model is uniquely solvable (it is even simple). One can check that for every solution $\boldsymbol{X}$ of $\mathcal{M}$, $X_1$ is not independent of $X_2$ given $\{X_3, X_4\}$. However, the variables $X_1$ and $X_2$ are $d$-separated given $\{X_3, X_4\}$ in $\mathcal{G}(\mathcal{M})$. Hence the global directed Markov property does not hold here.*

In constraint-based approaches to causal discovery, one usually assumes the converse of the directed global Markov property to hold [73, 51].

DEFINITION A.9 (d-Faithfulness). *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph and $\mathbb{P}_\mathcal{V}$ a probability distribution on $\boldsymbol{\mathcal{X}}_\mathcal{V} = \prod_{i \in \mathcal{V}} \mathcal{X}_i$, where each $\mathcal{X}_i$ is a standard probability space. The probability distribution $\mathbb{P}_\mathcal{V}$ is $d$-faithful to $\mathcal{G}$ if for all subsets $A, B, C \subseteq \mathcal{V}$ we have*

$$A \overset{d}{\underset{\mathcal{G}}{\perp}} B \,|\, C \quad \Longleftarrow \quad \boldsymbol{X}_A \underset{\mathbb{P}_\mathcal{V}}{\perp\!\!\!\perp} \boldsymbol{X}_B \,|\, \boldsymbol{X}_C,$$

*where we take the canonical projections $X_i : \boldsymbol{\mathcal{X}}_\mathcal{V} \to \mathcal{X}_i$ as random variables.*

In other words, the $d$-faithfulness assumption states that the graph explains, via $d$-separation, all the conditional independencies that are present in the observational distribution. Meek [41] showed that for multinomial and linear-Gaussian DAG (i.e., acyclic and causally sufficient SCMs) models, $d$-faithfulness holds for all parameter values up to a measure zero set. Up to our knowledge no such result has been shown for any subclass of SCMs that contains cycles, nor in more general acyclic settings.

A.2.2. *The general directed global Markov property*   In [18] the general directed global Markov property is introduced, that is based on $\sigma$-separation, an extension of $d$-separation. This notion of $\sigma$-separation was derived from the notion of $d$-separation in the acyclification of the graph. The acyclification of a graph generalizes the idea of the collapsed graph for directed graphs, developed by Spirtes [71], to HEDGes. In particular, this notion can be applied to directed mixed graphs, and thus to the graphs of SCMs. The main idea of the acyclification is that under the condition that the SCM is uniquely solvable w.r.t. each strongly connected component, we can replace the causal mechanisms of these strongly connected components by their measurable solution functions, which results in an acyclic SCM. This acyclification preserves the solutions, and $d$-separation in the acyclification can directly be translated into $\sigma$-separation in the original graph. This then leads to the general directed global Markov property. We will discuss this now in more detail.

EXAMPLE A.10 (Construction of an observationally equivalent acyclic SCM). *Consider the SCM $\mathcal{M}$ of Example A.8 which is uniquely solvable w.r.t. all its strongly connected components, i.e., the subsets $\{1\}$, $\{2\}$ and $\{3, 4\}$. Replacing the causal mechanisms of these strongly connected components by their measurable solution functions gives the SCM $\tilde{\mathcal{M}}$ that is the same as $\mathcal{M}$ except that its causal mechanism $\tilde{\boldsymbol{f}}$ is given by*

$$\tilde{f}_1(\boldsymbol{x}, \boldsymbol{e}) := e_1, \quad \tilde{f}_2(\boldsymbol{x}, \boldsymbol{e}) := e_2, \quad \tilde{f}_3(\boldsymbol{x}, \boldsymbol{e}) := \tfrac{x_1 e_4 + e_3}{1 - x_1 x_2}, \quad \tilde{f}_4(\boldsymbol{x}, \boldsymbol{e}) := \tfrac{x_2 e_3 + e_4}{1 - x_1 x_2}.$$

*By construction, $\mathcal{M}$ and $\tilde{\mathcal{M}}$ are observationally equivalent. Because $\tilde{\mathcal{M}}$ is acyclic (see Figure 7 on the right) we can apply the directed global Markov property to $\tilde{\mathcal{M}}$. The fact that $X_1$ and $X_2$ are not $d$-separated given $\{X_3, X_4\}$ in $\mathcal{G}(\tilde{\mathcal{M}})$ is in line with $X_1$ being dependent of $X_2$ given $\{X_3, X_4\}$ for every solution $\boldsymbol{X}$ of $\tilde{\mathcal{M}}$ (and hence of $\mathcal{M}$).*

One of the key insights in [18] is that this example can easily be generalized as follows.

DEFINITION A.11 (Acyclification of an SCM). *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be an SCM that is uniquely solvable w.r.t. each strongly connected component of $\mathcal{G}(\mathcal{M})$. For each $i \in \mathcal{I}$, let $g_i$ be the $i^{th}$ component of a measurable solution function $\boldsymbol{g}_{\mathrm{sc}(i)} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathrm{sc}(i)) \setminus \mathrm{sc}(i)} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathrm{sc}(i))} \to \boldsymbol{\mathcal{X}}_{\mathrm{sc}(i)}$ of $\mathcal{M}$ w.r.t. $\mathrm{sc}(i)$, where $\mathrm{pa}$ and $\mathrm{sc}$ denote the parents and strongly*
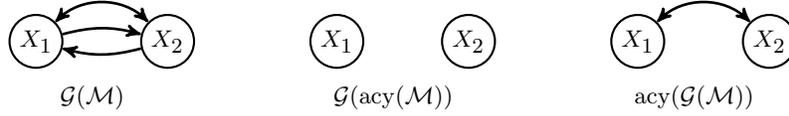
Fig 8: *The graphs of the original SCM $\mathcal{M}$ (left), of the acyclified SCM (center), and of the acyclification of the graph of $\mathcal{M}$ (right) corresponding to Example A.15.*

*connected components according to $\mathcal{G}^a(\mathcal{M})$ respectively. We call the SCM $\mathcal{M}^{\mathrm{acy}} := \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \hat{\boldsymbol{f}}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ with the* acyclified causal mechanism $\hat{\boldsymbol{f}} : \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{E}} \to \boldsymbol{\mathcal{X}}$ *given by*

$$\hat{f}_i(\boldsymbol{x}, \boldsymbol{e}) = g_i(\boldsymbol{x}_{\mathrm{pa(sc}(i))\backslash \mathrm{sc}(i)}, \boldsymbol{e}_{\mathrm{pa(sc}(i))}), \quad i \in \mathcal{I},$$

*an* acyclification *of $\mathcal{M}$. We denote by $\mathrm{acy}(\mathcal{M})$ the equivalence class of the acyclifications of $\mathcal{M}$.*

Note that $\mathrm{acy}(\mathcal{M})$ is well-defined: all acyclifications of an SCM $\mathcal{M}$ belong to the same equivalence class of SCMs.

PROPOSITION A.12. *Let $\mathcal{M}$ be an SCM that is uniquely solvable w.r.t. each strongly connected component of $\mathcal{G}(\mathcal{M})$. Then an acyclification $\mathcal{M}^{\mathrm{acy}}$ of $\mathcal{M}$ is acyclic and observationally equivalent to $\mathcal{M}$.*

We can also define a graphical acyclification for directed mixed graphs, which is a special case of the operation defined in [18] for HEDGes.

DEFINITION A.13 (Acyclification of a directed mixed graph). *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph. The* acyclification of $\mathcal{G}$ *maps $\mathcal{G}$ to the* acyclified graph $\mathcal{G}^{\mathrm{acy}} := (\mathcal{V}, \hat{\mathcal{E}}, \hat{\mathcal{B}})$ *with directed edges $j \to i \in \hat{\mathcal{E}}$ if and only if $j \in \mathrm{pa}_{\mathcal{G}}(\mathrm{sc}_{\mathcal{G}}(i)) \setminus \mathrm{sc}_{\mathcal{G}}(i)$ and bidirected edges $i \leftrightarrow j \in \hat{\mathcal{B}}$ if and only if there exist $i' \in \mathrm{sc}_{\mathcal{G}}(i)$ and $j' \in \mathrm{sc}_{\mathcal{G}}(j)$ with $i' = j'$ or $i' \leftrightarrow j' \in \mathcal{B}$.*

The following compatibility result is immediate from the definitions.

PROPOSITION A.14. *Let $\mathcal{M}$ be an SCM that is uniquely solvable w.r.t. each strongly connected component of $\mathcal{G}(\mathcal{M})$. Then $\mathcal{G}^a(\mathrm{acy}(\mathcal{M})) \subseteq \mathrm{acy}(\mathcal{G}^a(\mathcal{M}))$ and $\mathcal{G}(\mathrm{acy}(\mathcal{M})) \subseteq \mathrm{acy}(\mathcal{G}(\mathcal{M}))$.*

The following example illustrates that the graph of the acyclification of an SCM can be a strict subgraph of the acyclification of the graph of the SCM.

EXAMPLE A.15 (Graph of the acyclification of the SCM is a strict subgraph of the acyclification of its graph). *Consider the SCM $\mathcal{M} = \langle \mathbf{2}, \mathbf{1}, \mathbb{R}^2, \mathbb{R}, \boldsymbol{f}, \mathbb{P}_{\mathbb{R}} \rangle$ with the causal mechanism defined by*

$$f_1(\boldsymbol{x}, e) = x_2 - e, \quad f_2(\boldsymbol{x}, e) = \tfrac{1}{2}x_1 + e$$

*and $\mathbb{P}_{\mathbb{R}}$ the standard-normal measure on $\mathbb{R}$. The SCM $\mathcal{M}$ is uniquely solvable w.r.t. the (only) strongly connected component $\{1, 2\}$. An acyclification of $\mathcal{M}$ is the acyclified SCM $\mathcal{M}^{\mathrm{acy}}$ with the acyclified causal mechanism $\hat{\boldsymbol{f}}$ defined by*

$$\hat{f}_1(\boldsymbol{x}, e) = 0, \quad \hat{f}_2(\boldsymbol{x}, e) = e.$$

*The graph $\mathcal{G}(\mathrm{acy}(\mathcal{M}))$ is a strict subgraph of $\mathrm{acy}(\mathcal{G}(\mathcal{M}))$ as can be seen in Figure 8.*

Translating the notion of $d$-separation from the acyclified graph back to the original graph led to the notion of $\sigma$-separation.

DEFINITION A.16 ($\sigma$-separation [18]).  *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph and let $C \subseteq \mathcal{V}$ be a subset of nodes. A walk (path) $\pi = (i_0, \epsilon_1, i_1, \ldots, i_n)$ in $\mathcal{G}$ is said to be $C$-$\sigma$-blocked or $\sigma$-blocked by $C$ if*

1. *its first node $i_0 \in C$ or its last node $i_n \in C$, or*
2. *it contains a collider $i_k \notin \mathrm{an}_{\mathcal{G}}(C)$, or*
3. *it contains a non-endpoint non-collider $i_k \in C$ that points towards a neighboring node on $\pi$ that lies in a different strongly connected component of $\mathcal{G}$, i.e., such that $i_{k-1} \leftarrow i_k$ in $\pi$ and $i_{k-1} \notin \mathrm{sc}_{\mathcal{G}}(i_k)$, or $i_k \rightarrow i_{k+1}$ in $\pi$ and $i_{k+1} \notin \mathrm{sc}_{\mathcal{G}}(i_k)$.*

*The walk (path) $\pi$ is said to be $C$-$\sigma$-open if it is not $\sigma$-blocked by $C$. For two subsets of nodes $A, B \subseteq \mathcal{V}$, we say that $A$ is $\sigma$-separated from $B$ given $C$ in $\mathcal{G}$ if all paths between any node in $A$ and any node in $B$ are $\sigma$-blocked by $C$, and write*

$$A \overset{\sigma}{\underset{\mathcal{G}}{\perp}} B \,|\, C.$$

The only difference between $\sigma$-separation and $d$-separation is that $d$-separation does not have the extra condition on the non-collider that it has to point to a node in a different strongly connected component. It is therefore obvious that $\sigma$-separation reduces to $d$-separation for acyclic graphs, since $\mathrm{sc}_{\mathcal{G}}(i) = \{i\}$ for each $i \in \mathcal{V}$ in that case.

Although for proofs it is often easier to make use of walks, it suffices to formulate $\sigma$-separation in term of paths rather than walks because of the following result, which is analogous to a similar result for $d$-separation (see Lemma A.5).

LEMMA A.17.  *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph, $C \subseteq \mathcal{V}$ and $i, j \in \mathcal{V}$. There exists a $C$-$\sigma$-open walk between $i$ and $j$ in $\mathcal{G}$ if and only if there exists a $C$-$\sigma$-open path between $i$ and $j$ in $\mathcal{G}$.*

It is clear from the definitions that $\sigma$-separation implies $d$-separation. The other way around does not hold in general, as can be seen in the following example.

EXAMPLE A.18 ($d$-separation does not imply $\sigma$-separation).  *Consider the directed graph $\mathcal{G}$ as depicted in Figure 7 (left). Here $X_1$ is $d$-separated from $X_2$ given $\{X_3, X_4\}$, but $X_1$ is not $\sigma$-separated from $X_2$ given $\{X_3, X_4\}$.*

The following result in [18] relates $\sigma$-separation to $d$-separation.

PROPOSITION A.19.  *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph. Then for $A, B, C \subseteq \mathcal{V}$,*

$$A \overset{\sigma}{\underset{\mathcal{G}}{\perp}} B \,|\, C \iff A \overset{d}{\underset{\mathrm{acy}(\mathcal{G})}{\perp}} B \,|\, C.$$

By replacing in Definition A.6 "$d$-separation" by "$\sigma$-separation", one obtains the formulation of what Forré and Mooij [18] termed the general directed global Markov property.

DEFINITION A.20 (General directed global Markov property [18]). *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph and $\mathbb{P}_\mathcal{V}$ a probability distribution on $\boldsymbol{\mathcal{X}}_\mathcal{V} = \prod_{i \in \mathcal{V}} \mathcal{X}_i$, where each $\mathcal{X}_i$ is a standard probability space. The probability distribution $\mathbb{P}_\mathcal{V}$ satisfies the* general directed global Markov property *relative to $\mathcal{G}$ if for all subsets $A, B, C \subseteq \mathcal{V}$ we have*

$$A \stackrel{\sigma}{\underset{\mathcal{G}}{\perp}} B \mid C \quad \Longrightarrow \quad \boldsymbol{X}_A \underset{\mathbb{P}_\mathcal{V}}{\perp\!\!\!\perp} \boldsymbol{X}_B \mid \boldsymbol{X}_C \,,$$

*i.e., $(X_i)_{i \in A}$ and $(X_i)_{i \in B}$ are conditionally independent given $(X_i)_{i \in C}$ under $\mathbb{P}_\mathcal{V}$, where we take the canonical projections $X_i : \boldsymbol{\mathcal{X}}_\mathcal{V} \to \mathcal{X}_i$ as random variables.*

The fact that $\sigma$-separation implies $d$-separation means that the directed global Markov property implies the general directed global Markov property. In other words, the general directed global Markov property is weaker than the directed global Markov property. It is actually strictly weaker, as we saw in Example A.18.

The following fundamental result, also known as the $\sigma$-separation criterion, follows directly from the theory in [18].

THEOREM A.21 (General directed global Markov property for SCMs). *Let $\mathcal{M}$ be an SCM that is uniquely solvable w.r.t. each strongly connected component of $\mathcal{G}(\mathcal{M})$. Then its observational distribution $\mathbb{P}^{\boldsymbol{X}}$ exists, is unique and it satisfies the general directed global Markov property relative to $\mathcal{G}(\mathcal{M})$.*[15]

The proof is based on the reasoning that, for $A, B, C \subseteq \mathcal{I}$, if $A$ is $\sigma$-separated from $B$ given $C$ in $\mathcal{G}(\mathcal{M})$, then $A$ is $d$-separated from $B$ by $C$ in $\mathrm{acy}(\mathcal{G}(\mathcal{M}))$ and hence in $\mathcal{G}(\mathrm{acy}(\mathcal{M}))$, and since $\mathrm{acy}(\mathcal{M})$ is acyclic and observationally equivalent to $\mathcal{M}$, it follows from the directed global Markov property applied to $\mathrm{acy}(\mathcal{M})$ that $\boldsymbol{X}_A \perp\!\!\!\perp_{\mathbb{P}^{\boldsymbol{X}}} \boldsymbol{X}_B \mid \boldsymbol{X}_C$ for every solution $\boldsymbol{X}$ of $\mathcal{M}$. Note that the ancestral unique solvability condition for the discrete case is strictly weaker than the condition of unique solvability w.r.t. each strongly connected component in Theorem A.21. For the linear case, the condition of unique solvability is equivalent to the condition of unique solvability w.r.t. each strongly connected component (see Proposition C.4).

The results in Theorems A.7 and A.21 are not preserved under perfect intervention, because intervening on a strongly connected component could split it into several strongly connected components with different solvability properties. As the class of simple SCMs is preserved under perfect intervention and the twin operation (Proposition 8.2), we obtain the following corollary.

COROLLARY A.22 (Global Markov properties for simple SCMs). *Let $\mathcal{M}$ be a simple SCM. Then, the*

1. *observational distribution,*
2. *interventional distribution after perfect intervention on $I \subset \mathcal{I}$,*
3. *counterfactual distribution after perfect intervention on $\tilde{I} \subseteq \mathcal{I} \cup \mathcal{I}'$,*

---

[15]Since [18] also provides results under the weaker condition that an SCM is solvable (not necessarily uniquely) w.r.t. each strongly connected component of $\mathcal{G}(\mathcal{M})$, one might believe that Theorem A.21 could be generalized to stating that in that case, any of its observational distributions satisfies the general directed global Markov property. However, that is not true: consider for example the SCM $\mathcal{M} = \langle \boldsymbol{2}, \emptyset, \mathbb{R}^2, \boldsymbol{1}, \boldsymbol{f}, \mathbb{P_1} \rangle$ with $f_1(\boldsymbol{x}) = x_1$ and $f_2(\boldsymbol{x}) = x_2$. Then $\mathcal{M}$ is solvable w.r.t. each of its strongly connected components $\{1\}$ and $\{2\}$. The solution with $X_1 = X_2$ shows a dependence between $X_1$ and $X_2$ and thus $X_1 \perp\!\!\!\perp X_2$ does not hold. In general, all strongly connected components that admit multiple solutions may be dependent on any other variable(s) in the model.

*all exist, are unique and satisfy the general directed global Markov property relative to*
$\mathcal{G}(\mathcal{M})$, $\mathrm{do}(I)(\mathcal{G}(\mathcal{M}))$ *and* $\mathrm{do}(\tilde{I})(\mathrm{twin}(\mathcal{G}(\mathcal{M})))$ *respectively. Moreover, if $\mathcal{M}$ satisfies at least one of the three conditions (1), (2), (3) of Theorem A.7, then they also satisfies the directed global Markov property relative to* $\mathcal{G}(\mathcal{M})$, $\mathrm{do}(I)(\mathcal{G}(\mathcal{M}))$ *and* $\mathrm{do}(\tilde{I})(\mathrm{twin}(\mathcal{G}(\mathcal{M})))$ *respectively.*

Similarly to $d$-faithfulness, $\sigma$-faithfulness[16] is defined as follows.

DEFINITION A.23 ($\sigma$-Faithfulness). *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{B})$ be a directed mixed graph and $\mathbb{P}_{\mathcal{V}}$ a probability distribution on $\boldsymbol{\mathcal{X}}_{\mathcal{V}} = \prod_{i \in \mathcal{V}} \mathcal{X}_i$, where each $\mathcal{X}_i$ is a standard probability space. The probability distribution $\mathbb{P}_{\mathcal{V}}$ is $\sigma$-faithful to $\mathcal{G}$ if for all subsets $A, B, C \subseteq \mathcal{V}$ we have*

$$A \overset{\sigma}{\underset{\mathcal{G}}{\perp}} B \mid C \quad \Longleftarrow \quad \boldsymbol{X}_A \underset{\mathbb{P}_{\mathcal{V}}}{\perp\!\!\!\perp} \boldsymbol{X}_B \mid \boldsymbol{X}_C,$$

*where we take the canonical projections $X_i : \boldsymbol{\mathcal{X}}_{\mathcal{V}} \to \mathcal{X}_i$ as random variables.*

In other words, the graph explains, via $\sigma$-separation, all the conditional independencies that are present in the observational distribution. Although it has been conjectured [72] that under certain conditions $\sigma$-faithfulness should hold, it is not known if it holds generically for the class of simple SCMs.

**A.3. Simple SCMs** In this subsection, we give a summary of useful sufficient conditions for determining the different causal relationships according to a specific simple SCM $\mathcal{M}$, which follow directly from Proposition 7.1 and 7.3.

COROLLARY A.24 (Sufficient conditions for the presence of causal relationships for simple SCMs). *Let $\mathcal{M}$ be a simple SCM and $i, j \in \mathcal{I}$ such that $i \neq j$ and $I := \mathcal{I} \setminus \{i, j\}$. Then:*

1. *If there exist values $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$ and $\xi_i \neq \tilde{\xi}_i \in \mathcal{X}_i$ and a measurable set $\mathcal{B}_j \subseteq \mathcal{X}_j$ such that*

$$\mathbb{P}_{(\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)})_{\mathrm{do}(\{i\}, \xi_i)}}(X_j \in \mathcal{B}_j) \neq \mathbb{P}_{(\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)})_{\mathrm{do}(\{i\}, \tilde{\xi}_i)}}(X_j \in \mathcal{B}_j),$$

*then $i$ is a direct cause of $j$ according to $\mathcal{M}$, i.e., $i \to j \in \mathcal{G}(\mathcal{M})$;*

2. *If there exist values $\xi_i \neq \tilde{\xi}_i \in \mathcal{X}_i$ and a measurable set $\mathcal{B}_j \subseteq \mathcal{X}_j$ such that*

$$\mathbb{P}_{\mathcal{M}_{\mathrm{do}(\{i\}, \xi_i)}}(X_j \in \mathcal{B}_j) \neq \mathbb{P}_{\mathcal{M}_{\mathrm{do}(\{i\}, \tilde{\xi}_i)}}(X_j \in \mathcal{B}_j),$$

*then $i$ is a cause of $j$ according to $\mathcal{M}$, i.e., $i \to \cdots \to j$ in $\mathcal{G}(\mathcal{M})$;*

3. *If $j \notin \mathrm{an}_{\mathcal{G}(\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)})}(i)$ and there exist a value $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$ and a measurable set $\mathcal{B}_j \subseteq \mathcal{X}_j$ such that for every version of the regular conditional probability $\mathbb{P}_{\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}}(X_j \in \mathcal{B}_j \mid X_i = \xi_i)$ there exists a value $\xi_i \in \mathcal{X}_i$ such that*

$$\mathbb{P}_{\left(\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}\right)_{\mathrm{do}(\{i\}, \xi_i)}}(X_j \in \mathcal{B}_j) \neq \mathbb{P}_{\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}}(X_j \in \mathcal{B}_j \mid X_i = \xi_i),$$

*then $i$ and $j$ are confounded according to $\mathcal{M}$, i.e., $i \leftrightarrow j \in \mathcal{G}(\mathcal{M})$.*

**A.4. Modular SCMs** In this subsection we relate the class of (simple) SCMs to that of modular SCMs. Modular SCMs introduced by Forré and Mooij [18] are causal graphical models on which marginalizations and interventions are defined and they satisfy the general directed global Markov property. For a comprehensive account on modular SCMs we refer the reader to [18].

---

[16]In [63] it is called "collapsed graph faithfulness".

A.4.1. *Definition of a modular SCM*   In contrast to an SCM from which a graph can be derived, a modular SCM is defined in terms of a graphical object, which Forré and Mooij [18] call a directed graph with hyperedges (HEDG). The hyperedges of a HEDG are described in terms of a simplicial complex.

DEFINITION A.25 (Simplicial complex).   *Let $\mathcal{V}$ be a finite set. A* simplicial complex $\mathcal{H}$ *over $\mathcal{V}$ is a set of subsets of $\mathcal{V}$ such that*

1. *all single element sets $\{v\}$ are in $\mathcal{H}$ for $v \in \mathcal{V}$, and*
2. *if $\mathcal{F} \in \mathcal{H}$, then also all subsets $\tilde{\mathcal{F}} \subseteq F$ are elements of $\mathcal{H}$.*

DEFINITION A.26 (Directed graph with hyperedges (HEDGes) [18]).   *A directed graph with hyperedges (HEDG) is a triple $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{H})$, where $(\mathcal{V}, \mathcal{E})$ is a directed graph and $\mathcal{H}$ a simplicial complex over the set of nodes $\mathcal{V}$. The elements $\mathcal{F}$ of $\mathcal{H}$ are called* hyperedges *of $\mathcal{G}$. The elements $\mathcal{F}$ of $\mathcal{H}$ that are inclusion-maximal elements of $\mathcal{H}$ are called* maximal hyperedges *and are denoted by $\hat{\mathcal{H}}$.*

A HEDG $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{H})$ can be represented as a directed graph $\bar{\mathcal{G}} := (\mathcal{V}, \mathcal{E})$ consisting of nodes $\mathcal{V}$ and directed edges $\mathcal{E}$, with additional maximal hyperedges $\mathcal{F} \in \hat{\mathcal{H}}$ with $|\mathcal{F}| \geq 2$ (i.e., not corresponding to single element sets $\{v\} \in \hat{\mathcal{H}}$), that point to their target nodes $v \in \mathcal{F}$. For a HEDG $\mathcal{G}$ we define $\mathrm{pa}_{\mathcal{G}}$, $\mathrm{ch}_{\mathcal{G}}$, etc. in terms of the underlying directed graph $\bar{\mathcal{G}}$, i.e., $\mathrm{pa}_{\bar{\mathcal{G}}}$, $\mathrm{ch}_{\bar{\mathcal{G}}}$, etc. respectively.

A *loop* in a HEDG $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{H})$ is a subset $\mathcal{O} \subseteq \mathcal{V}$ that is a loop in the underlying directed graph $\bar{\mathcal{G}} = (\mathcal{V}, \mathcal{E})$. In other words, a loop of $\mathcal{G}$ is a set of nodes $\mathcal{O} \subseteq \mathcal{V}$ such that for every two nodes $v, w \in \mathcal{O}$ there are directed paths $v \to \cdots \to w$ and $w \to \cdots \to v$ in $\mathcal{G}$ for which all the intermediate nodes lie in $\mathcal{O}$ (if any exist). In particular, a loop may consist of a single element $\{v\}$ for $v \in \mathcal{V}$. The set of loops in $\mathcal{G}$ is denoted by $\mathcal{L}(\mathcal{G})$.

In order to define a modular SCM one needs the notion of a compatible system of solution functions, which assigns to each loop a separate solution function such that all these solution functions are "compatible" with each other.

DEFINITION A.27 (Compatible system of solution functions[17]).   *Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{H})$ be a HEDG. For every $v \in \mathcal{V}$ and maximal hyperedge $\mathcal{F}$ in $\hat{\mathcal{H}}$, let $\mathcal{X}_v$ and $\mathcal{E}_{\mathcal{F}}$ be standard measurable spaces. For a subset $\mathcal{O} \subseteq \mathcal{V}$ we define[18]*

$$\boldsymbol{\mathcal{X}}_{\mathcal{O}} := \prod_{v \in \mathcal{O}} \mathcal{X}_v \quad and \quad \widehat{\boldsymbol{\mathcal{E}}}_{\mathcal{O}} := \prod_{\substack{\mathcal{F} \in \hat{\mathcal{H}} \\ \mathcal{F} \cap \mathcal{O} \neq \emptyset}} \mathcal{E}_{\mathcal{F}}.$$

*Consider a family of measurable mappings $(\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}$ indexed by $\mathcal{L}(\mathcal{G})$ which are of the form*

$$\boldsymbol{g}_{\mathcal{O}} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}_{\mathcal{G}}(\mathcal{O}) \setminus \mathcal{O}} \times \widehat{\boldsymbol{\mathcal{E}}}_{\mathcal{O}} \to \boldsymbol{\mathcal{X}}_{\mathcal{O}}.$$

*We call the family of measurable mappings $(\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}$ a* compatible system of solution functions, *if for all $\mathcal{O}, \tilde{\mathcal{O}} \in \mathcal{L}(\mathcal{G})$ with $\tilde{\mathcal{O}} \subseteq \mathcal{O}$ and for all $\widehat{e}_{\mathcal{O}} \in \widehat{\boldsymbol{\mathcal{E}}}_{\mathcal{O}}$ and $\boldsymbol{x}_{\mathrm{pa}_{\mathcal{G}}(\mathcal{O}) \cup \mathcal{O}} \in \boldsymbol{\mathcal{X}}_{\mathrm{pa}_{\mathcal{G}}(\mathcal{O}) \cup \mathcal{O}}$ we have*

$$\boldsymbol{x}_{\mathcal{O}} = \boldsymbol{g}_{\mathcal{O}}(\boldsymbol{x}_{\mathrm{pa}_{\mathcal{G}}(\mathcal{O}) \setminus \mathcal{O}}, \widehat{e}_{\mathcal{O}}) \quad \implies \quad \boldsymbol{x}_{\tilde{\mathcal{O}}} = \boldsymbol{g}_{\tilde{\mathcal{O}}}(\boldsymbol{x}_{\mathrm{pa}_{\mathcal{G}}(\tilde{\mathcal{O}}) \setminus \tilde{\mathcal{O}}}, \widehat{\boldsymbol{e}}_{\tilde{\mathcal{O}}}).$$

---

[17]We deviate from the terminology in [18] where this is called a "compatible system of structural equations".

[18]We use the "hat" notation $\widehat{\boldsymbol{\mathcal{E}}}_{\mathcal{O}}$ to distinguish it from the ordinary subscript convention that $\boldsymbol{\mathcal{E}}_{\mathcal{O}} = \prod_{\mathcal{F} \in \mathcal{O}} \mathcal{E}_{\mathcal{F}}$ for some subset $\mathcal{O} \subseteq \hat{\mathcal{H}}$.

This structure of a compatible system of solution functions is at the heart of the defnition of a modular SCM.

DEFINITION A.28 (Modular structural causal model (mSCM) [18]). *A modular structural causal model (mSCM) is a tuple*

$$\widehat{\mathcal{M}} := \langle \mathcal{G}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, (\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle,$$

*where*

1. $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{H})$ *is a HEDG,*
2. $\boldsymbol{\mathcal{X}} = \prod_{v \in \mathcal{V}} \mathcal{X}_v$ *is the product of standard measurable spaces* $\mathcal{X}_v$,
3. $\boldsymbol{\mathcal{E}} = \prod_{\mathcal{F} \in \hat{\mathcal{H}}} \mathcal{E}_{\mathcal{F}}$ *is the product of standard measurable spaces* $\mathcal{E}_{\mathcal{F}}$,
4. $(\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}$ *is a compatible system of solution functions,*
5. $\mathbb{P}_{\boldsymbol{\mathcal{E}}} = \prod_{\mathcal{F} \in \hat{\mathcal{H}}} \mathbb{P}_{\mathcal{E}_{\mathcal{F}}}$ *is a product measure, where* $\mathbb{P}_{\mathcal{E}_{\mathcal{F}}}$ *is a probability measure on* $\mathcal{E}_{\mathcal{F}}$ *for each* $\mathcal{F} \in \hat{\mathcal{H}}$.

Let $\widehat{\mathcal{M}} = \langle \mathcal{G}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, (\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be a modular SCM and $\mathcal{O}_1, \dots, \mathcal{O}_r \in \mathcal{L}(\mathcal{G})$ the strongly connected components of $\mathcal{G}$ ordered according to a topological order of the DAG of strongly connected components of $\mathcal{G}$. Then, for any random variable $\boldsymbol{E} : \Omega \to \boldsymbol{\mathcal{E}}$ such that $\mathbb{P}^{\boldsymbol{E}} = \mathbb{P}_{\boldsymbol{\mathcal{E}}}$ one can inductively define the random variables $X_v := (\boldsymbol{g}_{\mathcal{O}_i})_v(\boldsymbol{X}_{\mathrm{pa}_{\mathcal{G}}(\mathcal{O}_i) \setminus \mathcal{O}_i}, \widehat{\boldsymbol{E}}_{\mathcal{O}_i})$ for all $v \in \mathcal{O}_i$ for all $i \geq 1$, starting at $X_v := (\boldsymbol{g}_{\mathcal{O}_1})_v(\widehat{\boldsymbol{E}}_{\mathcal{O}_1})$ for all $v \in \mathcal{O}_1$. Because $(\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}$ is a compatible system of solution functions, we have for every $\mathcal{O} \in \mathcal{L}(\mathcal{G})$

$$\boldsymbol{X}_{\mathcal{O}} = \boldsymbol{g}_{\mathcal{O}}(\boldsymbol{X}_{\mathrm{pa}_{\mathcal{G}}(\mathcal{O}) \setminus \mathcal{O}}, \widehat{\boldsymbol{E}}_{\mathcal{O}}).$$

We call the random variable $\boldsymbol{X}$ a *solution* of the modular SCM $\widehat{\mathcal{M}}$. Note that the solution $\boldsymbol{X}$ depends on the choice of the random variable $\boldsymbol{E} : \Omega \to \boldsymbol{\mathcal{E}}$.

The causal semantics of modular SCMs can be defined in terms of perfect interventions, which is defined as follows.

DEFINITION A.29 (Perfect intervention on an mSCM). *Consider a modular SCM* $\widehat{\mathcal{M}} = \langle \mathcal{G}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, (\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$, *a subset* $I \subseteq \mathcal{V}$ *of endogenous variables and a value* $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$. *The* perfect intervention $\mathrm{do}(I, \boldsymbol{\xi}_I)$ *maps* $\widehat{\mathcal{M}}$ *to the modular SCM*

$$\widehat{\mathcal{M}}_{\mathrm{do}(I, \boldsymbol{\xi}_I)} := \langle \mathcal{G}^{\mathrm{do}}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}^{\mathrm{do}}, (\boldsymbol{g}_{\mathcal{O}}^{\mathrm{do}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G}^{\mathrm{do}})}, \mathbb{P}_{\boldsymbol{\mathcal{E}}^{\mathrm{do}}} \rangle,$$

*where*

1. $\mathcal{G}^{\mathrm{do}} = (\mathcal{V}, \mathcal{E}^{\mathrm{do}}, \mathcal{H}^{\mathrm{do}})$, *where*

$$\mathcal{E}^{\mathrm{do}} = \mathcal{E} \setminus \{v \to w : v \in \mathcal{V}, w \in I\}$$

$$\mathcal{H}^{\mathrm{do}} = \{\mathcal{F} \setminus I : \mathcal{F} \in \mathcal{H}\} \cup \{\{v\} : v \in I\},$$

2. $\phi : \{\mathcal{F} \in \hat{\mathcal{H}} : \mathcal{F} \setminus I \neq \emptyset\} \to \hat{\mathcal{H}}^{\mathrm{do}} \setminus \{\{v\} : v \in I\}$ *is a mapping such that* $\phi(\mathcal{F}) \supseteq \mathcal{F} \setminus I$ *for all* $\mathcal{F} \in \hat{\mathcal{H}}$ *for which* $\mathcal{F} \setminus I \neq \emptyset$,
3. $\boldsymbol{\mathcal{E}}^{\mathrm{do}} = \prod_{\tilde{\mathcal{F}} \in \hat{\mathcal{H}}^{\mathrm{do}}} \mathcal{E}_{\tilde{\mathcal{F}}}^{\mathrm{do}}$, *where*

$$\mathcal{E}_{\tilde{\mathcal{F}}}^{\mathrm{do}} = \begin{cases} \mathcal{X}_v & \text{if } \tilde{\mathcal{F}} = \{v\} \text{ for } v \in I \\ \prod_{\mathcal{F} = \phi^{-1}(\tilde{\mathcal{F}})} \mathcal{E}_{\mathcal{F}} & \text{if } \tilde{\mathcal{F}} \in \hat{\mathcal{H}}^{\mathrm{do}} \setminus \{\{v\} : v \in I\}, \end{cases}$$

4. *for every* $\mathcal{O} \in \mathcal{L}(\mathcal{G}^{\mathrm{do}})$

$$\boldsymbol{g}_{\mathcal{O}}^{\mathrm{do}} = \begin{cases} \mathbb{I}_{\{v\}} & \textit{if } \mathcal{O} = \{v\} \textit{ for } v \in I \\ \boldsymbol{g}_{\mathcal{O}} & \textit{otherwise,} \end{cases}$$

*(note that if $\mathcal{O}$ is a loop in $\mathcal{G}^{\mathrm{do}}$, then it is a loop in $\mathcal{G}$),*

5. $\mathbb{P}_{\boldsymbol{\mathcal{E}}^{\mathrm{do}}} = \prod_{\tilde{\mathcal{F}} \in \hat{\mathcal{H}}^{\mathrm{do}}} \mathbb{P}_{\mathcal{E}_{\tilde{\mathcal{F}}}^{\mathrm{do}}}$, *where*

$$\mathbb{P}_{\mathcal{E}_{\tilde{\mathcal{F}}}^{\mathrm{do}}} = \begin{cases} \delta_{\xi_v} & \textit{if } \tilde{\mathcal{F}} = \{v\} \textit{ for } v \in I \\ \prod_{\mathcal{F} = \phi^{-1}(\tilde{\mathcal{F}})} \mathbb{P}_{\mathcal{E}_{\mathcal{F}}} & \textit{if } \tilde{\mathcal{F}} \in \hat{\mathcal{H}}^{\mathrm{do}} \setminus \{\{v\} : v \in I\}. \end{cases}$$

In contrast to SCMs, these perfect interventions on modular SCMs are directly defined on the underlying HEDG and depend on the choice of the mapping $\phi$.

A.4.2. *Relation between SCMs and modular SCMs* The solutions of a modular SCM can be described by an SCM that is loop-wisely solvable.

DEFINITION A.30 (Induced SCM). *Let $\widehat{\mathcal{M}} = \langle \mathcal{G}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, (\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be a modular SCM. Then, the mapping $\iota$ maps $\widehat{\mathcal{M}}$ to the induced SCM $\tilde{\mathcal{M}} := \langle \tilde{\mathcal{I}}, \tilde{\mathcal{J}}, \tilde{\boldsymbol{\mathcal{X}}}, \tilde{\boldsymbol{\mathcal{E}}}, \tilde{\boldsymbol{f}}, \mathbb{P}_{\tilde{\boldsymbol{\mathcal{E}}}} \rangle$, where*

1. $\tilde{\mathcal{I}} = \mathcal{V}$,
2. $\tilde{\mathcal{J}} = \hat{\mathcal{H}}$,
3. $\tilde{\boldsymbol{\mathcal{X}}} = \boldsymbol{\mathcal{X}}$,
4. $\tilde{\boldsymbol{\mathcal{E}}} = \boldsymbol{\mathcal{E}}$,
5. $\tilde{\boldsymbol{f}}$ *is given by* $\tilde{f}_v = (\boldsymbol{g}_{\{v\}})_v$ *for all* $v \in \mathcal{V}$,
6. $\mathbb{P}_{\tilde{\boldsymbol{\mathcal{E}}}} = \mathbb{P}_{\boldsymbol{\mathcal{E}}}$.

Every solution $\boldsymbol{X}$ of a modular SCM $\widehat{\mathcal{M}}$ is also a solution of the induced SCM $\iota(\widehat{\mathcal{M}})$.

Observe that for the modular SCM $\widehat{\mathcal{M}}$ we have that the induced subgraph $\mathcal{G}^a(\iota(\widehat{\mathcal{M}}))_{\tilde{\mathcal{I}}}$, of the augmented graph of the induced SCM $\mathcal{G}^a(\iota(\widehat{\mathcal{M}}))$ on $\tilde{\mathcal{I}}$, is a subgraph of the underlying HEDG $\mathcal{G}$, i.e., $\mathcal{G}^a(\iota(\widehat{\mathcal{M}}))_{\tilde{\mathcal{I}}} \subseteq \mathcal{G}$. This implies that, in general, the underlying HEDG $\mathcal{G}$ of $\widehat{\mathcal{M}}$ may have more loops than the loops in $\mathcal{G}(\iota(\widehat{\mathcal{M}}))$. For a subset $\mathcal{O} \subseteq \tilde{\mathcal{I}}$ we have for the exogenous parents of the induced SCM $\iota(\widehat{\mathcal{M}})$

$$\mathrm{pa}(\mathcal{O}) \cap \tilde{\mathcal{J}} \subseteq \{\mathcal{F} \in \tilde{\mathcal{J}} : \mathcal{F} \cap \mathcal{O} \neq \emptyset\},$$

where $\mathrm{pa}(\mathcal{O})$ denotes the set of parents of $\mathcal{O}$ in $\mathcal{G}^a(\iota(\widehat{\mathcal{M}}))$. Hence, in general, not all the hyperedges $\mathcal{F} \in \mathcal{H}$ such that $|\mathcal{F}| = 2$ (i.e., bidirected edges) are in the set of bidirected edges $\mathcal{B}$ of the graph of the induced SCM $\mathcal{G}(\iota(\widehat{\mathcal{M}})) = (\mathcal{V}, \mathcal{E}, \mathcal{B})$. We conclude that the graph of the induced SCM is, in general, a sparser graph than the HEDG of the modular SCM.

Next, we show that the compatible system of solution functions of a modular SCM induces a compatible system of solution functions on the induced SCM. For this we need the notion of loop-wise solvability for SCMs.

DEFINITION A.31 (Loop-wise (unique) solvability for SCMs). *We call an SCM $\mathcal{M}$*

1. loop-wisely solvable, *if $\mathcal{M}$ is solvable w.r.t. every loop $\mathcal{O} \in \mathcal{L}(\mathcal{G}(\mathcal{M}))$, and*
2. loop-wisely uniquely solvable, *if $\mathcal{M}$ is uniquely solvable w.r.t. every loop $\mathcal{O} \in \mathcal{L}(\mathcal{G}(\mathcal{M}))$.*
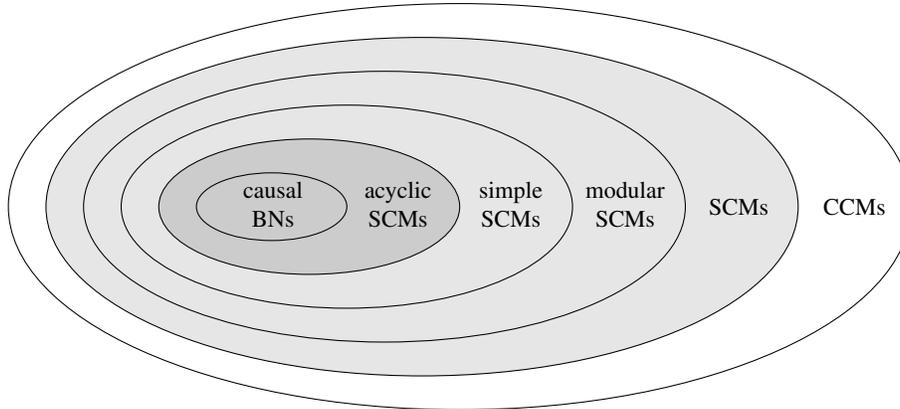
Fig 9: *Overview of causal graphical models. The "gray" and "dark gray" areas contain all the causal graphical models that can be modeled by an SCM and an acyclic SCM respectively.*

DEFINITION A.32 (Compatible system of solution functions for SCMs). *For a loop-wisely solvable SCM $\mathcal{M}$, we call a family of measurable solution functions $(\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G}(\mathcal{M}))}$, where $\boldsymbol{g}_{\mathcal{O}}$ is a measurable solution function of $\mathcal{M}$ w.r.t. $\mathcal{O}$, a* compatible system of solution functions, *if for all $\mathcal{O}, \tilde{\mathcal{O}} \in \mathcal{L}(\mathcal{G}(\mathcal{M}))$ with $\tilde{\mathcal{O}} \subseteq \mathcal{O}$ and for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$ we have*

$$\boldsymbol{x}_{\mathcal{O}} = \boldsymbol{g}_{\mathcal{O}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}) \quad \implies \quad \boldsymbol{x}_{\tilde{\mathcal{O}}} = \boldsymbol{g}_{\tilde{\mathcal{O}}}(\boldsymbol{x}_{\mathrm{pa}(\tilde{\mathcal{O}}) \setminus \tilde{\mathcal{O}}}, \boldsymbol{e}_{\mathrm{pa}(\tilde{\mathcal{O}})}).$$

The induced SCM of a modular SCM always has a compatible system of solution functions, by construction.

PROPOSITION A.33. *Let $\widehat{\mathcal{M}} = \langle \mathcal{G}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, (\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be a modular SCM. Then, the induced SCM $\tilde{\mathcal{M}} := \iota(\widehat{\mathcal{M}})$ is loop-wisely solvable. Moreover, it has a compatible system of solution functions $(\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G}(\tilde{\mathcal{M}}))}$, where $\boldsymbol{g}_{\mathcal{O}}$ is a measurable solution function of $\tilde{\mathcal{M}}$ w.r.t. $\mathcal{O}$.*

This shows that a modular SCM can be seen as an SCM together with an additional structure of a compatible system of solution functions, and is, in particular, loop-wisely solvable.

Moreover, the class of simple SCMs corresponds exactly with those SCMs that are loop-wisely uniquely solvable.

LEMMA A.34. *An SCM $\mathcal{M}$ is simple if and only if it is loop-wisely uniquely solvable.*

In particular, for simple SCMs, or loop-wisely uniquely solvable SCMs, there always exists a compatible system of solution functions.

PROPOSITION A.35. *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be a simple SCM. Then, every family of measurable solution functions $(\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G}(\mathcal{M}))}$, where $\boldsymbol{g}_{\mathcal{O}}$ is a measurable solution function of $\mathcal{M}$ w.r.t. $\mathcal{O}$, is a compatible system of solution functions.*

**A.5. Overview of causal graphical models**  Figure 9 gives an overview of the causal graphical models related to SCMs. The "gray" area contains all the causal graphical models that can be modeled by an SCM, by which we mean, that there exists an SCM that can describe all its observational and interventional distributions. The "dark gray" area contains all

the causal graphical models which can be modeled by an acyclic SCM. Acyclic SCMs generalize causal Bayesian networks (causal BNs) [51] to allow for latent confounders and to derive counterfactuals. Simple SCMs form a subclass of SCMs that extends acyclic SCMs to the cyclic setting, while preserving many of their convenient properties. Modular SCMs [18] can be seen as SCMs that have an additional structure of compatible system of solution functions and contain, in particular, the class of simple SCMs. Forré and Mooij [18] showed that modular SCMs satisfy various convenient properties, like marginalization and the general directed global Markov property. We show that for SCMs in general various of those properties still hold under certain solvability conditions. A generalization of SCMs, known as *causal constraints models (CCMs)*, has been proposed [3] in order to completely model the causal semantics of the equilibrium solutions of a dynamical system given the initial conditions. This class of CCMs is rich enough to model the causal semantics of SCMs, but does not come with a single graphical representation that provides both a Markov property and a causal interpretation [4].

## APPENDIX B: (UNIQUE) SOLVABILITY PROPERTIES

In this appendix we provide additional (unique) solvability properties for SCMs. In Appendix B.1 we provide a sufficient condition of solvability w.r.t. (strict) subsets. In Appendix B.2 we discuss how (unique) solvability is preserved under strict super- and subsets. In Appendix B.3 we discuss how (unique) solvability is preserved under unions and intersections. The proofs of the theoretical results in this appendix are given in Appendix E.

**B.1. Sufficient condition for solvability w.r.t. subsets** For solvability w.r.t. a (strict) subset of $\mathcal{I}$ there exists a sufficient condition that is similar to the sufficient (and necessary) condition (2) in Theorem 3.3 in the sense that it is formulated in terms of the solutions of (a subset of) the structural equations, but no measurability is required.

PROPOSITION B.1 (Sufficient condition for solvability w.r.t. a subset). *Let* $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ *be an SCM and* $\mathcal{O} \subseteq \mathcal{I}$ *a subset. If for* $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-*almost every* $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ *and for all* $\boldsymbol{x}_{\setminus \mathcal{O}} \in \boldsymbol{\mathcal{X}}_{\setminus \mathcal{O}}$ *the topological space*

$$\boldsymbol{\mathcal{S}}_{(\boldsymbol{e}, \boldsymbol{x}_{\setminus \mathcal{O}})} := \{\boldsymbol{x}_{\mathcal{O}} \in \boldsymbol{\mathcal{X}}_{\mathcal{O}} : \boldsymbol{x}_{\mathcal{O}} = \boldsymbol{f}_{\mathcal{O}}(\boldsymbol{x}, \boldsymbol{e})\},$$

*with the subspace topology induced by* $\boldsymbol{\mathcal{X}}_{\mathcal{O}}$ *is non-empty and* $\sigma$-*compact,*[19] *then* $\mathcal{M}$ *is solvable w.r.t.* $\mathcal{O}$.

For many purposes, this condition of $\sigma$-compactness suffices since it contains for example all countable discrete spaces, every interval of the real line, and moreover all the Euclidean spaces. In particular, it suffices to prove a sufficient and necessary condition for unique solvability w.r.t. a subset, in terms of the solutions of a subset of the structural equations (see Theorem 3.8). For larger solution spaces, we refer the reader to [30]. For the class of linear SCMs (see Definition C.1), we provide in Proposition C.2 a sufficient and necessary condition for solvability w.r.t. a (strict) subset of $\mathcal{I}$.

**B.2. (Unique) solvability w.r.t. strict super- and subsets** In general, (unique) solvability w.r.t. $\mathcal{O} \subseteq \mathcal{I}$ does not imply (unique) solvability w.r.t. a strict superset $\mathcal{O} \subsetneq \mathcal{V} \subseteq \mathcal{I}$ nor w.r.t. a strict subset $\mathcal{W} \subsetneq \mathcal{O}$, as can be seen in the following example.

---

[19] A topological space $\boldsymbol{\mathcal{X}}$ is called $\sigma$-*compact* if it is the union of a countable set of compact topological spaces.

EXAMPLE B.2 (Solvability is not preserved under strict sub- or supersets). *Consider the SCM $\mathcal{M} = \langle \mathbf{3}, \emptyset, \mathbb{R}^3, \mathbf{1}, \boldsymbol{f}, \mathbb{P_1} \rangle$ where the causal mechanism is given by*

$$f_1(\boldsymbol{x}) = x_1 \cdot (1 - \mathbf{1}_{\{1\}}(x_2)) + 1 \,, \ f_2(\boldsymbol{x}) = x_2 \,, \ f_3(\boldsymbol{x}) = x_3 \cdot (1 - \mathbf{1}_{\{-1\}}(x_2)) + 1 \,.$$

*This SCM is (uniquely) solvable w.r.t. the subsets $\{1, 2\}$, $\{2, 3\}$, however it is not (uniquely) solvable w.r.t. the subsets $\{1\}$, $\{3\}$ and $\{1, 2, 3\}$, and not uniquely solvable w.r.t. $\{2\}$.*

However, in Proposition 3.13 we show that solvability w.r.t. $\mathcal{O}$ implies solvability w.r.t. every ancestral subset in $\mathcal{G}(\mathcal{M})_\mathcal{O}$.

**B.3. (Unique) solvability w.r.t. unions and intersections** In general, (unique) solvability is not preserved under unions and intersections. The following example illustrates that (unique) solvability is in general not preserved under intersections.

EXAMPLE B.3 (Solvability is not preserved under intersections). *Consider the SCM $\mathcal{M} = \langle \mathbf{3}, \emptyset, \mathbb{R}^3, \mathbf{1}, \boldsymbol{f}, \mathbb{P_1} \rangle$ where the causal mechanism is given by*

$$f_1(\boldsymbol{x}) = 0 \,, \ f_2(\boldsymbol{x}) = x_2 \cdot (1 - \mathbf{1}_{\{0\}}(x_1 \cdot x_3)) + 1 \,, \ f_3(\boldsymbol{x}) = 0 \,.$$

*Then $\mathcal{M}$ is (uniquely) solvable w.r.t. $\{1, 2\}$ and $\{2, 3\}$, however it is not (uniquely) solvable w.r.t. their intersection.*

Example B.2 gives an example where (unique) solvability is not preserved under unions. Even, if we take the union of disjoint subsets, (unique) solvability is not preserved (see Example 3.7 with $\alpha = \beta = 1$). Although, in general, unique solvability is not preserved under unions, we show next that unique solvability is preserved under the union of ancestral subsets, under the following assumptions.

PROPOSITION B.4 (Combining measurable solution functions on different sets). *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P_{\mathcal{E}}} \rangle$ be an SCM, $\mathcal{O} \subseteq \mathcal{I}$ a subset and $\mathcal{A}, \tilde{\mathcal{A}} \subseteq \mathcal{O}$ two ancestral subsets in $\mathcal{G}(\mathcal{M})_\mathcal{O}$. If $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{A}$, $\tilde{\mathcal{A}}$ and $\mathcal{A} \cap \tilde{\mathcal{A}}$, then $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{A} \cup \tilde{\mathcal{A}}$.*

A consequence of this property is that in order to check whether an SCM is ancestrally uniquely solvable w.r.t. $\mathcal{O}$, it suffices to check that it is uniquely solvable w.r.t. the ancestral subsets for each node in $\mathcal{O}$.

COROLLARY B.5. *Let $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \boldsymbol{f}, \mathbb{P_{\mathcal{E}}} \rangle$ be an SCM and $\mathcal{O} \subseteq \mathcal{I}$ a subset. Then $\mathcal{M}$ is ancestrally uniquely solvable w.r.t. $\mathcal{O}$ if and only if $\mathcal{M}$ is uniquely solvable w.r.t. $\mathrm{an}_{\mathcal{G}(\mathcal{M})_\mathcal{O}}(i)$ for every $i \in \mathcal{O}$.*

## APPENDIX C: LINEAR SCMS

In this appendix we provide some results about (unique) solvability and marginalization for linear SCMs. Linear SCMs form a special class of SCMs that has seen much attention in the literature [see, e.g., 5, 27]. The proofs of the theoretical results in this appendix are given in Appendix E.

DEFINITION C.1 (Linear SCM). *We call an SCM $\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \mathbb{R}^\mathcal{I}, \mathbb{R}^\mathcal{J}, \boldsymbol{f}, \mathbb{P_{\mathbb{R}^\mathcal{J}}} \rangle$ linear if each component of the causal mechanism is a linear combination of the endogenous and exogenous variables, that is*

$$f_i(\boldsymbol{x}, \boldsymbol{e}) = \sum_{j \in \mathcal{I}} B_{ij} x_j + \sum_{k \in \mathcal{J}} \Gamma_{ik} e_k \,,$$

where $i \in \mathcal{I}$, $B \in \mathbb{R}^{\mathcal{I} \times \mathcal{I}}$ and $\Gamma \in \mathbb{R}^{\mathcal{I} \times \mathcal{J}}$ are matrices, and $\mathbb{P}_{\mathbb{R}^{\mathcal{J}}}$ is a product probability measure[20] on $\mathbb{R}^{\mathcal{J}}$.

For a subset $\mathcal{O} \subseteq \mathcal{I}$ we also use the shorthand vector-notation

$$\boldsymbol{f}_{\mathcal{O}}(\boldsymbol{x}, \boldsymbol{e}) = B_{\mathcal{O}\mathcal{I}} \boldsymbol{x} + \Gamma_{\mathcal{O}\mathcal{J}} \boldsymbol{e}.$$

A non-zero coefficient $B_{ij}$ for $i, j \in \mathcal{I}$ such that $i \neq j$ corresponds with a directed edge $j \to i$ in the (augmented) graph, and a coefficient $B_{ii} = 1$ for $i \in \mathcal{I}$ corresponds with a self-cycle $i \to i$ in the (augmented) graph of the SCM. A non-zero coefficient $\Gamma_{ij}$ for $i \in \mathcal{I}$, $j \in \mathcal{J}$ with $\mathbb{P}_{\mathcal{E}_j}$ a non-degenerate probability distribution over $\mathbb{R}$ corresponds with a directed edge $j \to i$ in the augmented graph, and a non-zero entry $(\Gamma\Gamma^T)_{ij}$ for $i \in \mathcal{I}$, $j \in \mathcal{J}$ such that there exists a $k \in \mathcal{J}$ with $\Gamma_{ik}, \Gamma_{kj} \neq 0$ and $\mathbb{P}_{\mathcal{E}_k}$ a non-degenerate probability distribution over $\mathbb{R}$ corresponds with a bidirected edge $i \leftrightarrow j$ in the graph of the SCM.

For linear SCMs the solvability condition w.r.t. a subset, Definition 3.1, translates into a matrix condition. In order to state this condition we need to define the pseudoinverse (or the Moore-Penrose inverse) $A^+$ of a real matrix $A$ [54, 24]. The *pseudoinverse of the matrix $A$ is defined by* $A^+ := V\Sigma^+ U^*$, where $A = U\Sigma V^*$ is the singular value decomposition of $A$ and $\Sigma^+$ is obtained by replacing each non-zero entry on the diagonal of $\Sigma$ by its reciprocal [24]. One of its useful properties is that $AA^+A = A$.

PROPOSITION C.2 (Sufficient and necessary condition for solvability w.r.t. a subset for linear SCMs). *Let $\mathcal{M}$ be a linear SCM and $\mathcal{L} \subseteq \mathcal{I}$ and $\mathcal{O} = \mathcal{I} \setminus \mathcal{L}$. Then $\mathcal{M}$ is solvable w.r.t. $\mathcal{L}$ if and only if for the matrix $A_{\mathcal{L}\mathcal{L}} = \mathbb{I}_{\mathcal{L}} - B_{\mathcal{L}\mathcal{L}}$, for $\mathbb{P}_{\mathcal{E}}$-almost every $\boldsymbol{e} \in \mathcal{E}$ and for all $\boldsymbol{x}_{\mathcal{O}} \in \boldsymbol{\mathcal{X}}_{\mathcal{O}}$ the identity*

$$A_{\mathcal{L}\mathcal{L}} A_{\mathcal{L}\mathcal{L}}^+ (B_{\mathcal{L}\mathcal{O}} \boldsymbol{x}_{\mathcal{O}} + \Gamma_{\mathcal{L}\mathcal{J}} \boldsymbol{e}) = B_{\mathcal{L}\mathcal{O}} \boldsymbol{x}_{\mathcal{O}} + \Gamma_{\mathcal{L}\mathcal{J}} \boldsymbol{e}$$

*is satisfied, where $A_{\mathcal{L}\mathcal{L}}^+$ is the pseudoinverse of $A_{\mathcal{L}\mathcal{L}}$. Moreover, if $\mathcal{M}$ is solvable w.r.t. $\mathcal{L}$, then for every vector $\boldsymbol{v} \in \mathbb{R}^{\mathcal{L}}$ the mapping $\boldsymbol{g}_{\mathcal{L}}^{\boldsymbol{v}} : \mathbb{R}^{\mathcal{O}} \times \mathbb{R}^{\mathcal{J}} \to \mathbb{R}^{\mathcal{L}}$ given by*

$$\boldsymbol{g}_{\mathcal{L}}^{\boldsymbol{v}}(\boldsymbol{x}_{\mathcal{O}}, \boldsymbol{e}) = A_{\mathcal{L}\mathcal{L}}^+ (B_{\mathcal{L}\mathcal{O}} \boldsymbol{x}_{\mathcal{O}} + \Gamma_{\mathcal{L}\mathcal{J}} \boldsymbol{e}) + [\mathbb{I}_{\mathcal{L}} - A_{\mathcal{L}\mathcal{L}}^+ A_{\mathcal{L}\mathcal{L}}] \boldsymbol{v},$$

*is a measurable solution function for $\mathcal{M}$ w.r.t. $\mathcal{L}$.*

For linear SCMs the unique solvability condition w.r.t. a subset translates into a matrix invertibility condition, as was already shown in [27].

PROPOSITION C.3 (Sufficient and necessary condition for unique solvability w.r.t. a subset for linear SCMs). *Let $\mathcal{M}$ be a linear SCM, $\mathcal{L} \subseteq \mathcal{I}$ and $\mathcal{O} = \mathcal{I} \setminus \mathcal{L}$. Then $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$ if and only if the matrix $A_{\mathcal{L}\mathcal{L}} = \mathbb{I}_{\mathcal{L}} - B_{\mathcal{L}\mathcal{L}}$ is invertible. Moreover, if $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$, then the mapping $\boldsymbol{g}_{\mathcal{L}} : \mathbb{R}^{\mathcal{O}} \times \mathbb{R}^{\mathcal{J}} \to \mathbb{R}^{\mathcal{L}}$ given by*

$$\boldsymbol{g}_{\mathcal{L}}(\boldsymbol{x}_{\mathcal{O}}, \boldsymbol{e}) = A_{\mathcal{L}\mathcal{L}}^{-1} (B_{\mathcal{L}\mathcal{O}} \boldsymbol{x}_{\mathcal{O}} + \Gamma_{\mathcal{L}\mathcal{J}} \boldsymbol{e}),$$

*is a measurable solution function for $\mathcal{M}$ w.r.t. $\mathcal{L}$.*

Note that if $A_{\mathcal{L}\mathcal{L}}$ is invertible, then $A_{\mathcal{L}\mathcal{L}}^+ = A_{\mathcal{L}\mathcal{L}}^{-1}$ (see Lemma 1.3 in [54]), and the matrix condition of Proposition C.2 is always satisfied and all the measurable solution functions $\boldsymbol{g}_{\mathcal{L}}^{\boldsymbol{v}}$ of Proposition C.2 are (up to a $\mathbb{P}_{\mathcal{E}}$-null set) equal to the solution function $\boldsymbol{g}_{\mathcal{L}}$ of Proposition C.3.

---

[20] Note that we do not assume that the probability measure $\mathbb{P}_{\mathbb{R}^{\mathcal{J}}}$ is Gaussian.

REMARK. *A sufficient condition for $A_{\mathcal{L}\mathcal{L}}$ to be invertible is that the spectral radius of $B_{\mathcal{L}\mathcal{L}}$ is less than one. If that is the case, then $A_{\mathcal{L}\mathcal{L}}^{-1} = \sum_{n=0}^{\infty}(B_{\mathcal{L}\mathcal{L}})^n$. Note that the non-zero non-diagonal entries of the matrix $B_{\mathcal{L}\mathcal{L}}$ represent the directed edges in the induced subgraph $\mathcal{G}(\mathcal{M})_{\mathcal{L}}$. In particular, if the diagonal entries of the matrix $B_{\mathcal{L}\mathcal{L}}$ are zero, then for $n \in \mathbb{N}$, the coefficients of the matrix $(B_{\mathcal{L}\mathcal{L}})^n$ in the sum represent the sum of the product of the edge weights $B_{ij}$ over directed paths of length $n$ in the induced subgraph $\mathcal{G}(\mathcal{M})_{\mathcal{L}}$.*

From Proposition 3.13 we know that an SCM is solvable w.r.t. $\mathcal{L}$ if and only if it is ancestrally solvable w.r.t. $\mathcal{L}$. In particular, this result also holds for linear SCMs. We saw in Example 3.14 that a similar result for unique solvability does not hold, that is, in general, it does not hold that unique solvability w.r.t. $\mathcal{L}$ implies ancestral unique solvability w.r.t. $\mathcal{L}$. For the class of linear SCMs we do have the following positive result.

PROPOSITION C.4 (Equivalent unique solvability conditions for linear SCMs). *For a linear SCM $\mathcal{M}$ and a subset $\mathcal{L} \subseteq \mathcal{I}$ the following are equivalent:*

1. *$\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$;*
2. *$\mathcal{M}$ is ancestrally uniquely solvable w.r.t. $\mathcal{L}$;*
3. *$\mathcal{M}$ is uniquely solvable w.r.t. each strongly connected component in $\mathcal{G}(\mathcal{M})_{\mathcal{L}}$.*

Under the condition of unique solvability w.r.t. a subset $\mathcal{L}$ we can define the marginalization w.r.t. $\mathcal{L}$ of a linear SCM by mere substitution.

PROPOSITION C.5 (Marginalization of a linear SCM). *Let $\mathcal{M}$ be a linear SCM and $\mathcal{L} \subseteq \mathcal{I}$ a subset of endogenous variables such that $\mathbb{I}_{\mathcal{L}} - B_{\mathcal{L}\mathcal{L}}$ is invertible. Then there exists a marginalization $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ that is linear and with marginal causal mechanism $\tilde{\boldsymbol{f}} : \mathbb{R}^{\mathcal{O}} \times \mathbb{R}^{\mathcal{J}} \to \mathbb{R}^{\mathcal{O}}$ given by*

$$\tilde{\boldsymbol{f}}(\boldsymbol{x}_{\mathcal{O}}, \boldsymbol{e}) = [B_{\mathcal{O}\mathcal{O}} + B_{\mathcal{O}\mathcal{L}}A_{\mathcal{L}\mathcal{L}}^{-1}B_{\mathcal{L}\mathcal{O}}]\boldsymbol{x}_{\mathcal{O}} + [B_{\mathcal{O}\mathcal{L}}A_{\mathcal{L}\mathcal{L}}^{-1}\Gamma_{\mathcal{L}\mathcal{J}} + \Gamma_{\mathcal{O}\mathcal{J}}]\boldsymbol{e},$$

*where $A_{\mathcal{L}\mathcal{L}} = \mathbb{I}_{\mathcal{L}} - B_{\mathcal{L}\mathcal{L}}$. Moreover, this marginalization respects the latent projection, i.e., $\big(\mathcal{G}^a \circ \mathrm{marg}(\mathcal{L})\big)(\mathcal{M}) \subseteq \big(\mathrm{marg}(\mathcal{L}) \circ \mathcal{G}^a\big)(\mathcal{M})$.*

From Theorem 5.6 we know that $\mathcal{M}$ and its marginalization $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ over $\mathcal{L}$ are observationally, interventionally and counterfactually equivalent w.r.t. $\mathcal{O}$. A similar result can also be found in [27]. In contrast to non-linear SCMs, this class of linear SCMs has the convenient property that every marginalization of a model of this class respects the latent projection. Moreover, the subclass of simple linear SCMs is even closed under marginalization.

## APPENDIX D: EXAMPLES

In this appendix we provide additional examples. In Appendix D.1 we provide some examples of SCMs that describe the equilibrium states of certain feedback systems governed by (random) differential equations [6] that motivated our study of cyclic SCMs. In Appendix D.2 we provide some examples that illustrate how observational, interventional and counterfactual equivalence differ.

**D.1. SCMs as equilibrium models** In many systems occurring in the real world feedback loops between observed variables are present. For example, in economics, the price of a product may be a function of the demanded or supplied quantities, and vice versa; or in physics, two masses that are connected by a spring may exert forces on each other. Such systems are often described by a system of (random) differential equations. In [6] it was shown
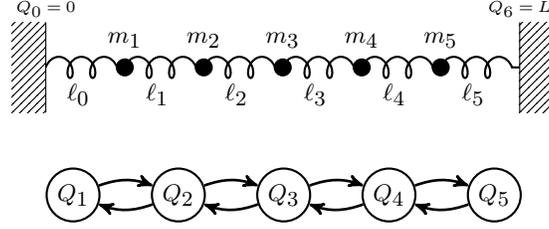
Fig 10: *Damped coupled harmonic oscillator (top) and the graph of the SCM $\mathcal{M}$ that describes the positions of the masses at equilibrium (bottom) of Example D.1 for $d = 5$.*

that SCMs are capable of modeling the causal semantics of the equilibrium states of such systems. For illustration purposes we provide the following toy example of interacting masses that are attached to springs.

EXAMPLE D.1 (Damped coupled harmonic oscillator). *Consider a one-dimensional system of $d$ point masses $m_i \in \mathbb{R}$ ($i = 1, \dots, d$) with positions $Q_i$, which are coupled by springs, with spring constants $k_i > 0$ and equilibrium lengths $\ell_i > 0$ ($i = 0, \dots, d$), under influence of friction with friction coefficients $b_i \in \mathbb{R}$ ($i = 1, \dots, d$) and with fixed endpoints $Q_0 = 0$ and $Q_{d+1} = L > 0$ (see Figure 10 (top)). The equations of motion of this system are provided by the following differential equations*

$$\frac{d^2 Q_i}{dt^2} = \frac{k_i}{m_i}(Q_{i+1} - Q_i - \ell_i) + \frac{k_{i-1}}{m_i}(Q_{i-1} - Q_i + \ell_{i-1}) - \frac{b_i}{m_i}\frac{dQ_i}{dt} \qquad (i = 1, \dots, d).$$

*The dynamics of the masses, in terms of the position, velocity and acceleration, is described by a single and separate equation of motion for each mass. Under friction, i.e., $b_i > 0$ ($i = 1, \dots, d$), there is a unique equilibrium position, where the sum of forces vanishes for each mass. If one starts out of equilibrium, for example, by moving one or several masses out of equilibrium, then the masses will start to oscillate and converge to their unique equilibrium position. At equilibrium (i.e., for $t \to \infty$) the velocity $\frac{dQ_i}{dt}$ and acceleration $\frac{d^2 Q_i}{dt^2}$ of the masses vanish (i.e., $\frac{dQ_i}{dt}, \frac{d^2 Q_i}{dt^2} \to 0$), and thus the following equation holds at equilibrium*

$$0 = \frac{k_i}{m_i}(Q_{i+1} - Q_i - \ell_i) + \frac{k_{i-1}}{m_i}(Q_{i-1} - Q_i + \ell_{i-1}),$$

*for each mass ($i = 1, \dots, d$). Hence, for each mass $i = 1, \dots, d$ its equilibrium position $Q_i$ is given by*

$$Q_i = \frac{k_i(Q_{i+1} - \ell_i) + k_{i-1}(Q_{i-1} + \ell_{i-1})}{k_i + k_{i-1}}.$$

*By considering the $\ell_i$ and $k_i$ and $L$ as fixed parameters, we arrive at a linear SCM (see [6] for more details about constructing an SCM from a dynamical system)*

$$\mathcal{M} = \langle \{1, \dots, d\}, \emptyset, \mathbb{R}^d, \mathbf{1}, \boldsymbol{f}, \mathbb{P}_{\mathbf{1}} \rangle,$$

*where the causal mechanism $\boldsymbol{f}$ is given by*

$$f_i(\boldsymbol{q}) = \frac{k_i(q_{i+1} - \ell_i) + k_{i-1}(q_{i-1} + \ell_{i-1})}{k_i + k_{i-1}}.$$

*Alternatively, (some of) the parameters could be treated as exogenous variables instead. Its graph is depicted in Figure 10 (bottom). This SCM allows us to describe the equilibrium behavior of the system under perfect intervention. For example, when forcing the mass $j$ to a fixed position $Q_j = \xi_j$ with $0 \le \xi_j \le L$, the equilibrium positions of the masses correspond to the solutions of the intervened model $\mathcal{M}_{\mathrm{do}(\{j\}, \xi_j)}$. It is an easy exercise to show that $\mathcal{M}$ is a simple SCM by using Proposition C.3.*
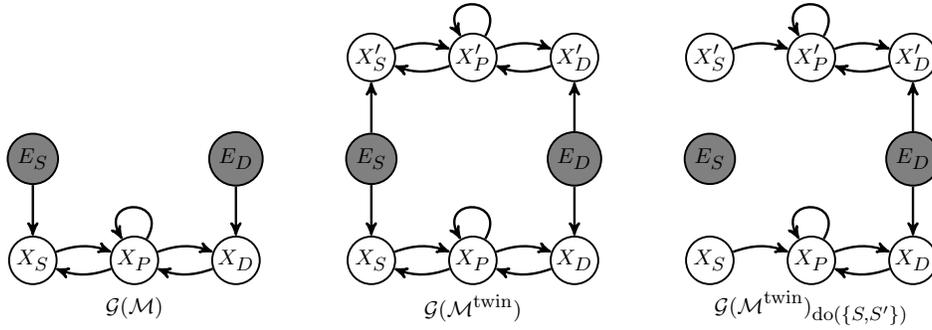
Fig 11: *The augmented graph of the SCM $\mathcal{M}$ (left), its twin SCM $\mathcal{M}^{\text{twin}}$ (center) and the intervened twin SCM $(\mathcal{M}^{\text{twin}})_{\text{do}(\{S,S'\},(s,s'))}$ (right) of Examples D.2 and D.3.*

Next, we show that the well-known market equilibrium model from economics, which has been thoroughly discussed in the literature [see e.g., 65], can be described by a (non-simple) SCM. This example illustrates how self-cycles enrich the class of SCMs.

EXAMPLE D.2 (Price, supply and demand). *Let $X_D$ denote the demand and $X_S$ the supply of a quantity of a product. The price of the product is denoted by $X_P$. The following system of differential equations describes how the demanded and supplied quantities are determined by the price, and how price adjustments occur in the market:*

$$X_D = \beta_D X_P + E_D$$

$$X_S = \beta_S X_P + E_S$$

$$\frac{dX_P}{dt} = X_D - X_S,$$

*where $E_D$ and $E_S$ are exogenous random influences on the demand and supply respectively, $\beta_D < 0$ is the reciprocal of the slope of the demand curve, and $\beta_S > 0$ is the reciprocal of the slope of the supply curve. At the situation known as a "market equilibrium", the price is determined implicitly by the condition that demanded and supplied quantities should be equal, since $\frac{dX_P}{dt} = 0$ at equilibrium. Applying the results in [6] gives rise to a linear SCM $\mathcal{M} = \langle \{P, S, D\}, \{S, D\}, \mathbb{R}^3, \mathbb{R}^2, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\varepsilon}} \rangle$ at equilibrium with the causal mechanism defined by*

$$f_D(\boldsymbol{x}, \boldsymbol{e}) := \beta_D x_P + e_D$$

$$f_S(\boldsymbol{x}, \boldsymbol{e}) := \beta_S x_P + e_S$$

$$f_P(\boldsymbol{x}, \boldsymbol{e}) := x_P + (x_D - x_S).$$

*Note how we use a self-cycle for $P$ in order to implement the equilibrium equation $X_D = X_S$ as the causal mechanism for the price $P$.[21] Moreover, $\mathcal{M}$ is uniquely solvable. Its augmented graph is depicted in Figure 11 (left).*

Next, we provide an example of how counterfactuals can be sensibly formulated for cyclic SCMs, namely for the price, supply and demand model at equilibrium.

---

[21] Richardson and Robins [65] argue that this market equilibrium model cannot be modeled as an SCM. We observe that it can, as long as one allows for self-cycles.
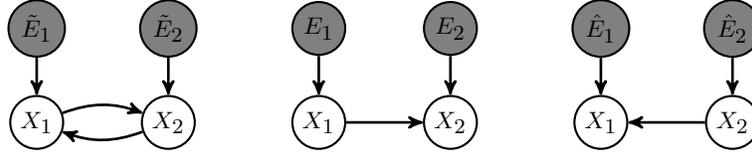
Fig 12: *The augmented graph of the SCM $\tilde{\mathcal{M}}$ (left) of Example 3.7, and of the SCMs $\mathcal{M}$ (center) and $\hat{\mathcal{M}}$ (right) of Example D.4 (and of Example 4.2). The SCMs $\tilde{\mathcal{M}}$, $\mathcal{M}$ and $\hat{\mathcal{M}}$ are all observationally equivalent, but not interventionally equivalent.*

EXAMPLE D.3 (Price, supply and demand at equilibrium). *Consider the price, supply and demand model at equilibrium of Example D.2 given by the SCM $\mathcal{M}$. As an example of a counterfactual query, consider*

$$\mathbb{P}(X_P' \,|\, \mathrm{do}(X_S = s, X_{S'} = s'), X_P = p)\,,$$

*which denotes the conditional distribution of $X_P'$ given $X_P = p$ of a solution of the intervened twin model $\mathcal{M}^{\mathrm{twin}}_{\mathrm{do}(\{S,S'\},(s,s'))}$. In words: how would—ceteris paribus—price have been distributed, had we intervened to set supplied quantities equal to $s'$, given that actually we intervened to set supplied quantities equal to $s$ and observed that this led to price $p$? A straightforward calculation shows that this counterfactual distribution of price is the Dirac measure on $x_P' = p + (s' - s)/\beta_D$. The augmented graphs of the SCM, its twin graph, and its intervened twin graph are depicted in Figure 11.*

**D.2. Equivalences**   In this subsection we provide some examples that illustrate how observational, interventional and counterfactual equivalence differ.

The following example illustrates that observational equivalence does not imply equivalence and interventional equivalence (see also Example 4.2).

EXAMPLE D.4 (Observational equivalence does not imply (interventional) equivalence). *Consider the SCM $\tilde{\mathcal{M}}$ of Example 3.7 and let $\mathcal{M} = \langle \mathbf{2}, \mathbf{2}, \mathbb{R}^2, \mathbb{R}^2, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be the SCM with the causal mechanism*

$$f_1(\boldsymbol{x}, \boldsymbol{e}) = e_1\,, \quad f_2(\boldsymbol{x}, \boldsymbol{e}) = \gamma x_1 + e_2\,,$$

*where*

$$\gamma = \frac{\beta \sigma_1^2 + \alpha \sigma_2^2}{\sigma_1^2 + \alpha^2 \sigma_2^2}\,,$$

*and $\mathbb{P}_{\boldsymbol{\mathcal{E}}} = \mathbb{P}^{\boldsymbol{E}}$ with $E_1 \sim \mathcal{N}(\bar{\mu}_1, \bar{\sigma}_1^2)$, $E_2 \sim \mathcal{N}(\bar{\mu}_2, \bar{\sigma}_2^2)$ and $E_1 \perp\!\!\!\perp E_2$, where*

$$\bar{\mu}_1 = c[\mu_1 + \alpha \mu_2], \qquad\qquad \bar{\sigma}_1^2 = c^2[\sigma_1^2 + \alpha^2 \sigma_2^2],$$

$$\bar{\mu}_2 = c[(\beta - \gamma)\mu_1 + (1 - \alpha\gamma)\mu_2], \quad \bar{\sigma}_2^2 = c^2[(\beta - \gamma)^2 \sigma_1^2 + (1 - \alpha\gamma)^2 \sigma_2^2]$$

*with $c = (1 - \alpha\beta)^{-1}$. The augmented graphs of $\tilde{\mathcal{M}}$ and $\mathcal{M}$ are depicted in Figure 12. The SCMs $\tilde{\mathcal{M}}$ and $\mathcal{M}$ are observationally equivalent, as one can check by explicit calculation. Similarly, one can define an SCM $\hat{\mathcal{M}}$ with augmented graph as depicted in Figure 12 that is observationally equivalent to both $\tilde{\mathcal{M}}$ and $\mathcal{M}$. Because each of the SCMs has a different augmented graph, we conclude that none of the SCMs $\tilde{\mathcal{M}}$, $\mathcal{M}$ and $\hat{\mathcal{M}}$ are equivalent to each other. Although $\tilde{\mathcal{M}}$, $\mathcal{M}$ and $\hat{\mathcal{M}}$ are observationally equivalent, none of them is interventionally equivalent to each other, as one can easily check.*

In general, interventional equivalence does not imply counterfactual equivalence. Even interventionally equivalent SCMs with the same causal mechanism (that differ only in their exogenous distribution) may not be counterfactually equivalent. For example, the SCMs $\mathcal{M}_\rho$ and $\mathcal{M}_{\rho'}$ with $\rho \neq \rho'$ in the following example (due to Dawid [11]) are interventionally but not counterfactually equivalent.
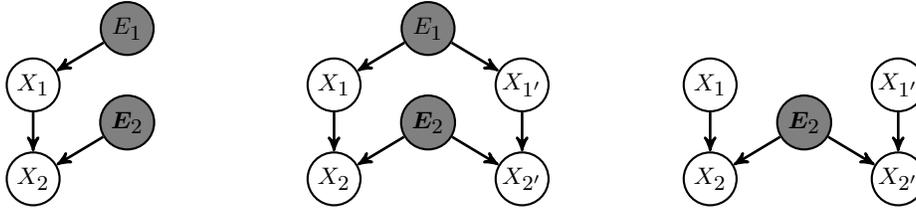
Fig 13: *The augmented graph of the SCM $\mathcal{M}_\rho$ (left), its twin SCM $\mathcal{M}_\rho^{\text{twin}}$ (center) and the intervened twin SCM $(\mathcal{M}_\rho^{\text{twin}})_{\text{do}(\{1',1\},(1,0))}$ (right) of Example D.5.*

EXAMPLE D.5 (Counterfactual density unidentifiable from observational and interventional densities [11]). *Let $\rho \in \mathbb{R}$ and*

$$\mathcal{M}_\rho = \langle \mathbf{2}, \mathbf{2}, \{0,1\} \times \mathbb{R}, \{0,1\} \times \mathbb{R}^2, \boldsymbol{f}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$$

*be the SCM with causal mechanism given by*

$$f_1(\boldsymbol{x}, \boldsymbol{e}) = e_1\,, \quad f_2(\boldsymbol{x}, \boldsymbol{e}) = e_{21}(1 - x_1) + e_{22}x_1$$

*and $\mathbb{P}_{\boldsymbol{\mathcal{E}}} = \mathbb{P}^{(E_1, \boldsymbol{E}_2)}$ with $E_1 \sim \text{Bernoulli}(1/2)$,*

$$\boldsymbol{E}_2 := \begin{pmatrix} E_{21} \\ E_{22} \end{pmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$$

*normally distributed and $E_1 \perp\!\!\!\perp \boldsymbol{E}_2$. In an epidemiological setting, this SCM could be used to model whether a patient was treated or not ($X_1$) and the corresponding outcome for that patient ($X_2$).*

*Suppose in the actual world we did not assign treatment to a patient ($X_1 = 0$) and the outcome was $X_2 = c \in \mathbb{R}$. Consider the counterfactual query "What would the outcome have been, if we had assigned treatment to this patient?". We can answer this question by introducing a parallel counterfactual world that is modeled by the twin SCM $\mathcal{M}_\rho^{\text{twin}}$, as depicted in Figure 13. The counterfactual query then asks for $p(X_{2'} = x_{2'} \mid \text{do}(X_{1'} = 1, X_1 = 0), X_2 = c)$. One can calculate that*

$$\begin{pmatrix} X_{2'} \\ X_2 \end{pmatrix} \mid \text{do}(X_{1'} = 1, X_1 = 0) \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$$

*and hence $X_{2'} \mid \text{do}(X_{1'} = 1, X_1 = 0), X_2 = c \sim \mathcal{N}(-\rho c, 1 - \rho^2)$. Note that the answer to the counterfactual query depends on a quantity $\rho$ that we cannot identify from the observational density $p(X_1, X_2)$ or the interventional densities $p(X_2 \mid \text{do}(X_1 = 0))$ and $p(X_2 \mid \text{do}(X_1 = 1))$, none of which depends on $\rho$. Therefore, even data from randomized controlled trials combined with observational data would not suffice to determine the value of this particular counterfactual query.*

## APPENDIX E: PROOFS

This appendix contains the proofs of all the theoretical results in the appendices A, B and C, and the main text. Some of the proofs will rely on the measure theoretic terminology and results of Appendix F.

### E.1. Proofs of the appendices

*Appendix A*

PROOF OF LEMMA A.5. It suffices to show that for every $C$-$d$-open walk between $i$ and $j$ in $\mathcal{G}$, there exists a $C$-$d$-open path between $i$ and $j$ in $\mathcal{G}$. Take a $C$-$d$-open walk $\pi = (i = i_0, \ldots, i_n = j)$. If a node $\ell$ occurs more than once in $\pi$, let $i_j$ be the first occurrence of $\ell$ in $\pi$ and $i_k$ the last occurrence of $\ell$ in $\pi$. We now construct a new walk $\pi'$ from $\pi$ by removing the subwalk between $i_j$ and $i_k$ of $\pi$ from $\pi$. It is easy to check that the new walk $\pi'$ is still $C$-$d$-open. If $\ell$ is an endpoint on $\pi'$, then $i_j$ or $i_k$ must be endpoint of $\pi$, and hence $\ell \notin C$. If $\ell$ is a non-endpoint non-collider on $\pi'$, then also $i_j$ or $i_k$ must have been a non-endpoint non-collider on $\pi$, and hence $\ell \notin C$. If $\ell$ is a collider on $\pi'$, then either (i) $i_j$ or $i_k$ are both colliders on $\pi$, and hence $\ell$ is ancestor of $C$ in $\mathcal{G}$, or (ii) on the subwalk between $i_j$ and $i_k$ that was removed, there must be a directed path in $\mathcal{G}$ from $i_j$ or $i_k$ to a collider in $\mathrm{an}_{\mathcal{G}}(C)$, and hence, $\ell$ is in $\mathrm{an}_{\mathcal{G}}(C)$. The other nodes on $\pi'$ cannot be responsible for $C$-$d$-blocking the walk, since they also occur (together with their adjacent edges) on $\pi$ and they do not $C$-$d$-block $\pi$.

In $\pi'$, the number of nodes that occur multiple times is at least one less than in $\pi$. Repeat this procedure until no repeated nodes are left. □

PROOF OF THEOREM A.7. The first case is a well-known result. An elementary proof is obtained by noting that an acyclic system of structural equations trivially satisfies the local directed Markov property, and then apply [35, Proposition 4], followed by applying the stability of $d$-separation with respect to (graphical) marginalization [18, Lemma 2.2.15]. Alternatively, the result also follows from sequential application of Theorems 3.8.2, 3.8.11, 3.7.7, 3.7.2 and 3.3.3 (using Remark 3.3.4) in [18].

The discrete case is proved by the series of results Theorem 3.8.12, Remark 3.7.2, Theorem 3.6.6 and 3.5.2 in [18].

The linear case is proved in Example 3.8.17 in [18]. To connect the assumptions made there with the ones we state here, observe that under the linear transformation rule for Lebesgue measures, the image measure of $\mathbb{P}_{\mathcal{E}}$ under the linear mapping $\mathbb{R}^{\mathcal{J}} \to \mathbb{R}^{\mathcal{I}} : e \mapsto \Gamma_{\mathcal{I}\mathcal{J}} e$ gives a measure on $\mathcal{X} = \mathbb{R}^{\mathcal{I}}$ with a density w.r.t. the Lebesgue measure on $\mathbb{R}^{\mathcal{I}}$, as long as the image of the linear mapping is the entire $\mathbb{R}^{\mathcal{I}}$. This is guaranteed if each causal mechanism has a non-trivial dependence on some exogenous variable(s), i.e., for each $i \in \mathcal{I}$ there is some $j \in \mathcal{J}$ with $\Gamma_{ij} \neq 0$. □

PROOF OF PROPOSITION A.12. This follows directly from the fact that the strongly connected components of $\mathcal{G}^a(\mathcal{M})$ form a DAG by Lemma A.2 and that the directed edges in $\mathcal{G}^a(\mathrm{acy}(\mathcal{M}))$ by construction respect every topological ordering of that DAG. Both SCMs are observationally equivalent by construction. □

PROOF OF PROPOSITION A.14. This follows immediately from the Definitions A.11 and A.13. □

PROOF OF LEMMA A.17. It suffices to show that for every $C$-$\sigma$-open walk between $i$ and $j$ in $\mathcal{G}$, there exists a $C$-$\sigma$-open path between $i$ and $j$ in $\mathcal{G}$. Let $\pi = (i = i_0, \ldots, i_n = j)$ be a $C$-$\sigma$-open walk in $\mathcal{G}$. If a node $\ell$ occurs more than once in $\pi$, let $i_j$ be the first node in $\pi$ and $i_k$ the last node in $\pi$ that are in the same strongly connected component as $\ell$. Since $i_j$ and $i_k$ are in the same strongly connected component, there are directed paths $i_j \to \cdots \to i_k$ and $i_k \to \cdots \to i_j$ in $\mathcal{G}$. We now construct a new walk $\pi'$ from $\pi$ by replacing the subwalk between $i_j$ and $i_k$ of $\pi$ by a particular directed path between $i_j$ and $i_k$: (i) If $k = n$, or if $k < n$ and $i_k \to i_{k+1}$ on $\pi$, we replace it by a shortest directed path $i_j \to \cdots \to i_k$, otherwise

(ii) we replace it by a shortest directed path $i_j \leftarrow \cdots \leftarrow i_k$. We now show that the new walk $\pi'$ is still $C$-$\sigma$-open.

$\pi'$ cannot become $C$-$\sigma$-blocked through one of the initial nodes $i_0 \ldots i_{j-1}$ or one of the final nodes $i_{k+1} \ldots i_n$ on $\pi'$, since these nodes occur in the same local configuration on $\pi$ and do not $C$-$\sigma$-block $\pi$ by assumption. Furthermore, $\pi'$ cannot become $C$-$\sigma$-blocked through one of the nodes strictly between $i_j$ and $i_k$ on $\pi'$ (if there are any), since these nodes are all non-endpoint non-colliders that only point to nodes in the same strongly connected component on $\pi'$. Because $\pi$ is $C$-$\sigma$-open, $i_k \notin C$ if $k = n$ or if $i_k \rightarrow i_{k+1}$ on $\pi$. This holds in particular in case (i). Similarly, $i_j \notin C$ if $j = 0$ or $i_{j-1} \leftarrow i_j$ on $\pi$.

In case (i), $\pi'$ is not $C$-$\sigma$-blocked by $i_k$ because $i_k$ is a non-collider on $\pi'$ but $i_k \notin C$. Also $i_j$ does not $C$-$\sigma$-block $\pi'$. Assume $i_j \neq i_k$ (otherwise there is nothing to prove). If $j = 0$, or if $j > 0$ and $i_{j-1} \leftarrow i_j$ on $\pi'$, then the same holds for $\pi$ and hence $i_j \notin C$; $i_j$ is then a non-collider on $\pi'$, but $i_j \notin C$. If $j > 0$ and $i_{j-1} \leftrightarrow i_j$ or $i_{j-1} \rightarrow i_j$ on $\pi'$ then $i_j$ is a non-endpoint non-collider on $\pi'$ that does not point to a node in another strongly connected component.

Now consider case (ii). If $j = 0$ or $i_{j-1} \leftarrow i_j$ on $\pi'$ then this case is analogous to case (i). So assume $j > 0$ and $i_{j-1} \rightarrow i_j$ or $i_{j-1} \leftrightarrow i_j$ on $\pi'$. If $i_j$ is an endpoint of $\pi'$, then $i_j = i_k$ and $k = n$ and therefore $i_k \notin C$, and hence $i_j$ and $i_k$ do not $C$-$\sigma$-block $\pi'$. Otherwise, $i_j$ must be a collider on $\pi'$ (whether $i_j = i_k$ or not). Then on the subwalk of $\pi$ between $i_j$ and $i_k$ there must be a directed path from $i_j$ to a collider that is ancestor of $C$, which implies that $i_j$ is itself ancestor of $C$, and hence $i_j$ does not $C$-$\sigma$-block $\pi'$. Also $i_k$ cannot $C$-$\sigma$-block $\pi'$. Assume $i_j \neq i_k$ (otherwise there is nothing to prove). Since $i_k \leftarrow i_{k+1}$ or $i_k \leftrightarrow i_{k+1}$ on $\pi'$, $i_k$ is a non-endpoint non-collider on $\pi'$ that does not point to a node in another strongly connected component.

Now in $\pi'$, the number of nodes that occurs more than once is at least one less than in $\pi$. Repeat this procedure until no nodes occur more than once. $\square$

PROOF OF PROPOSITION A.19. This follows directly as a special case of Corollary 2.8.4 in [18]. $\square$

PROOF OF THEOREM A.21. An SCM $\mathcal{M}$ that is uniquely solvable w.r.t. each strongly connected component is uniquely solvable and hence, by Theorem 3.8, all its solutions have the same observational distribution. The last statement follows from the series of results Theorem 3.8.2, 3.8.11, Lemma 3.7.7 and Remark 3.7.2 in [18]. Alternatively, we give here a shorter proof: Under the stated conditions one can always construct the acyclification $\mathrm{acy}(\mathcal{M})$ which is observationally equivalent to $\mathcal{M}$ and is acyclic (see Proposition A.12) and hence we can apply Theorem A.7 to $\mathrm{acy}(\mathcal{M})$. Together with Proposition A.14 and A.19 this gives

$$A \underset{\mathcal{G}(\mathcal{M})}{\overset{\sigma}{\perp}} B \,|\, C \iff A \underset{\mathrm{acy}(\mathcal{G}(\mathcal{M}))}{\overset{d}{\perp}} B \,|\, C \implies A \underset{\mathcal{G}(\mathrm{acy}(\mathcal{M}))}{\overset{d}{\perp}} B \,|\, C \implies \boldsymbol{X}_A \underset{\mathbb{P}_{\mathcal{M}}^{\boldsymbol{X}}}{\perp} \boldsymbol{X}_B \,|\, \boldsymbol{X}_C,$$

for $A, B, C \subseteq \mathcal{I}$ and $\boldsymbol{X}$ a solution of $\mathcal{M}$. $\square$

PROOF OF COROLLARY A.22. First observe that simplicity is preserved under both perfect intervention and the twin operation (see Proposition 8.2). Now the first statement follows from Theorem A.21 if one takes into account the identities of Proposition 2.15 and 2.20. Similarly, the last statement follows from Theorem A.7. $\square$

PROOF OF PROPOSITION A.33. Let $\tilde{\mathcal{M}} =: \langle \mathcal{V}, \hat{\mathcal{H}}, \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{E}}, \tilde{\boldsymbol{f}}, \mathbb{P}_{\boldsymbol{\mathcal{E}}} \rangle$ be the induced SCM. Observe that every loop $\mathcal{O} \in \mathcal{L}(\mathcal{G}(\tilde{\mathcal{M}}))$ is a loop in $\mathcal{L}(\mathcal{G})$. Fix $\check{x} \in \boldsymbol{\mathcal{X}}$ and $\check{e} \in \boldsymbol{\mathcal{E}}$. For every $\mathcal{O} \in \mathcal{L}(\mathcal{G}(\tilde{\mathcal{M}}))$, define

$$I_{\mathcal{O}} := (\mathrm{pa}_{\mathcal{G}}(\mathcal{O}) \setminus \mathcal{O}) \setminus (\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}) \subseteq \tilde{\mathcal{I}}$$

and

$$J_{\mathcal{O}} := \{\mathcal{F} \in \tilde{\mathcal{J}} \, : \, \mathcal{F} \cap \mathcal{O} \neq \emptyset\} \setminus \mathrm{pa}(\mathcal{O}) \subseteq \tilde{\mathcal{J}} \, .$$

Now, define the family of measurable mappings $(\tilde{\boldsymbol{g}}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G}(\tilde{\mathcal{M}}))}$, where the mapping $\tilde{\boldsymbol{g}}_{\mathcal{O}}$ : $\boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{O})} \to \boldsymbol{\mathcal{X}}_{\mathcal{O}}$ is given by

$$\tilde{\boldsymbol{g}}_{\mathcal{O}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}) := \boldsymbol{g}_{\mathcal{O}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \check{\boldsymbol{x}}_{I_{\mathcal{O}}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}, \check{\boldsymbol{e}}_{J_{\mathcal{O}}})$$

where $\boldsymbol{x}_{\mathrm{pa}_{\mathcal{G}}(\mathcal{O}) \setminus \mathcal{O}} = (\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \check{\boldsymbol{x}}_{I_{\mathcal{O}}})$ and $\widehat{\boldsymbol{e}}_{\mathcal{O}} = (\boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}, \check{\boldsymbol{e}}_{J_{\mathcal{O}}})$. Observe that from the definition of the parents (see Definition 2.7) it follows that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$ we have

$$\boldsymbol{x}_{\mathcal{O}} = \tilde{\boldsymbol{f}}_{\mathcal{O}}(\boldsymbol{x}_{\setminus I_{\mathcal{O}}}, \check{\boldsymbol{x}}_{I_{\mathcal{O}}}, \boldsymbol{e}_{\setminus J_{\mathcal{O}}}, \check{\boldsymbol{e}}_{J_{\mathcal{O}}}) \quad \Longleftrightarrow \quad \boldsymbol{x}_{\mathcal{O}} = \tilde{\boldsymbol{f}}_{\mathcal{O}}(\boldsymbol{x}, \boldsymbol{e}) \, .$$

This, together with the fact that the family of mappings $(\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G})}$ is a compatible system of solution functions, implies that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$ we have

$$\boldsymbol{x}_{\mathcal{O}} = \tilde{\boldsymbol{g}}_{\mathcal{O}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}) \quad \Longrightarrow \quad \boldsymbol{x}_{\mathcal{O}} = \tilde{\boldsymbol{f}}_{\mathcal{O}}(\boldsymbol{x}, \boldsymbol{e}) \, .$$

Hence, $\iota(\widehat{\mathcal{M}})$ is loop-wisely solvable and thus $(\tilde{\boldsymbol{g}}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G}(\tilde{\mathcal{M}}))}$ is a family of measurable solution functions. In particular, for all $\mathcal{O}, \tilde{\mathcal{O}} \in \mathcal{L}(\mathcal{G}(\tilde{\mathcal{M}}))$ with $\tilde{\mathcal{O}} \subseteq \mathcal{O}$ and for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$ we have

$$\boldsymbol{x}_{\mathcal{O}} = \tilde{\boldsymbol{g}}_{\mathcal{O}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}) \quad \Longrightarrow \quad \boldsymbol{x}_{\tilde{\mathcal{O}}} = \tilde{\boldsymbol{g}}_{\tilde{\mathcal{O}}}(\boldsymbol{x}_{\mathrm{pa}(\tilde{\mathcal{O}}) \setminus \tilde{\mathcal{O}}}, \boldsymbol{e}_{\mathrm{pa}(\tilde{\mathcal{O}})}) \, .$$

From this we conclude that $(\tilde{\boldsymbol{g}}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G}(\tilde{\mathcal{M}}))}$ is a compatible system of solution functions. $\qquad \square$

PROOF OF LEMMA A.34. Suppose $\mathcal{M}$ is loop-wisely uniquely solvable and consider a subset $\mathcal{O} \subseteq \mathcal{I}$. Consider the induced subgraph $\mathcal{G}^a(\mathcal{M})_{\mathcal{O}}$ of $\mathcal{G}^a(\mathcal{M})$ on the nodes $\mathcal{O}$. Then every strongly connected component of $\mathcal{G}^a(\mathcal{M})_{\mathcal{O}}$ is an element of $\mathcal{L}(\mathcal{G}(\mathcal{M}))$. Let $\mathcal{C}$ be such a strongly connected component in $\mathcal{G}^a(\mathcal{M})_{\mathcal{O}}$, and let $\boldsymbol{g}_{\mathcal{C}} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{C}) \setminus \mathcal{C}} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{C})} \to \boldsymbol{\mathcal{X}}_{\mathcal{C}}$ be a measurable solution function for $\mathcal{M}$ w.r.t. $\mathcal{C}$. Since $\mathcal{G}^a(\mathcal{M})_{\mathcal{O}}$ partitions into strongly connected components, we can recursively (by following a topological ordering of the DAG $\mathcal{G}^a(\mathcal{M})_{\mathcal{O}}^{\mathrm{sc}}$ from Lemma A.2) insert these mappings into each other to obtain a mapping $\boldsymbol{g}_{\mathcal{O}} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{O})} \to \boldsymbol{\mathcal{X}}_{\mathcal{O}}$ that makes $\mathcal{M}$ uniquely solvable w.r.t. $\mathcal{O}$. $\qquad \square$

PROOF OF PROPOSITION A.35. Let $(\boldsymbol{g}_{\mathcal{O}})_{\mathcal{O} \in \mathcal{L}(\mathcal{G}(\mathcal{M}))}$ be any family of measurable solution functions, where $\boldsymbol{g}_{\mathcal{O}}$ is measurable solution function of $\mathcal{M}$ w.r.t. $\mathcal{O}$. Then, for $\mathcal{O}, \tilde{\mathcal{O}} \in \mathcal{L}(\mathcal{G}(\mathcal{M}))$ such that $\tilde{\mathcal{O}} \subseteq \mathcal{O}$, we have that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$

$$\boldsymbol{x}_{\mathcal{O}} = \boldsymbol{f}_{\mathcal{O}}(\boldsymbol{x}, \boldsymbol{e}) \quad \Longrightarrow \quad \boldsymbol{x}_{\tilde{\mathcal{O}}} = \boldsymbol{f}_{\tilde{\mathcal{O}}}(\boldsymbol{x}, \boldsymbol{e}) \, .$$

This implies that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$

$$\boldsymbol{x}_{\mathcal{O}} = \boldsymbol{g}_{\mathcal{O}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}) \quad \Longrightarrow \quad \boldsymbol{x}_{\tilde{\mathcal{O}}} = \boldsymbol{g}_{\tilde{\mathcal{O}}}(\boldsymbol{x}_{\mathrm{pa}(\tilde{\mathcal{O}}) \setminus \tilde{\mathcal{O}}}, \boldsymbol{e}_{\mathrm{pa}(\tilde{\mathcal{O}})}) \, .$$

$\qquad \square$

PROOF OF COROLLARY A.24. This follows directly from Proposition 7.1 and 7.3. $\qquad \square$

*Appendix B*

PROOF OF PROPOSITION B.1. Let $\tilde{f} : \mathcal{E} \times \mathcal{X} \to \mathcal{X}$ be the causal mechanism of a structurally minimal SCM that is equivalent to $\mathcal{M}$ (see Proposition 2.12). In particular, for any $\epsilon_{\backslash \mathrm{pa}(\mathcal{O})} \in \mathcal{E}_{\backslash \mathrm{pa}(\mathcal{O})}$ and $\xi_{\backslash \mathrm{pa}(\mathcal{O})} \in \mathcal{X}_{\backslash \mathrm{pa}(\mathcal{O})}$, we have that for all $x \in \mathcal{X}$ and all $e \in \mathcal{E}$, $\tilde{f}(x, e) = \tilde{f}(x_{\mathrm{pa}(\mathcal{O})}, \xi_{\backslash \mathrm{pa}(\mathcal{O})}, e_{\mathrm{pa}(\mathcal{O})}, \epsilon_{\backslash \mathrm{pa}(\mathcal{O})})$. This means that we may also consider $\tilde{f}$ as a mapping $\tilde{f} : \mathcal{X}_{\mathrm{pa}(\mathcal{O})} \times \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}$.

Consider the set

$$\tilde{\mathcal{S}} := \{ (e_{\mathrm{pa}(\mathcal{O})}, x_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}}, x_{\mathcal{O}}) \in \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \times \mathcal{X}_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}} \times \mathcal{X}_{\mathcal{O}} : x_{\mathcal{O}} = \tilde{f}_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O})}, e_{\mathrm{pa}(\mathcal{O})}) \} .$$

By similar reasoning as in the proof of Theorem 3.3, $\tilde{\mathcal{S}}$ is measurable.

By assumption, for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x_{\backslash \mathcal{O}} \in \mathcal{X}_{\backslash \mathcal{O}}$ the space $\{ x_{\mathcal{O}} \in \mathcal{X}_{\mathcal{O}} : x_{\mathcal{O}} = f_{\mathcal{O}}(x, e) \}$ is non-empty and $\sigma$-compact. By applying Lemma F.10 to the canonical projection $pr_{\mathcal{E}_{\mathrm{pa}(\mathcal{O})}} : \mathcal{E} \to \mathcal{E}_{\mathrm{pa}(\mathcal{O})}$ and using the equivalence of $f$ and $\tilde{f}$, we obtain that for $\mathbb{P}_{\mathcal{E}_{\mathrm{pa}(\mathcal{O})}}$-almost every $e_{\mathrm{pa}(\mathcal{O})} \in \mathcal{E}_{\mathrm{pa}(\mathcal{O})}$ and for all $x_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}} \in \mathcal{X}_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}}$ the space

$$\tilde{\mathcal{S}}_{(e_{\mathrm{pa}(\mathcal{O})}, x_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}})} := \{ x_{\mathcal{O}} \in \mathcal{X}_{\mathcal{O}} : x_{\mathcal{O}} = \tilde{f}_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O})}, e_{\mathrm{pa}(\mathcal{O})}) \}$$

is non-empty and $\sigma$-compact.

The second measurable selection theorem, Theorem F.9, now implies that there exists a measurable $g_{\mathcal{O}} : \mathcal{X}_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}_{\mathcal{O}}$ such that for $\mathbb{P}_{\mathcal{E}_{\mathrm{pa}(\mathcal{O})}}$-almost every $e_{\mathrm{pa}(\mathcal{O})} \in \mathcal{E}_{\mathrm{pa}(\mathcal{O})}$ and for all $x_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}} \in \mathcal{X}_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}}$

$$g_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}) = \tilde{f}_{\mathcal{O}}\big(x_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}}, g_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}), e_{\mathrm{pa}(\mathcal{O})}\big).$$

Once more applying Lemma F.10, we obtain that for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$x_{\mathcal{O}} = g_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O}) \backslash \mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}) \implies x_{\mathcal{O}} = f_{\mathcal{O}}(x, e).$$

Hence $\mathcal{M}$ is solvable w.r.t. $\mathcal{O}$. $\qquad\square$

PROOF OF PROPOSITION B.4. Without loss of generality, we assume that $\mathcal{M}$ is structurally minimal (see Proposition 2.12). Define $\mathcal{C} := \mathcal{A} \cap \tilde{\mathcal{A}}$ and $\mathcal{D} := \mathcal{A} \cup \tilde{\mathcal{A}}$. Let $g_{\mathcal{A}}$, $g_{\tilde{\mathcal{A}}}$ be measurable solution functions for $\mathcal{M}$ w.r.t. $\mathcal{A}$ and $\tilde{\mathcal{A}}$, respectively. Note that $\mathrm{pa}(\mathcal{C}) \backslash \mathcal{C} \subseteq \mathrm{pa}(\mathcal{A}) \backslash \mathcal{A}$ and similarly $\mathrm{pa}(\mathcal{C}) \backslash \mathcal{C} \subseteq \mathrm{pa}(\tilde{\mathcal{A}}) \backslash \tilde{\mathcal{A}}$. Indeed, for $c \in \mathrm{pa}(\mathcal{C})$: if $c \in \mathcal{O}$ then $c \in \mathcal{C}$ because $\mathcal{A}$ and $\tilde{\mathcal{A}}$ are both ancestral in $\mathcal{G}(\mathcal{M})_{\mathcal{O}}$, while if $c \notin \mathcal{O}$ then $c \notin \mathcal{A}$ and $c \notin \tilde{\mathcal{A}}$. Hence by Lemma E.1, for $\mathbb{P}_{\mathcal{E}}$-almost all $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$(g_{\mathcal{A}})_{\mathcal{C}}(x_{\mathrm{pa}(\mathcal{A}) \backslash \mathcal{A}}, e_{\mathrm{pa}(\mathcal{A})}) = (g_{\tilde{\mathcal{A}}})_{\mathcal{C}}(x_{\mathrm{pa}(\tilde{\mathcal{A}}) \backslash \tilde{\mathcal{A}}}, e_{\mathrm{pa}(\tilde{\mathcal{A}})}) .$$

Hence for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$x_{\mathcal{D}} = f_{\mathcal{D}}(x, e)$$

$$\iff \begin{cases} x_{\mathcal{A}\setminus\mathcal{C}} &= f_{\mathcal{A}\setminus\mathcal{C}}(x, e) \\ x_{\mathcal{C}} &= f_{\mathcal{C}}(x, e) \\ x_{\mathcal{C}} &= f_{\mathcal{C}}(x, e) \\ x_{\tilde{\mathcal{A}}\setminus\mathcal{C}} &= f_{\tilde{\mathcal{A}}\setminus\mathcal{C}}(x, e) \end{cases}$$

$$\iff \begin{cases} x_{\mathcal{A}\setminus\mathcal{C}} &= (g_{\mathcal{A}})_{\mathcal{A}\setminus\mathcal{C}}(x_{\mathrm{pa}(\mathcal{A})\setminus\mathcal{A}}, e_{\mathrm{pa}(\mathcal{A})}) \\ x_{\mathcal{C}} &= (g_{\mathcal{A}})_{\mathcal{C}}(x_{\mathrm{pa}(\mathcal{A})\setminus\mathcal{A}}, e_{\mathrm{pa}(\mathcal{A})}) \\ x_{\mathcal{C}} &= (g_{\tilde{\mathcal{A}}})_{\mathcal{C}}(x_{\mathrm{pa}(\tilde{\mathcal{A}})\setminus\tilde{\mathcal{A}}}, e_{\mathrm{pa}(\tilde{\mathcal{A}})}) \\ x_{\tilde{\mathcal{A}}\setminus\mathcal{C}} &= (g_{\tilde{\mathcal{A}}})_{\tilde{\mathcal{A}}\setminus\mathcal{C}}(x_{\mathrm{pa}(\tilde{\mathcal{A}})\setminus\tilde{\mathcal{A}}}, e_{\mathrm{pa}(\tilde{\mathcal{A}})}) \end{cases}$$

$$\iff \begin{cases} x_{\mathcal{A}} &= g_{\mathcal{A}}(x_{\mathrm{pa}(\mathcal{A})\setminus\mathcal{A}}, e_{\mathrm{pa}(\mathcal{A})}) \\ x_{\tilde{\mathcal{A}}} &= g_{\tilde{\mathcal{A}}}(x_{\mathrm{pa}(\tilde{\mathcal{A}})\setminus\tilde{\mathcal{A}}}, e_{\mathrm{pa}(\tilde{\mathcal{A}})}). \end{cases}$$

Now $\mathrm{pa}(\mathcal{A}) \setminus \mathcal{A} \subseteq \mathrm{pa}(\mathcal{D}) \setminus \mathcal{D}$, and similarly, $\mathrm{pa}(\tilde{\mathcal{A}}) \setminus \tilde{\mathcal{A}} \subseteq \mathrm{pa}(\mathcal{D}) \setminus \mathcal{D}$. Hence, we conclude that the mapping $h_{\mathcal{D}} : \mathcal{X}_{\mathrm{pa}(\mathcal{D})\setminus\mathcal{D}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{D})} \to \mathcal{X}_{\mathcal{D}}$ defined by

$$h_{\mathcal{D}}(x_{\mathrm{pa}(\mathcal{D})\setminus\mathcal{D}}, e_{\mathrm{pa}(\mathcal{D})}) :=$$

$$\left( (g_{\mathcal{A}})_{\mathcal{A}\setminus\mathcal{C}}(x_{\mathrm{pa}(\mathcal{A})\setminus\mathcal{A}}, e_{\mathrm{pa}(\mathcal{A})}), (g_{\mathcal{A}})_{\mathcal{C}}(x_{\mathrm{pa}(\mathcal{A})\setminus\mathcal{A}}, e_{\mathrm{pa}(\mathcal{A})}), (g_{\tilde{\mathcal{A}}})_{\tilde{\mathcal{A}}\setminus\mathcal{C}}(x_{\mathrm{pa}(\tilde{\mathcal{A}})\setminus\tilde{\mathcal{A}}}, e_{\mathrm{pa}(\tilde{\mathcal{A}})}) \right)$$

is a measurable solution function for $\mathcal{M}$ w.r.t. $\mathcal{D}$, and that $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{D}$. $\qquad\square$

PROOF OF COROLLARY B.5. It suffices to show the implication to the left. We have to show that $\mathcal{M}$ is uniquely solvable w.r.t. each ancestral subset of $\mathcal{G}(\mathcal{M})_{\mathcal{O}}$. The proof proceeds via induction with respect to the size of the ancestral subset. For ancestral subsets of size 0, the claim is trivially true. Ancestral subsets of size 1 must be of the form $\{i\} = \mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(i)$ for $i \in \mathcal{O}$ and hence the claim is true by assumption. Assume that the claim holds for all ancestral subsets of size $\leq n$. Let $\mathcal{A}$ be an ancestral subset of $\mathcal{G}(\mathcal{M})_{\mathcal{O}}$ of size $n + 1$. If $\mathcal{A} = \mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(i)$ for some $i \in \mathcal{O}$ then the claim holds for $\mathcal{A}$ by assumption. Otherwise, $\mathcal{A} = \bigcup_{i\in\mathcal{A}} \mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(i)$ is a union of ancestral subsets of size $\leq n$. Choose distinct elements $\{i_1, \ldots, i_k\} \subseteq \mathcal{A}$ where $k$ is the smallest integer such that $\bigcup_{j=1}^{k} \mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(i_j) = \mathcal{A}$. By applying Proposition B.4 to $\bigcup_{j=1}^{k-1} \mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(i_j)$ and $\mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(i_k)$, thereby noting that the intersection of these two sets is an ancestral subset of size $\leq n$ and making use of the induction hypothesis, we arrive at the conclusion that $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{A}$. $\qquad\square$

*Appendix C*

PROOF OF PROPOSITION C.2. Let $e \in \mathcal{E}$ and $x_{\mathcal{O}} \in \mathcal{X}_{\mathcal{O}}$. For $x_{\mathcal{L}} \in \mathcal{X}$,

$$x_{\mathcal{L}} = f_{\mathcal{L}}(x, e)$$

$$\iff x_{\mathcal{L}} = B_{\mathcal{L}\mathcal{L}} x_{\mathcal{L}} + B_{\mathcal{L}\mathcal{O}} x_{\mathcal{O}} + \Gamma_{\mathcal{L}\mathcal{J}} e$$

$$\iff A_{\mathcal{L}\mathcal{L}} x_{\mathcal{L}} = B_{\mathcal{L}\mathcal{O}} x_{\mathcal{O}} + \Gamma_{\mathcal{L}\mathcal{J}} e$$

$$\iff \begin{cases} A_{\mathcal{L}\mathcal{L}} A_{\mathcal{L}\mathcal{L}}^{+}(B_{\mathcal{L}\mathcal{O}} x_{\mathcal{O}} + \Gamma_{\mathcal{L}\mathcal{J}} e) = B_{\mathcal{L}\mathcal{O}} x_{\mathcal{O}} + \Gamma_{\mathcal{L}\mathcal{J}} e \\ \exists_{v \in \mathcal{X}_{\mathcal{L}}} : x_{\mathcal{L}} = A_{\mathcal{L}\mathcal{L}}^{+}(B_{\mathcal{L}\mathcal{O}} x_{\mathcal{O}} + \Gamma_{\mathcal{L}\mathcal{J}} e) + [\mathbb{I}_{\mathcal{L}} - A_{\mathcal{L}\mathcal{L}}^{+} A_{\mathcal{L}\mathcal{L}}] v, \end{cases}$$

where the last equivalence follows from [Theorem 2, 54]. $\qquad\square$

PROOF OF PROPOSITION C.3. $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$ if and only if for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $\boldsymbol{x}_{\mathcal{O}} \in \mathcal{X}_{\mathcal{O}}$ the linear system of equations

$$\boldsymbol{x}_{\mathcal{L}} = \boldsymbol{f}_{\mathcal{L}}(\boldsymbol{x}, e)$$

$$\iff \boldsymbol{x}_{\mathcal{L}} = B_{\mathcal{LL}}\boldsymbol{x}_{\mathcal{L}} + B_{\mathcal{LO}}\boldsymbol{x}_{\mathcal{O}} + \Gamma_{\mathcal{LJ}}e$$

$$\iff A_{\mathcal{LL}}\boldsymbol{x}_{\mathcal{L}} = B_{\mathcal{LO}}\boldsymbol{x}_{\mathcal{O}} + \Gamma_{\mathcal{LJ}}e$$

has a unique solution $\boldsymbol{x}_{\mathcal{L}} \in \mathcal{X}_{\mathcal{L}}$. Hence, $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$ if and only if $A_{\mathcal{LL}}$ is invertible. $\square$

PROOF OF PROPOSITION C.4. It suffices to show (1) $\implies$ (2) and (1) $\iff$ (3). We start by showing that (1) $\implies$ (2). Let $\mathcal{V} \subseteq \mathcal{L}$ and denote $\mathcal{U} := \mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{L}}}(\mathcal{V})$, then we need to show that $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{U}$. From Proposition C.3 we know that $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$ if and only if the matrix $A_{\mathcal{LL}} = \mathbb{I}_{\mathcal{L}} - B_{\mathcal{LL}}$ is invertible. The matrix $A_{\mathcal{LL}}$ is invertible if and only if the rows of $A_{\mathcal{LL}}$ are all linearly independent. In particular, the rows of $A_{\mathcal{UL}}$ are all linearly independent. Because $A_{\mathcal{UL}} = [A_{\mathcal{UU}} \, Z_{\mathcal{UL}}]$, where $Z_{\mathcal{UL}}$ is the zero matrix, we know that the rows of $A_{\mathcal{UU}} = \mathbb{I}_{\mathcal{U}} - B_{\mathcal{UU}}$ are also all linearly independent, and hence $A_{\mathcal{UU}}$ is invertible.

Next, we show that (1) $\iff$ (3). Observe that the strongly connected components of $\mathcal{G}(\mathcal{M})_{\mathcal{L}}$ form a partition of the set $\mathcal{L}$ and that the directed mixed graph $\mathcal{G}(\mathcal{M})_{\mathcal{L}}$ and the directed graph $\mathcal{G}^a(\mathcal{M})_{\mathcal{L}}$ have the same strongly connected components. Because, by Lemma A.2, the graph of strongly connected components $\mathcal{G}^{\mathrm{sc}}$ of the directed graph $\mathcal{G}^a(\mathcal{M})_{\mathcal{L}}$ is a DAG, the square matrix $B_{\mathcal{LL}}$ can be permuted to an upper triangular block matrix $\tilde{B}_{\mathcal{LL}}$, where for each diagonal block $\tilde{B}_{\mathcal{VV}}$ of $\tilde{B}_{\mathcal{LL}}$ the set of nodes $\mathcal{V}$ is a strongly connected component in $\mathcal{G}(\mathcal{M})_{\mathcal{L}}$.

Without loss of generality we assume now that $B_{\mathcal{LL}}$ is an upper triangular block matrix. From Proposition C.3 it follows that $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$ if and only if the matrix $A_{\mathcal{LL}} = \mathbb{I}_{\mathcal{L}} - B_{\mathcal{LL}}$ is invertible. Because $B_{\mathcal{LL}}$ is an upper triangular block matrix, we know that $A_{\mathcal{LL}}$ is an upper triangular block matrix, where for each diagonal block $A_{\mathcal{VV}}$ of $A_{\mathcal{LL}}$ the set of nodes $\mathcal{V}$ is a strongly connected component in $\mathcal{G}(\mathcal{M})_{\mathcal{L}}$. Since an upper triangular block matrix $A_{\mathcal{LL}}$ is invertible if and only if every diagonal block in $A_{\mathcal{LL}}$ is invertible, we have that $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$ if and only if $\mathcal{M}$ is uniquely solvable w.r.t. each strongly connected component in $\mathcal{G}(\mathcal{M})_{\mathcal{L}}$. $\square$

PROOF OF PROPOSITION C.5. By the definition of marginalization and Proposition C.3 the marginal causal mechanism $\tilde{\boldsymbol{f}}$ is given by

$$\tilde{\boldsymbol{f}}(\boldsymbol{x}_{\mathcal{O}}, e) := \boldsymbol{f}_{\mathcal{O}}(\boldsymbol{x}_{\mathcal{O}}, \boldsymbol{g}_{\mathcal{L}}(\boldsymbol{x}_{\mathcal{O}}, e), e)$$

$$= B_{\mathcal{OO}}\boldsymbol{x}_{\mathcal{O}} + B_{\mathcal{OL}}\boldsymbol{g}_{\mathcal{L}}(\boldsymbol{x}_{\mathcal{O}}, e) + \Gamma_{\mathcal{OJ}}e$$

$$= [B_{\mathcal{OO}} + B_{\mathcal{OL}}A_{\mathcal{LL}}^{-1}B_{\mathcal{LO}}]\boldsymbol{x}_{\mathcal{O}} + [B_{\mathcal{OL}}A_{\mathcal{LL}}^{-1}\Gamma_{\mathcal{LJ}} + \Gamma_{\mathcal{OJ}}]e.$$

From Proposition C.4 and 5.12 it follows that the marginalization respects the latent projection. $\square$

## E.2. Proofs of the main text

*Section 2*

PROOF OF PROPOSITION 2.12. Let $i \in \mathcal{I}$. Note that Definition 2.7 can alternatively be formulated as follows: for $k \in \mathcal{I} \cup \mathcal{J}$, $k \notin \mathrm{pa}(i)$ if and only if there exists a measurable mapping $\hat{f}_i : \boldsymbol{\mathcal{X}} \times \mathcal{E} \to \mathcal{X}_i$ such that for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$,

$$x_i = f_i(\boldsymbol{x}, e) \iff x_i = \hat{f}_i(\boldsymbol{x}, e)$$

and either $k \in \mathcal{I}$ and there exists $\hat{x}_k \in \mathcal{X}_k$ such that $\hat{f}_i(\boldsymbol{x}, \boldsymbol{e}) = \hat{f}_i(\boldsymbol{x}_{\backslash k}, \hat{x}_k, \boldsymbol{e})$ for all $\boldsymbol{x} \in \mathcal{X}, \boldsymbol{e} \in \mathcal{E}$, or $k \in \mathcal{J}$ and there exists $\hat{e}_k \in \mathcal{E}_k$ such that $\hat{f}_i(\boldsymbol{x}, \boldsymbol{e}) = \hat{f}_i(\boldsymbol{x}, \boldsymbol{e}_{\backslash k}, \hat{e}_k)$ for all $\boldsymbol{x} \in \mathcal{X}, \boldsymbol{e} \in \mathcal{E}$. By repeatedly applying (this formulation of) Definition 2.7 to all $k \notin \mathrm{pa}(i)$, we obtain the existence of a measurable mapping $\tilde{f}_i : \mathcal{X} \times \mathcal{E} \to \mathcal{X}_i$ and $\hat{\boldsymbol{x}}_{\backslash \mathrm{pa}(i)} \in \mathcal{X}_{\backslash \mathrm{pa}(i)}$, $\hat{\boldsymbol{e}}_{\backslash \mathrm{pa}(i)} \in \mathcal{E}_{\backslash \mathrm{pa}(i)}$ such that for $\mathbb{P}_{\mathcal{E}}$-almost every $\boldsymbol{e} \in \mathcal{E}$ and for all $\boldsymbol{x} \in \mathcal{X}$,

$$x_i = f_i(\boldsymbol{x}, \boldsymbol{e}) \iff x_i = \tilde{f}_i(\boldsymbol{x}, \boldsymbol{e}),$$

and for all $\boldsymbol{e} \in \mathcal{E}$ and all $\boldsymbol{x} \in \mathcal{X}$,

$$\tilde{f}_i(\boldsymbol{x}, \boldsymbol{e}) = \tilde{f}_i(\boldsymbol{x}_{\mathrm{pa}(i)}, \hat{\boldsymbol{x}}_{\backslash \mathrm{pa}(i)}, \boldsymbol{e}_{\mathrm{pa}(i)}, \hat{\boldsymbol{e}}_{\backslash \mathrm{pa}(i)}).$$

Define the SCM $\tilde{\mathcal{M}}$ as $\mathcal{M}$ except that its causal mechanism is $\tilde{\boldsymbol{f}}$ instead of $\boldsymbol{f}$. Then $\tilde{\mathcal{M}}$ is structurally minimal and equivalent to $\mathcal{M}$. $\qquad \square$

PROOF OF PROPOSITION 2.15. The $\mathrm{do}(I, \boldsymbol{\xi}_I)$ operation on $\mathcal{M}$ completely removes the functional dependence on $\boldsymbol{x}$ and $\boldsymbol{e}$ from the $f_i$ components for $i \in I$ and hence the corresponding incoming directed and bidirected edges on nodes in $I$ from the (augmented) graph. $\qquad \square$

PROOF OF PROPOSITION 2.16. The first statement follows from Definitions 2.13 and 2.14. For the second statement, note that a perfect intervention can only remove parental relations, and therefore will never introduce a cycle. $\qquad \square$

PROOF OF PROPOSITION 2.20. This follows directly from Definition 2.18 and 2.19. $\qquad \square$

PROOF OF PROPOSITION 2.21. The additional edges introduced by the twin operation cannot lead to a directed cycle involving both copied and original nodes, because there are no edges pointing from copied nodes to original nodes (i.e., of the form $i' \to v$ with $i' \in I'$ and $v \in \mathcal{V}$). Directed cycles involving only original nodes are absent by assumption, and directed cycles involving only copied nodes as well since they would correspond with a directed cycle in the original directed graph. $\qquad \square$

PROOF OF PROPOSITION 2.22. It suffices to proof the property for directed graphs, since the property for SCMs follows directly from Definitions 2.13 and 2.18.

Applying the intervention $\mathrm{do}(I)$ on the graph $\mathcal{G}$ removes all the incoming edges from the nodes in $I$. Now, if we perform the twin operation w.r.t. $\mathcal{I}$ on this graph $\mathrm{do}(I)(\mathcal{G})$, then we copy the same edges as if we had twinned the graph $\mathcal{G}$ w.r.t. $\mathcal{I}$, except those edges that do point to one of the nodes in $I$. Hence, if we apply the intervention $\mathrm{do}(I \cup I')$ on the graph $\mathrm{twin}(\mathcal{I})(\mathcal{G})$, which removes all incoming edges of both $I$ and its copy $I'$, then we clearly obtain the same graph. $\qquad \square$

*Section 3*

PROOF OF THEOREM 3.3. First we define the solution space $\boldsymbol{\mathcal{S}}(\mathcal{M})$ of $\mathcal{M}$ by

$$\boldsymbol{\mathcal{S}}(\mathcal{M}) := \{(\boldsymbol{e}, \boldsymbol{x}) \in \boldsymbol{\mathcal{E}} \times \boldsymbol{\mathcal{X}} : \boldsymbol{x} = \boldsymbol{f}(\boldsymbol{x}, \boldsymbol{e})\}.$$

This is a measurable set, since $\boldsymbol{\mathcal{S}}(\mathcal{M}) = \boldsymbol{h}^{-1}(\Delta)$, where $\boldsymbol{h} : \boldsymbol{\mathcal{E}} \times \boldsymbol{\mathcal{X}} \to \boldsymbol{\mathcal{X}} \times \boldsymbol{\mathcal{X}}$ is the measurable mapping defined by $\boldsymbol{h}(\boldsymbol{e}, \boldsymbol{x}) = (\boldsymbol{x}, \boldsymbol{f}(\boldsymbol{x}, \boldsymbol{e}))$ and $\Delta$ is the set defined by $\{(\boldsymbol{x}, \boldsymbol{x}) : \boldsymbol{x} \in \boldsymbol{\mathcal{X}}\}$, which is measurable since $\boldsymbol{\mathcal{X}}$ is Hausdorff. Note that

$$\boldsymbol{\mathcal{A}} := \boldsymbol{pr}_{\boldsymbol{\mathcal{E}}}(\boldsymbol{\mathcal{S}}(\mathcal{M})) = \{\boldsymbol{e} \in \boldsymbol{\mathcal{E}} : \exists \boldsymbol{x} \in \boldsymbol{\mathcal{X}} \text{ s.t. } \boldsymbol{x} = \boldsymbol{f}(\boldsymbol{x}, \boldsymbol{e})\},$$

is an analytic set because the projection $pr_{\mathcal{E}} : \mathcal{X} \times \mathcal{E} \to \mathcal{E}$ is a measurable mapping between standard measurable spaces (Lemma F.3).

Suppose that (1) holds, that is, $\mathcal{M}$ has a solution. Then there exists a pair of random variables $(E, X) : \Omega \to \mathcal{E} \times \mathcal{X}$ such that $X = f(X, E)$ $\mathbb{P}$-a.s.. Note that

$$\{\omega \in \Omega : X(\omega) = f\big(X(\omega), E(\omega)\big)\} \subseteq \{\omega \in \Omega : \exists x \in \mathcal{X} \text{ s.t. } x = f\big(x, E(\omega)\big)\}$$

$$\subseteq E^{-1}\Big(\{e \in \mathcal{E} : \exists x \in \mathcal{X} \text{ s.t. } x = f(x, e)\}\Big)$$

$$= E^{-1}(\mathcal{A}).$$

By Lemma F.6, $\mathcal{A}$ is $\mathbb{P}^E$-measurable because it is analytic, and we can write $\mathcal{A} = \mathcal{B} \,\dot\cup\, \mathcal{N}$ with $\mathcal{B} \subseteq \mathcal{E}$ measurable and $\mathcal{N}$ a $\mathbb{P}^E$-null set. Hence $E^{-1}(\mathcal{A}) = E^{-1}(\mathcal{B}) \cup E^{-1}(\mathcal{N})$ where $E^{-1}(\mathcal{N})$ is a $\mathbb{P}$-null set. Therefore,

$$E^{-1}(\mathcal{B}) \supseteq \{\omega \in \Omega : X(\omega) = f\big(X(\omega), E(\omega)\big)\} \setminus E^{-1}(\mathcal{N})$$

which implies that $\mathbb{P}(E^{-1}(\mathcal{B})) = 1$. Hence, $\mathcal{E} \setminus \mathcal{A}$ is a $\mathbb{P}_{\mathcal{E}}$-null set. In other words, for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ the structural equations $x = f(x, e)$ have a solution $x \in \mathcal{X}$, i.e., (2) holds.

Suppose that (2) holds. Then $\mathcal{E} \setminus pr_{\mathcal{E}}(\mathcal{S}(\mathcal{M}))$ is a $\mathbb{P}_{\mathcal{E}}$-null set. By application of the measurable selection theorem F.8, there exists a measurable $g : \mathcal{E} \to \mathcal{X}$ such that for $\mathbb{P}_{\mathcal{E}}$-almost all $e \in \mathcal{E}$, $g(e) = f(g(e), e)$. Hence, there exists a measurable mapping $g : \mathcal{E} \to \mathcal{X}$ such that for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$x = g(e) \quad \implies \quad x = f(x, e),$$

which we call property (A). Let $\tilde{f} : \mathcal{E} \times \mathcal{X} \to \mathcal{X}$ be the causal mechanism of a structurally minimal SCM that is equivalent to $\mathcal{M}$ (see Proposition 2.12). In particular, for any $\epsilon_{\setminus \mathrm{pa}(\mathcal{I})} \in \mathcal{E}_{\setminus \mathrm{pa}(\mathcal{I})}$, we have that $\tilde{f}(x, e) = \tilde{f}(x, e_{\mathrm{pa}(\mathcal{I})}, \epsilon_{\setminus \mathrm{pa}(\mathcal{I})})$ for all $x \in \mathcal{X}$ and all $e \in \mathcal{E}$. This means that we may also consider $\tilde{f}$ as a mapping $\tilde{f} : \mathcal{X} \times \mathcal{E}_{\mathrm{pa}(\mathcal{I})} \to \mathcal{X}$. By applying Lemma F.10 to the canonical projection $pr_{\mathcal{E}_{\mathrm{pa}(\mathcal{I})}} : \mathcal{E} \to \mathcal{E}_{\mathrm{pa}(\mathcal{I})}$ and using the equivalence of $f$ and $\tilde{f}$, we obtain that for $\mathbb{P}_{\mathcal{E}_{\mathrm{pa}(\mathcal{I})}}$-almost all $e_{\mathrm{pa}(\mathcal{I})} \in \mathcal{E}_{\mathrm{pa}(\mathcal{I})}$ there exists $x \in \mathcal{X}$ with $x = \tilde{f}(x, e_{\mathrm{pa}(\mathcal{I})})$. By applying the implication (2) $\implies$ (A) to $\mathcal{E}_{\mathrm{pa}(\mathcal{I})}$ and $\tilde{f}$, we conclude the existence of a measurable $g : \mathcal{E}_{\mathrm{pa}(\mathcal{I})} \to \mathcal{X}$ such that for $\mathbb{P}_{\mathcal{E}_{\mathrm{pa}(\mathcal{I})}}$-almost all $e_{\mathrm{pa}(\mathcal{I})} \in \mathcal{E}_{\mathrm{pa}(\mathcal{I})}$, $g(e_{\mathrm{pa}(\mathcal{I})}) = \tilde{f}(g(e_{\mathrm{pa}(\mathcal{I})}), e_{\mathrm{pa}(\mathcal{I})})$. Once more using Lemma F.10, we obtain that for $\mathbb{P}_{\mathcal{E}}$-almost all $e \in \mathcal{E}$, $g(e_{\mathrm{pa}(\mathcal{I})}) = f(g(e_{\mathrm{pa}(\mathcal{I})}), e)$. In other words, (3) holds.

Lastly, suppose that (3) holds, that is there exists a measurable solution function $g : \mathcal{E}_{\mathrm{pa}(\mathcal{I})} \to \mathcal{X}$. Then the measurable mappings $E : \mathcal{E} \to \mathcal{E}$ and $X : \mathcal{E} \to \mathcal{X}$, defined by $E(e) := e$ and $X(e) := g(e_{\mathrm{pa}(\mathcal{I})})$ respectively, define a pair of random variables $(X, E)$ such that $X = f(X, E)$ holds a.s. and hence $(X, E)$ is a solution. Hence (1) holds. $\qquad \square$

PROOF OF PROPOSITION 3.6. Let $\tilde{f} : \mathcal{E} \times \mathcal{X} \to \mathcal{X}$ be the causal mechanism of a structurally minimal SCM $\tilde{\mathcal{M}}$ that is equivalent to $\mathcal{M}$ (see Proposition 2.12). For a subset $\mathcal{O} \subseteq \mathcal{I}$ consider the induced subgraph $\mathcal{G}^a(\mathcal{M})_{\mathcal{O}}$ of the augmented graph $\mathcal{G}^a(\mathcal{M})$ on $\mathcal{O}$. Then the acyclicity of $\mathcal{G}^a(\mathcal{M})$ implies that the induced subgraph $\mathcal{G}^a(\mathcal{M})_{\mathcal{O}}$ is acyclic, and hence there exists a topological ordering on the nodes $\mathcal{O}$. We can substitute the components $\tilde{f}_i$ of the causal mechanism $\tilde{f}$ for $i \in \mathcal{O}$ into each other along this topological ordering. This gives a measurable solution function $g_{\mathcal{O}} : \mathcal{X}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}_{\mathcal{O}}$ for $\tilde{\mathcal{M}}$, and hence for $\mathcal{M}$. It is clear from the acyclic structure that this mapping $g_{\mathcal{O}}$ is independent of the choice of the topological ordering and is the only solution function for $\mathcal{M}$. Therefore, $\tilde{\mathcal{M}}$ is uniquely solvable w.r.t. $\mathcal{O}$, and so is $\mathcal{M}$. $\qquad \square$

PROOF OF PROPOSITION 3.9. This follows immediately from Definition 2.8 and 3.5. □

PROOF OF THEOREM 3.8. Suppose that (1) holds. By Proposition B.1 there exists a measurable solution function $g_{\mathcal{O}} : \mathcal{X}_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}_{\mathcal{O}}$ for $\mathcal{M}$ w.r.t. $\mathcal{O}$. Then for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x_{\setminus\mathcal{O}} \in \mathcal{X}_{\setminus\mathcal{O}}$ we have that $g_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})})$ is a solution of $x_{\mathcal{O}} = f_{\mathcal{O}}(x, e)$. Hence, because of (1), for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x_{\setminus\mathcal{O}} \in \mathcal{X}_{\setminus\mathcal{O}}$ we have that $x_{\mathcal{O}} = f_{\mathcal{O}}(x, e)$ implies $x_{\mathcal{O}} = g_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})})$. Thus, $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{O}$, that is, (2) holds.

Suppose that (2) holds. Let $g_{\mathcal{O}} : \mathcal{X}_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}_{\mathcal{O}}$ be a measurable solution function for $\mathcal{M}$ w.r.t. $\mathcal{O}$. Then, for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$x_{\mathcal{O}} = g_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}) \quad \Longleftrightarrow \quad x_{\mathcal{O}} = f_{\mathcal{O}}(x, e).$$

This implies (1).

For the last statement, assume that $\mathcal{M}$ is uniquely solvable. Let $g : \mathcal{E}_{\mathrm{pa}(\mathcal{I})} \to \mathcal{X}$ be a measurable solution function. Then there exists a measurable set $B \subseteq \mathcal{E}$ with $\mathbb{P}_{\mathcal{E}}(B) = 1$ and for all $e \in B$,

$$\forall x \in \mathcal{X} : x = f(x, e) \implies x = g(e_{\mathrm{pa}(\mathcal{I})}).$$

The existence of a solution for $\mathcal{M}$ follows directly from Theorem 3.3. Each solution $(X, E) : \Omega \to \mathcal{X} \times \mathcal{E}$ of $\mathcal{M}$ satisfies $X(\omega) = f(X(\omega), E(\omega))$ $\mathbb{P}$-a.s.. In addition, it satisfies $E(\omega) \in B$ $\mathbb{P}$-a.s., since $\mathbb{P} \circ E^{-1} = \mathbb{P}_{\mathcal{E}}$. Hence, it satisfies $X(\omega) = g(E(\omega)_{\mathrm{pa}(\mathcal{I})})$ $\mathbb{P}$-a.s.. Thus for every solution $(X, E)$ the associated observational distribution is the push-forward of $\mathbb{P}_{\mathcal{E}}$ under $g \circ pr_{\mathrm{pa}(\mathcal{I})}$. □

PROOF OF PROPOSITION 3.11. Let $g_{\mathcal{O}} : \mathcal{X}_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}_{\mathcal{O}}$ be a measurable solution function for $\mathcal{M}$ w.r.t. $\mathcal{O}$. Then the mapping $\tilde{g}_{\mathcal{O}\cup I} : \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}_{\mathcal{O}\cup I}$ defined by $\tilde{g}_{\mathcal{O}\cup I}(e_{\mathrm{pa}(\mathcal{O})}) := (g_{\mathcal{O}}(\xi_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}), \xi_I)$ is a measurable solution function for the SCM $\mathcal{M}_{\mathrm{do}(I,\xi_I)}$ w.r.t. $\mathcal{O} \cup I$. If $\mathcal{M}$ is (uniquely) solvable w.r.t. $\mathcal{O}$, then it follows that $\mathcal{M}_{\mathrm{do}(I,\xi_I)}$ is (uniquely) solvable w.r.t. $\mathcal{O} \cup I$. □

PROOF OF PROPOSITION 3.13. It suffices to show that solvability of $\mathcal{M}$ w.r.t. $\mathcal{O}$ implies ancestral solvability w.r.t. $\mathcal{O}$. Solvability of $\mathcal{M}$ w.r.t. $\mathcal{O}$ implies that there exists a measurable mapping $g_{\mathcal{O}} : \mathcal{X}_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{O})} \to \mathcal{X}_{\mathcal{O}}$ such that for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$x_{\mathcal{O}} = g_{\mathcal{O}}(x_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}) \quad \Longrightarrow \quad x_{\mathcal{O}} = f_{\mathcal{O}}(x, e).$$

Let $\tilde{f} : \mathcal{E} \times \mathcal{X} \to \mathcal{X}$ be the causal mechanism of a structurally minimal SCM $\tilde{\mathcal{M}}$ that is equivalent to $\mathcal{M}$ (see Proposition 2.12). Let $\mathcal{P} := \mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(\mathcal{A})$ for some $\mathcal{A} \subseteq \mathcal{O}$. Then for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$\begin{cases} x_{\mathcal{P}} &= (g_{\mathcal{O}})_{\mathcal{P}}(x_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}) \\ x_{\mathcal{O}\setminus\mathcal{P}} &= (g_{\mathcal{O}})_{\mathcal{O}\setminus\mathcal{P}}(x_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}) \end{cases} \Longrightarrow \begin{cases} x_{\mathcal{P}} &= \tilde{f}_{\mathcal{P}}(x_{\mathrm{pa}(\mathcal{P})}, e_{\mathrm{pa}(\mathcal{P})}) \\ x_{\mathcal{O}\setminus\mathcal{P}} &= \tilde{f}_{\mathcal{O}\setminus\mathcal{P}}(x_{\mathrm{pa}(\mathcal{O}\setminus\mathcal{P})}, e_{\mathrm{pa}(\mathcal{O}\setminus\mathcal{P})}). \end{cases}$$

Since $\mathrm{pa}(\mathcal{P}) \setminus \mathcal{P} \subseteq \mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}$, we have that in particular for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$x_{\mathcal{P}} = (g_{\mathcal{O}})_{\mathcal{P}}(x_{\mathrm{pa}(\mathcal{O})\setminus\mathcal{O}}, e_{\mathrm{pa}(\mathcal{O})}) \quad \Longrightarrow \quad x_{\mathcal{P}} = \tilde{f}_{\mathcal{P}}(x_{\mathrm{pa}(\mathcal{P})}, e_{\mathrm{pa}(\mathcal{P})}).$$

This implies that the mapping $(g_{\mathcal{O}})_{\mathcal{P}}$ cannot depend on elements different from $\mathrm{pa}(\mathcal{P})$. Moreover, it follows from the definition of $\mathcal{P}$ that $(\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}) \cap \mathrm{pa}(\mathcal{P}) = \mathrm{pa}(\mathcal{P}) \setminus \mathcal{P}$ and

thus we have $\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O} = (\mathrm{pa}(\mathcal{P}) \setminus \mathcal{P}) \cup (\mathrm{pa}(\mathcal{O}) \setminus (\mathcal{O} \cup \mathrm{pa}(\mathcal{P})))$. Now, pick an element $\hat{\boldsymbol{x}}_{\mathrm{pa}(\mathcal{O}) \setminus (\mathcal{O} \cup \mathrm{pa}(\mathcal{P}))} \in \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{O}) \setminus (\mathcal{O} \cup \mathrm{pa}(\mathcal{P}))}$ and define the mapping $\tilde{\boldsymbol{g}}_{\mathcal{P}} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{P}) \setminus \mathcal{P}} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{P})} \to \boldsymbol{\mathcal{X}}_{\mathcal{P}}$ by

$$\tilde{\boldsymbol{g}}_{\mathcal{P}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{P}) \setminus \mathcal{P}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{P})}) := (\boldsymbol{g}_{\mathcal{O}})_{\mathcal{P}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{P}) \setminus \mathcal{P}}, \hat{\boldsymbol{x}}_{\mathrm{pa}(\mathcal{O}) \setminus (\mathcal{O} \cup \mathrm{pa}(\mathcal{P}))}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}).$$

Then, for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$

$$\boldsymbol{x}_{\mathcal{P}} = \tilde{\boldsymbol{g}}_{\mathcal{P}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{P}) \setminus \mathcal{P}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{P})}) \quad \Longleftrightarrow \quad \boldsymbol{x}_{\mathcal{P}} = (\boldsymbol{g}_{\mathcal{O}})_{\mathcal{P}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}).$$

Together this gives that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$

$$\boldsymbol{x}_{\mathcal{P}} = \tilde{\boldsymbol{g}}_{\mathcal{P}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{P}) \setminus \mathcal{P}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{P})}) \quad \Longrightarrow \quad \boldsymbol{x}_{\mathcal{P}} = \tilde{\boldsymbol{f}}_{\mathcal{P}}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{P})}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{P})}).$$

which is equivalent to the statement that $\mathcal{M}$ is solvable w.r.t. $\mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{O}}}(\mathcal{A})$. $\qquad \square$

*Section 4*

LEMMA E.1. *Let $\mathcal{M}$ be an SCM that is uniquely solvable w.r.t. two subsets $A, B \subseteq \mathcal{I}$ that satisfy $A \subseteq B$ and $\mathrm{pa}(A) \setminus A \subseteq \mathrm{pa}(B) \setminus B$. Let $\boldsymbol{g}_A : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(A) \setminus A} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(A)} \to \boldsymbol{\mathcal{X}}_A$ and $\boldsymbol{g}_B : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(B) \setminus B} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(B)} \to \boldsymbol{\mathcal{X}}_B$ be measurable solution functions for $\mathcal{M}$ w.r.t. A and B, respectively. Then for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$*

$$\boldsymbol{g}_A(\boldsymbol{x}_{\mathrm{pa}(A) \setminus A}, \boldsymbol{e}_{\mathrm{pa}(A)}) = (\boldsymbol{g}_B)_A(\boldsymbol{x}_{\mathrm{pa}(B) \setminus B}, \boldsymbol{e}_{\mathrm{pa}(B)}).$$

PROOF. Without loss of generality, we assume that $\mathcal{M}$ is structurally minimal (see Proposition 2.12). Let $\bar{\boldsymbol{\mathcal{E}}} \subseteq \boldsymbol{\mathcal{E}}$ be a measurable set with $\mathbb{P}_{\boldsymbol{\mathcal{E}}}(\bar{\boldsymbol{\mathcal{E}}}) = 1$ such that for all $\boldsymbol{e} \in \bar{\boldsymbol{\mathcal{E}}}$ for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$:

$$\boldsymbol{x}_A = \boldsymbol{g}_A(\boldsymbol{x}_{\mathrm{pa}(A) \setminus A}, \boldsymbol{e}_{\mathrm{pa}(A)}) \iff \boldsymbol{x}_A = \boldsymbol{f}_A(\boldsymbol{x}_{\mathrm{pa}(A)}, \boldsymbol{e}_{\mathrm{pa}(A)})$$

and

$$\boldsymbol{x}_B = \boldsymbol{g}_B(\boldsymbol{x}_{\mathrm{pa}(B) \setminus B}, \boldsymbol{e}_{\mathrm{pa}(B)}) \iff \boldsymbol{x}_B = \boldsymbol{f}_B(\boldsymbol{x}_{\mathrm{pa}(B)}, \boldsymbol{e}_{\mathrm{pa}(B)}).$$

Now let $\boldsymbol{e} \in \bar{\boldsymbol{\mathcal{E}}}$ and let $\boldsymbol{x}_{A \cup \mathrm{pa}(B) \setminus B} \in \boldsymbol{\mathcal{X}}_{A \cup \mathrm{pa}(B) \setminus B}$. Then

$$\boldsymbol{x}_A = (\boldsymbol{g}_B)_A(\boldsymbol{x}_{\mathrm{pa}(B) \setminus B}, \boldsymbol{e}_{\mathrm{pa}(B)})$$

$$\Longrightarrow \begin{cases} \boldsymbol{x}_A = (\boldsymbol{g}_B)_A(\boldsymbol{x}_{\mathrm{pa}(B) \setminus B}, \boldsymbol{e}_{\mathrm{pa}(B)}) \\ \exists \boldsymbol{x}_{B \setminus A} \in \boldsymbol{\mathcal{X}}_{B \setminus A} : \quad \boldsymbol{x}_{B \setminus A} = (\boldsymbol{g}_B)_{B \setminus A}(\boldsymbol{x}_{\mathrm{pa}(B) \setminus B}, \boldsymbol{e}_{\mathrm{pa}(B)}) \end{cases}$$

$$\Longrightarrow \exists \boldsymbol{x}_{B \setminus A} \in \boldsymbol{\mathcal{X}}_{B \setminus A} : \quad \boldsymbol{x}_B = \boldsymbol{g}_B(\boldsymbol{x}_{\mathrm{pa}(B) \setminus B}, \boldsymbol{e}_{\mathrm{pa}(B)})$$

$$\Longrightarrow \exists \boldsymbol{x}_{B \setminus A} \in \boldsymbol{\mathcal{X}}_{B \setminus A} : \quad \boldsymbol{x}_B = \boldsymbol{f}_B(\boldsymbol{x}_{\mathrm{pa}(B)}, \boldsymbol{e}_{\mathrm{pa}(B)})$$

$$\Longrightarrow \exists \boldsymbol{x}_{B \setminus A} \in \boldsymbol{\mathcal{X}}_{B \setminus A} : \quad \boldsymbol{x}_A = \boldsymbol{f}_A(\boldsymbol{x}_{\mathrm{pa}(A)}, \boldsymbol{e}_{\mathrm{pa}(A)})$$

$$\Longrightarrow \boldsymbol{x}_A = \boldsymbol{f}_A(\boldsymbol{x}_{\mathrm{pa}(A)}, \boldsymbol{e}_{\mathrm{pa}(A)})$$

$$\Longrightarrow \boldsymbol{x}_A = \boldsymbol{g}_A(\boldsymbol{x}_{\mathrm{pa}(A) \setminus A}, \boldsymbol{e}_{\mathrm{pa}(A)}),$$

where the exists-quantifier could be omitted because the expression it binds to does not depend on $\boldsymbol{x}_{B \setminus A}$ (from the assumptions it follows that $(A \cup \mathrm{pa}(A)) \cap (B \setminus A) = \emptyset$). Hence, for all $\boldsymbol{e} \in \bar{\boldsymbol{\mathcal{E}}}$ and all $\boldsymbol{x}_{A \cup \mathrm{pa}(B) \setminus B} \in \boldsymbol{\mathcal{X}}_{A \cup \mathrm{pa}(B) \setminus B}$

$$\boldsymbol{x}_A = (\boldsymbol{g}_B)_A(\boldsymbol{x}_{\mathrm{pa}(B) \setminus B}, \boldsymbol{e}_{\mathrm{pa}(B)}) \Longrightarrow \boldsymbol{x}_A = \boldsymbol{g}_A(\boldsymbol{x}_{\mathrm{pa}(A) \setminus A}, \boldsymbol{e}_{\mathrm{pa}(A)}).$$

Hence, for all $e \in \bar{\mathcal{E}}$ and all $x_{A \cup \mathrm{pa}(B) \setminus B} \in \mathcal{X}_{A \cup \mathrm{pa}(B) \setminus B}$

$$(g_B)_A(x_{\mathrm{pa}(B) \setminus B}, e_{\mathrm{pa}(B)}) = g_A(x_{\mathrm{pa}(A) \setminus A}, e_{\mathrm{pa}(A)}).$$

Since this expression does not depend on $x_{(B \setminus A) \cup \mathcal{I} \setminus (B \cup \mathrm{pa}(B))}$, from Lemma F.11.(2) we conclude that for all $e \in \bar{\mathcal{E}}$ and all $x \in \mathcal{X}$

$$(g_B)_A(x_{\mathrm{pa}(B) \setminus B}, e_{\mathrm{pa}(B)}) = g_A(x_{\mathrm{pa}(A) \setminus A}, e_{\mathrm{pa}(A)}).$$

$\square$

LEMMA E.2. *An SCM $\mathcal{M}$ is observationally equivalent to $\mathcal{M}^{\mathrm{twin}}$ w.r.t. $\mathcal{O} \subseteq \mathcal{I}$.*

PROOF. Let $(X, E)$ be a solution of $\mathcal{M}$, then $((X, X), E)$ is a solution of $\mathcal{M}^{\mathrm{twin}}$. Conversely, let $((X, X'), E)$ be a solution of $\mathcal{M}^{\mathrm{twin}}$, then $(X, E)$ is a solution of $\mathcal{M}$. $\square$

PROOF OF PROPOSITION 4.6. First we show that equivalence implies counterfactual equivalence w.r.t. $\mathcal{O}$. The twin operation preserves the equivalence relation on SCMs and since equivalent SCMs are interventionally equivalent w.r.t. every subset, the two equivalent twin SCMs have to be interventionally equivalent w.r.t. $\mathcal{O} \cup \mathcal{O}'$ for every $\mathcal{O} \subseteq \mathcal{I}$ with $\mathcal{O}'$ the copy of $\mathcal{O}$ in $\mathcal{I}'$.

Now, let $\mathcal{M}$ and $\tilde{\mathcal{M}}$ be counterfactually equivalent w.r.t. $\mathcal{O}$. Then $\mathcal{M}^{\mathrm{twin}}$ and $\tilde{\mathcal{M}}^{\mathrm{twin}}$ are interventionally equivalent w.r.t. $\mathcal{O} \cup \mathcal{O}'$. Thus for $I \subseteq \mathcal{O}$, $I' \subseteq \mathcal{O}'$ the copy of $I$ and $\boldsymbol{\xi}_{I'} = \boldsymbol{\xi}_I \in \mathcal{X}_I$, $\mathcal{M}^{\mathrm{twin}}_{\mathrm{do}(I \cup I', \boldsymbol{\xi}_{I \cup I'})}$ and $\tilde{\mathcal{M}}^{\mathrm{twin}}_{\mathrm{do}(I \cup I', \boldsymbol{\xi}_{I \cup I'})}$ are observationally equivalent w.r.t. $\mathcal{O} \cup \mathcal{O}'$. In particular, they are observationally equivalent w.r.t. $\mathcal{O}$. From Proposition 2.22 we have that $\mathcal{M}^{\mathrm{twin}}_{\mathrm{do}(I \cup I', \boldsymbol{\xi}_{I \cup I'})} = (\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)})^{\mathrm{twin}}$ and $\tilde{\mathcal{M}}^{\mathrm{twin}}_{\mathrm{do}(I \cup I', \boldsymbol{\xi}_{I \cup I'})} = (\tilde{\mathcal{M}}_{\mathrm{do}(I, \boldsymbol{\xi}_I)})^{\mathrm{twin}}$, and together with Lemma E.2 this gives that $\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$ and $\tilde{\mathcal{M}}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$ are observationally equivalent w.r.t. $\mathcal{O}$. $\square$

*Section 5*

LEMMA E.3. *Let $\mathcal{M}$ be an SCM. Let $B \subseteq \mathcal{I}$ and $A \subseteq \mathcal{I} \cup \mathcal{J}$ such that $(\mathrm{pa}(B) \setminus B) \subseteq A$ and $B \cap A = \emptyset$. Assume that $g_B : \mathcal{X}_A \times \mathcal{E}_A \to \mathcal{X}_B$ is a measurable function such that for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$*

$$x_B = f_B(x_{\mathrm{pa}(B)}, e_{\mathrm{pa}(B)}) \iff x_B = g_B(x_A, e_A).$$

*Then $\mathcal{M}$ is uniquely solvable w.r.t. $B$.*

PROOF. Assume that for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$x_B = f_B(x_{\mathrm{pa}(B)}, e_{\mathrm{pa}(B)}) \iff x_B = g_B(x_A, e_A).$$

Let $C := A \setminus (\mathrm{pa}(B) \setminus B)$, then by Lemma F.11.(7) we have that there exists $\hat{e}_C \in \mathcal{E}_C$ and $\hat{x}_C \in \mathcal{X}_C$ such that for $\mathbb{P}_{\mathcal{E}_{\mathcal{J} \setminus C}}$-almost every $e_{\mathcal{J} \setminus C} \in \mathcal{E}_{\mathcal{J} \setminus C}$ and for all $x_{\mathcal{I} \setminus C} \in \mathcal{X}_{\mathcal{I} \setminus C}$

$$x_B = f_B(x_{\mathrm{pa}(B)}, e_{\mathrm{pa}(B)}) \iff x_B = g_B(x_{\mathrm{pa}(B) \setminus B}, \hat{x}_C, e_{\mathrm{pa}(B)}, \hat{e}_C).$$

Defining the mapping $h_B : \mathcal{X}_{\mathrm{pa}(B) \setminus B} \times \mathcal{E}_{\mathrm{pa}(B)} \to \mathcal{X}_B$ by

$$h_B(x_{\mathrm{pa}(B) \setminus B}, e_{\mathrm{pa}(B)}) := g_B(x_{\mathrm{pa}(B) \setminus B}, \hat{x}_C, e_{\mathrm{pa}(B)}, \hat{e}_C),$$

where we picked $\hat{e}_C \in \mathcal{E}_C$ and $\hat{x}_C \in \mathcal{X}_C$ such that the above equivalence holds, and applying Lemma F.11.(6) we get that for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$x_B = f_B(x_{\mathrm{pa}(B)}, e_{\mathrm{pa}(B)}) \iff x_B = h_B(x_{\mathrm{pa}(B) \setminus B}, e_{\mathrm{pa}(B)})$$

holds. Thus, $\mathcal{M}$ is uniquely solvable w.r.t. $B$. $\square$

PROOF OF PROPOSITION 5.4. From unique solvability of $\mathcal{M}$ w.r.t. $\mathcal{L}_1$ it follows that there exists a mapping $\boldsymbol{g}_{\mathcal{L}_1} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{L}_1)\setminus(\mathcal{L}_1)} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{L}_1)} \to \boldsymbol{\mathcal{X}}_{\mathcal{L}_1}$ such that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$

$$\boldsymbol{x}_{\mathcal{L}_1} = \boldsymbol{g}_{\mathcal{L}_1}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1)\setminus\mathcal{L}_1}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1)}) \quad \Longleftrightarrow \quad \boldsymbol{x}_{\mathcal{L}_1} = \boldsymbol{f}_{\mathcal{L}_1}(\boldsymbol{x}, \boldsymbol{e}).$$

Let $\widehat{\mathrm{pa}}$ denotes the parents in $\mathcal{G}^a(\mathcal{M}_{\mathrm{marg}(\mathcal{L}_1)})$. Note that $\widehat{\mathrm{pa}}(\mathcal{L}_2) \setminus \mathcal{L}_2 \subseteq \mathrm{pa}(\mathcal{L}_1 \cup \mathcal{L}_2) \setminus (\mathcal{L}_1 \cup \mathcal{L}_2)$. Let $\tilde{\boldsymbol{f}}$ denote the marginal causal mechanism of a structurally minimal SCM that is equivalent to the marginalization $\mathcal{M}_{\mathrm{marg}(\mathcal{L}_1)}$ constructed from $\boldsymbol{g}_{\mathcal{L}_1}$ (see Proposition 2.12).

$\Longrightarrow$ : If $\mathcal{M}_{\mathrm{marg}(\mathcal{L}_1)}$ is uniquely solvable w.r.t. $\mathcal{L}_2$, then there exists a mapping $\tilde{\boldsymbol{g}}_{\mathcal{L}_2}$ : $\boldsymbol{\mathcal{X}}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)\setminus\mathcal{L}_2} \times \boldsymbol{\mathcal{E}}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)} \to \boldsymbol{\mathcal{X}}_{\mathcal{L}_2}$ such that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x}_{\mathcal{I}\setminus\mathcal{L}_1} \in \boldsymbol{\mathcal{X}}_{\mathcal{I}\setminus\mathcal{L}_1}$

$$\boldsymbol{x}_{\mathcal{L}_2} = \tilde{\boldsymbol{g}}_{\mathcal{L}_2}(\boldsymbol{x}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)\setminus\mathcal{L}_2}, \boldsymbol{e}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}) \Longleftrightarrow \boldsymbol{x}_{\mathcal{L}_2} = \boldsymbol{f}_{\mathcal{L}_2}(\boldsymbol{g}_{\mathcal{L}_1}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1)\setminus\mathcal{L}_1}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1)}), \boldsymbol{x}_{\mathcal{I}\setminus\mathcal{L}_1}, \boldsymbol{e}).$$

Define the mapping $\boldsymbol{h} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)\setminus(\mathcal{L}_1\cup\mathcal{L}_2)} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)} \to \boldsymbol{\mathcal{X}}_{\mathcal{L}_1\cup\mathcal{L}_2}$ by

$$(\boldsymbol{h}_{\mathcal{L}_1}, \boldsymbol{h}_{\mathcal{L}_2})(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)\setminus(\mathcal{L}_1\cup\mathcal{L}_2)}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)}) :=$$

$$\left( \boldsymbol{g}_{\mathcal{L}_1}\big((\tilde{\boldsymbol{g}}_{\mathcal{L}_2})_{\mathrm{pa}(\mathcal{L}_1)}(\boldsymbol{x}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)\setminus\mathcal{L}_2}, \boldsymbol{e}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}), \boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1)\setminus(\mathcal{L}_1\cup\mathcal{L}_2)}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1)}\big), \tilde{\boldsymbol{g}}_{\mathcal{L}_2}(\boldsymbol{x}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)\setminus\mathcal{L}_2}, \boldsymbol{e}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}) \right).$$

Then for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$

$$\begin{cases} \boldsymbol{x}_{\mathcal{L}_1} &= \boldsymbol{f}_{\mathcal{L}_1}(\boldsymbol{x}, \boldsymbol{e}) \\ \boldsymbol{x}_{\mathcal{L}_2} &= \boldsymbol{f}_{\mathcal{L}_2}(\boldsymbol{x}, \boldsymbol{e}) \end{cases}$$

$$\Longleftrightarrow \begin{cases} \boldsymbol{x}_{\mathcal{L}_1} &= \boldsymbol{g}_{\mathcal{L}_1}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1)\setminus\mathcal{L}_1}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1)}) \\ \boldsymbol{x}_{\mathcal{L}_2} &= \boldsymbol{f}_{\mathcal{L}_2}(\boldsymbol{x}, \boldsymbol{e}) \end{cases}$$

$$\Longleftrightarrow \begin{cases} \boldsymbol{x}_{\mathcal{L}_1} &= \boldsymbol{g}_{\mathcal{L}_1}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1)\setminus\mathcal{L}_1}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1)}) \\ \boldsymbol{x}_{\mathcal{L}_2} &= \boldsymbol{f}_{\mathcal{L}_2}(\boldsymbol{g}_{\mathcal{L}_1}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1)\setminus\mathcal{L}_1}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1)}), \boldsymbol{x}_{\mathcal{I}\setminus\mathcal{L}_1}, \boldsymbol{e}) \end{cases}$$

$$\Longleftrightarrow \begin{cases} \boldsymbol{x}_{\mathcal{L}_1} &= \boldsymbol{g}_{\mathcal{L}_1}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1)\setminus\mathcal{L}_1}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1)}) \\ \boldsymbol{x}_{\mathcal{L}_2} &= \tilde{\boldsymbol{g}}_{\mathcal{L}_2}(\boldsymbol{x}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)\setminus\mathcal{L}_2}, \boldsymbol{e}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}) \end{cases}$$

$$\Longleftrightarrow \begin{cases} \boldsymbol{x}_{\mathcal{L}_1} &= \boldsymbol{g}_{\mathcal{L}_1}\big((\tilde{\boldsymbol{g}}_{\mathcal{L}_2})_{\mathrm{pa}(\mathcal{L}_1)}(\boldsymbol{x}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)\setminus\mathcal{L}_2}, \boldsymbol{e}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}), \boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1)\setminus(\mathcal{L}_1\cup\mathcal{L}_2)}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1)}\big) \\ \boldsymbol{x}_{\mathcal{L}_2} &= \tilde{\boldsymbol{g}}_{\mathcal{L}_2}(\boldsymbol{x}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)\setminus\mathcal{L}_2}, \boldsymbol{e}_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}) \end{cases}$$

$$\Longleftrightarrow \begin{cases} \boldsymbol{x}_{\mathcal{L}_1} &= \boldsymbol{h}_{\mathcal{L}_1}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)\setminus(\mathcal{L}_1\cup\mathcal{L}_2)}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)}) \\ \boldsymbol{x}_{\mathcal{L}_2} &= \boldsymbol{h}_{\mathcal{L}_2}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)\setminus(\mathcal{L}_1\cup\mathcal{L}_2)}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)}), \end{cases}$$

where in the first equivalence we used unique solvability w.r.t. $\mathcal{L}_1$ of $\mathcal{M}$, in the second we used substitution, in the third we used unique solvability w.r.t. $\mathcal{L}_2$ of $\mathcal{M}_{\mathrm{marg}(\mathcal{L}_1)}$, in the fourth we used again substitution and in the last equivalence we used the definition of $\boldsymbol{h}$. From this we conclude that $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}_1 \cup \mathcal{L}_2$. Hence, by definition it follows that $\mathrm{marg}(\mathcal{L}_2) \circ \mathrm{marg}(\mathcal{L}_1)(\mathcal{M}) = \mathrm{marg}(\mathcal{L}_1 \cup \mathcal{L}_2)(\mathcal{M})$.

$\Longleftarrow$ : If $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}_1 \cup \mathcal{L}_2$, then there exists a mapping $\boldsymbol{h}$ : $\boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)\setminus(\mathcal{L}_1\cup\mathcal{L}_2)} \times \boldsymbol{\mathcal{E}}_{\mathcal{L}_1\cup\mathcal{L}_2} \to \boldsymbol{\mathcal{X}}_{\mathcal{L}_1\cup\mathcal{L}_2}$ such that for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$

$$\boldsymbol{x}_{\mathcal{L}_1\cup\mathcal{L}_2} = \boldsymbol{h}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)\setminus(\mathcal{L}_1\cup\mathcal{L}_2)}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L}_1\cup\mathcal{L}_2)}) \quad \Longleftrightarrow \quad \boldsymbol{x}_{\mathcal{L}_1\cup\mathcal{L}_2} = \boldsymbol{f}_{\mathcal{L}_1\cup\mathcal{L}_2}(\boldsymbol{x}, \boldsymbol{e}).$$

Then, for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ for all $x \in \mathcal{X}$

$$\begin{cases} x_{\mathcal{L}_1} &= h_{\mathcal{L}_1}(x_{\mathrm{pa}(\mathcal{L}_1 \cup \mathcal{L}_2) \setminus (\mathcal{L}_1 \cup \mathcal{L}_2)}, e_{\mathrm{pa}(\mathcal{L}_1 \cup \mathcal{L}_2)}) \\ x_{\mathcal{L}_2} &= h_{\mathcal{L}_2}(x_{\mathrm{pa}(\mathcal{L}_1 \cup \mathcal{L}_2) \setminus (\mathcal{L}_1 \cup \mathcal{L}_2)}, e_{\mathrm{pa}(\mathcal{L}_1 \cup \mathcal{L}_2)}) \end{cases}$$

$$\iff \begin{cases} x_{\mathcal{L}_1} &= f_{\mathcal{L}_1}(x, e) \\ x_{\mathcal{L}_2} &= f_{\mathcal{L}_2}(x, e) \end{cases}$$

$$\iff \begin{cases} x_{\mathcal{L}_1} &= g_{\mathcal{L}_1}(x_{\mathrm{pa}(\mathcal{L}_1) \setminus \mathcal{L}_1}, e_{\mathrm{pa}(\mathcal{L}_1)}) \\ x_{\mathcal{L}_2} &= f_{\mathcal{L}_2}(x, e) \end{cases}$$

$$\iff \begin{cases} x_{\mathcal{L}_1} &= g_{\mathcal{L}_1}(x_{\mathrm{pa}(\mathcal{L}_1) \setminus \mathcal{L}_1}, e_{\mathrm{pa}(\mathcal{L}_1)}) \\ x_{\mathcal{L}_2} &= f_{\mathcal{L}_2}(g_{\mathcal{L}_1}(x_{\mathrm{pa}(\mathcal{L}_1) \setminus \mathcal{L}_1}, e_{\mathrm{pa}(\mathcal{L}_1)}), x_{\mathcal{I} \setminus \mathcal{L}_1}, e) \end{cases}$$

$$\iff \begin{cases} x_{\mathcal{L}_1} &= g_{\mathcal{L}_1}(x_{\mathrm{pa}(\mathcal{L}_1) \setminus \mathcal{L}_1}, e_{\mathrm{pa}(\mathcal{L}_1)}) \\ x_{\mathcal{L}_2} &= \tilde{f}_{\mathcal{L}_2}(x_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}, e_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}). \end{cases}$$

This gives for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ for all $x_{\mathcal{I} \setminus \mathcal{L}_1} \in \mathcal{X}_{\mathcal{I} \setminus \mathcal{L}_1}$

$$x_{\mathcal{L}_2} = h_{\mathcal{L}_2}(x_{\mathrm{pa}(\mathcal{L}_1 \cup \mathcal{L}_2) \setminus (\mathcal{L}_1 \cup \mathcal{L}_2)}, e_{\mathrm{pa}(\mathcal{L}_1 \cup \mathcal{L}_2)})$$

$$\iff x_{\mathcal{L}_2} = \tilde{f}_{\mathcal{L}_2}(x_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}, e_{\widehat{\mathrm{pa}}(\mathcal{L}_2)}).$$

Now apply Lemma E.3 to conclude that $\mathcal{M}_{\mathrm{marg}(\mathcal{L}_1)}$ is uniquely solvable w.r.t. $\mathcal{L}_2$. $\qquad\square$

PROOF OF PROPOSITION 5.5. The commutation relation with the perfect intervention follows straightforwardly from the definitions of perfect intervention and marginalization and the fact that if $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$, then $\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$ is also uniquely solvable w.r.t. $\mathcal{L}$, since the structural equations for the variables $\mathcal{L}$ are the same for $\mathcal{M}$ and $\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$.

The commutation relation with the twin operation follows straightforwardly from the definition of the twin operation and marginalization and the fact that if $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$, then $\mathrm{twin}(\mathcal{M})$ is uniquely solvable w.r.t. $\mathcal{L} \cup \mathcal{L}'$, where $\mathcal{L}'$ is the copy of $\mathcal{L}$ in $\mathcal{I}'$. $\qquad\square$

LEMMA E.4. *Given an SCM $\mathcal{M}$ and a subset $\mathcal{L} \subseteq \mathcal{I}$ such that $\mathcal{M}$ is uniquely solvable w.r.t. $\mathcal{L}$. Then $\mathcal{M}$ and $\mathrm{marg}(\mathcal{L})(\mathcal{M})$ are observationally equivalent w.r.t. $\mathcal{I} \setminus \mathcal{L}$.*

PROOF. Let $\mathcal{O} := \mathcal{I} \setminus \mathcal{L}$. From unique solvability w.r.t. $\mathcal{L}$ it follows that for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$

$$\begin{cases} x_{\mathcal{L}} &= f_{\mathcal{L}}(x, e) \\ x_{\mathcal{O}} &= f_{\mathcal{O}}(x, e) \end{cases}$$

$$\iff \begin{cases} x_{\mathcal{L}} &= g_{\mathcal{L}}(x_{\mathrm{pa}(\mathcal{L}) \setminus \mathcal{L}}, e_{\mathrm{pa}(\mathcal{L})}) \\ x_{\mathcal{O}} &= f_{\mathcal{O}}(g_{\mathcal{L}}(x_{\mathrm{pa}(\mathcal{L}) \setminus \mathcal{L}}, e_{\mathrm{pa}(\mathcal{L})}), x_{\mathcal{O}}, e) \end{cases}$$

$$\iff \begin{cases} x_{\mathcal{L}} &= g_{\mathcal{L}}(x_{\mathrm{pa}(\mathcal{L}) \setminus \mathcal{L}}, e_{\mathrm{pa}(\mathcal{L})}) \\ x_{\mathcal{O}} &= \tilde{f}(x_{\mathcal{O}}, e), \end{cases}$$

where $\tilde{f}$ is the marginal causal mechanism of $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ constructed from a measurable solution function $g_{\mathcal{L}} : \mathcal{X}_{\mathrm{pa}(\mathcal{L}) \setminus \mathcal{L}} \times \mathcal{E}_{\mathrm{pa}(\mathcal{L})} \to \mathcal{X}_{\mathcal{L}}$ for $\mathcal{M}$ w.r.t. $\mathcal{L}$. Hence, a solution $(X, E)$ of $\mathcal{M}$ satisfies $X_{\mathcal{O}} = \tilde{f}(X_{\mathcal{O}}, E)$ a.s.. Conversely, if $(\tilde{X}_{\mathcal{O}}, E)$ is a solution of the marginal SCM $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ then with $\tilde{X}_{\mathcal{L}} := g_{\mathcal{L}}(\tilde{X}_{\mathrm{pa}(\mathcal{L}) \setminus \mathcal{L}}, E_{\mathrm{pa}(\mathcal{L})})$, the random variables $(X, E) := (\tilde{X}_{\mathcal{O}}, \tilde{X}_{\mathcal{L}}, E)$ are a solution of $\mathcal{M}$. $\qquad\square$

PROOF OF THEOREM 5.6. The observational equivalence follows from Lemma E.4. Using both Lemma E.4 and Proposition 5.5 we can prove the interventional equivalence. Observe that from Proposition 5.5 we know that for a subset $I \subseteq \mathcal{I} \setminus \mathcal{L}$ and a value $\boldsymbol{\xi}_I \in \mathcal{X}_I$, $(\mathrm{marg}(\mathcal{L}) \circ \mathrm{do}(I, \boldsymbol{\xi}_I))(\mathcal{M})$ exists. By Lemma E.4 we know that $\mathrm{do}(I, \boldsymbol{\xi}_I)(\mathcal{M})$ and $(\mathrm{marg}(\mathcal{L}) \circ \mathrm{do}(I, \boldsymbol{\xi}_I))(\mathcal{M})$ are observationally equivalent w.r.t. $\mathcal{O}$ and hence by applying again Proposition 5.5, $\mathrm{do}(I, \boldsymbol{\xi}_I)(\mathcal{M})$ and $(\mathrm{do}(I, \boldsymbol{\xi}) \circ \mathrm{marg}(\mathcal{L}))(\mathcal{M})$ are observationally equivalent w.r.t. $\mathcal{O}$. This implies that $\mathcal{M}$ and $\mathrm{marg}(\mathcal{L})(\mathcal{M})$ are interventionally equivalent w.r.t. $\mathcal{O}$. Lastly, we need to show that $\mathrm{twin}(\mathcal{M})$ and $(\mathrm{twin} \circ \mathrm{marg}(\mathcal{L}))(\mathcal{M})$ are interventionally equivalent w.r.t. $(\mathcal{I} \cup \mathcal{I}') \setminus (\mathcal{L} \cup \mathcal{L}')$, where $\mathcal{L}'$ is the copy of $\mathcal{L}$ in $\mathcal{I}'$. From Proposition 5.5 $(\mathrm{twin} \circ \mathrm{marg}(\mathcal{L}))(\mathcal{M})$ is equivalent to $(\mathrm{marg}(\mathcal{L} \cup \mathcal{L}') \circ \mathrm{twin})(\mathcal{M})$ and since we proved that $(\mathrm{marg}(\mathcal{L} \cup \mathcal{L}') \circ \mathrm{twin})(\mathcal{M})$ and $\mathrm{twin}(\mathcal{M})$ are interventionally equivalent w.r.t. $(\mathcal{I} \cup \mathcal{I}') \setminus (\mathcal{L} \cup \mathcal{L}')$ the result follows. $\qquad\square$

PROOF OF PROPOSITION 5.9. A similar proof as for Theorem 1 in [15] works. $\qquad\square$

PROOF OF PROPOSITION 5.10. First we prove the commutation relation of the perfect intervention. Observe that applying the $\mathrm{do}(I)$ operation to the latent projection $\mathrm{marg}(\mathcal{L})(\mathcal{G})$ removes all the incoming edges on the nodes $I$. Such an incoming edge at a node in $I$ in $\mathrm{marg}(\mathcal{L})(\mathcal{G})$ corresponds to a path in $\mathcal{G}$ that points to that node. But since $\mathrm{do}(I)(\mathcal{G})$ is just $\mathcal{G}$ with all the incoming edges on $I$ removed, the graph $(\mathrm{marg}(\mathcal{L}) \circ \mathrm{do}(I))(\mathcal{G})$ also has all the incoming edges on the nodes $I$ removed.

Next, we will prove the commutation relation of the twin operation. We will denote the copy in $\mathcal{I}'$ of any node $i \in \mathcal{I}$ by $i'$, i.e., $\mathcal{I}' = \{i' : i \in \mathcal{I}\}$. The edges in $(\mathrm{twin}(\mathcal{I} \setminus \mathcal{L}) \circ \mathrm{marg}(\mathcal{L}))(\mathcal{G})$ can be partitioned into three cases:

$$\begin{cases} v \to w & v \in \mathcal{J} \cup \mathcal{I} \setminus \mathcal{L}, w \in \mathcal{J} \cup \mathcal{I} \setminus \mathcal{L}, v \to w \in \mathrm{marg}(\mathcal{L})(\mathcal{G}), \\ v \to w' & v \in \mathcal{J}, w \in \mathcal{I} \setminus \mathcal{L}, v \to w \in \mathrm{marg}(\mathcal{L})(\mathcal{G}), \\ v' \to w' & v \in \mathcal{I} \setminus \mathcal{L}, w \in \mathcal{I} \setminus \mathcal{L}, v \to w \in \mathrm{marg}(\mathcal{L})(\mathcal{G}), \end{cases}$$

where $\mathcal{J} := \mathcal{V} \setminus \mathcal{I}$.

Note that in $\mathrm{twin}(\mathcal{I})(\mathcal{G})$, there are no directed edges of the form $v' \to w$ by definition. Therefore, the edges in $(\mathrm{marg}(\mathcal{L} \cup \mathcal{L}') \circ \mathrm{twin}(\mathcal{I}))(\mathcal{G})$ can be partitioned into three cases:

$$\begin{cases} v \to w & v \in \mathcal{J} \cup \mathcal{I} \setminus \mathcal{L}, w \in \mathcal{J} \cup \mathcal{I} \setminus \mathcal{L}, v \to \ell_1 \to \cdots \to \ell_n \to w \in \mathrm{twin}(\mathcal{I})(\mathcal{G}), \\ v \to w' & v \in \mathcal{J}, w \in \mathcal{I} \setminus \mathcal{L}, v \to \ell_1' \to \cdots \to \ell_n' \to w' \in \mathrm{twin}(\mathcal{I})(\mathcal{G}), \\ v' \to w' & v \in \mathcal{I} \setminus \mathcal{L}, w \in \mathcal{I} \setminus \mathcal{L}, v' \to \ell_1' \to \cdots \to \ell_n' \to w' \in \mathrm{twin}(\mathcal{I})(\mathcal{G}), \end{cases}$$

where all $\ell_1, \ldots, \ell_n \in \mathcal{L}$ and $\ell_1', \ldots, \ell_n' \in \mathcal{L}'$. Thus, the non-endpoint nodes on the directed paths in $\mathrm{twin}(\mathcal{I})(\mathcal{G})$ must either all lie in $\mathcal{L}$ or in $\mathcal{L}'$. With the definition of $\mathrm{twin}(\mathcal{I})(\mathcal{G})$ we can rewrite this as follows:

$$\begin{cases} v \to w & v \in \mathcal{J} \cup \mathcal{I} \setminus \mathcal{L}, w \in \mathcal{J} \cup \mathcal{I} \setminus \mathcal{L}, v \to \ell_1 \to \cdots \to \ell_n \to w \in \mathcal{G}, \\ v \to w' & v \in \mathcal{J}, w \in \mathcal{I} \setminus \mathcal{L}, v \to \ell_1 \to \cdots \to \ell_n \to w \in \mathcal{G}, \\ v' \to w' & v \in \mathcal{I} \setminus \mathcal{L}, w \in \mathcal{I} \setminus \mathcal{L}, v \to \ell_1 \to \cdots \to \ell_n \to w \in \mathcal{G}, \end{cases}$$

where all intermediate $\ell_1, \ldots, \ell_n$ must lie in $\mathcal{L}$. This corresponds exactly with the edges in $(\mathrm{twin}(\mathcal{I} \setminus \mathcal{L}) \circ \mathrm{marg}(\mathcal{L}))(\mathcal{G})$.

$\qquad\square$

PROOF OF PROPOSITION 5.12. Without loss of generality, we assume that $\mathcal{M}$ is structurally minimal (see Proposition 2.12). Let $\boldsymbol{g}_{\mathcal{L}}$ be a measurable solution function for $\mathcal{M}$ w.r.t.

$\mathcal{L}$ and denote by $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ the marginal SCM constructed from $\boldsymbol{g}_{\mathcal{L}}$. For $j \in \mathcal{I} \setminus \mathcal{L}$, define $A_j := \mathrm{an}_{\mathcal{G}(\mathcal{M})_{\mathcal{L}}}(\mathrm{pa}(j) \cap \mathcal{L}) \subseteq \mathcal{L}$ and let $\tilde{\boldsymbol{g}}_{A_j}$ be a measurable solution function for $\mathcal{M}$ w.r.t. $A_j$. Because $A_j \subseteq \mathcal{L}$ and $\mathrm{pa}(A_j) \setminus A_j \subseteq \mathrm{pa}(\mathcal{L}) \setminus \mathcal{L}$, by Lemma E.1, for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $\boldsymbol{x} \in \mathcal{X}$

$$(\boldsymbol{g}_{\mathcal{L}})_{A_j}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}) \setminus \mathcal{L}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L})}) = \tilde{\boldsymbol{g}}_{A_j}(\boldsymbol{x}_{\mathrm{pa}(A_j) \setminus A_j}, \boldsymbol{e}_{\mathrm{pa}(A_j)}) \, .$$

Therefore, the component $\tilde{f}_j$ of the marginal causal mechanism $\tilde{\boldsymbol{f}}$ of $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ satisfies for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $\boldsymbol{x} \in \mathcal{X}$

$$\tilde{f}_j(\boldsymbol{x}_{\mathcal{I} \setminus \mathcal{L}}, \boldsymbol{e}) := f_j\big((\boldsymbol{g}_{\mathcal{L}})_{\mathrm{pa}(j)}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{L}) \setminus \mathcal{L}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{L})}), \boldsymbol{x}_{\mathrm{pa}(j) \setminus \mathcal{L}}, \boldsymbol{e}_{\mathrm{pa}(j)}\big)$$
$$= f_j\big((\tilde{\boldsymbol{g}}_{A_j})_{\mathrm{pa}(j) \cap \mathcal{L}}(\boldsymbol{x}_{\mathrm{pa}(A_j) \setminus A_j}, \boldsymbol{e}_{\mathrm{pa}(A_j)}), \boldsymbol{x}_{\mathrm{pa}(j) \setminus \mathcal{L}}, \boldsymbol{e}_{\mathrm{pa}(j)}\big) \, .$$

Hence, the endogenous parents of $j$ in $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ are a subset of $\big((\mathrm{pa}(A_j) \setminus A_j) \cup (\mathrm{pa}(j) \setminus \mathcal{L})\big) \cap \mathcal{I}$ and the exogenous parents of $j$ in $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ are a subset of $(\mathrm{pa}(A_j) \cup \mathrm{pa}(j)) \cap \mathcal{J}$. Hence, all parents of $j$ in $\mathcal{M}_{\mathrm{marg}(\mathcal{L})}$ are a subset of those $k \in (\mathcal{I} \setminus \mathcal{L}) \cup \mathcal{J}$ such that there exists a path $k \to \ell_1 \to \cdots \to \ell_n \to j \in \mathcal{G}^a(\mathcal{M})$ for $n \geq 0$ and $\ell_1, \ldots, \ell_n \in \mathcal{L}$. Therefore, the augmented graph $\mathcal{G}^a\big(\mathrm{marg}(\mathcal{L})(\mathcal{M})\big)$ is a subgraph of the latent projection $\mathrm{marg}(\mathcal{L})\big(\mathcal{G}^a(\mathcal{M})\big)$. Hence,

$$\mathcal{G}\big(\mathrm{marg}(\mathcal{L})(\mathcal{M})\big) = \mathrm{marg}(\mathcal{J})\Big(\mathcal{G}^a\big(\mathrm{marg}(\mathcal{L})(\mathcal{M})\big)\Big)$$
$$\subseteq \mathrm{marg}(\mathcal{J})\Big(\mathrm{marg}(\mathcal{L})\big(\mathcal{G}^a(\mathcal{M})\big)\Big)$$
$$= \mathrm{marg}(\mathcal{L})\Big(\mathrm{marg}(\mathcal{J})\big(\mathcal{G}^a(\mathcal{M})\big)\Big)$$
$$= \mathrm{marg}(\mathcal{L})\big(\mathcal{G}(\mathcal{M})\big)$$

and we conclude that also the graph $\mathcal{G}\big(\mathrm{marg}(\mathcal{L})(\mathcal{M})\big)$ is a subgraph of the latent projection $\mathrm{marg}(\mathcal{L})\big(\mathcal{G}(\mathcal{M})\big)$. $\qquad\square$

*Section 6*

PROOF OF THEOREM 6.3. This follows directly from Theorem A.7 and A.21. $\qquad\square$

*Section 7*

PROOF OF PROPOSITION 7.1. We define $\tilde{\mathcal{M}} := \mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$, $\widetilde{\mathrm{pa}} := \mathrm{pa}_{\mathcal{G}^a(\tilde{\mathcal{M}})}$ and $\mathcal{A} := \mathrm{an}_{\mathcal{G}(\tilde{\mathcal{M}})_{\setminus i}}(j)$. Suppose that $i \to j \notin \mathrm{marg}(\mathcal{I} \setminus \mathcal{O})(\mathcal{G}(\mathcal{M}))$ and assume that the two induced distributions do not coincide. Because $i \to j \notin \mathrm{marg}(\mathcal{I} \setminus \mathcal{O})(\mathcal{G}(\mathcal{M}))$ it follows that $(\widetilde{\mathrm{pa}}(\mathcal{A}) \setminus \mathcal{A}) \cap \mathcal{I} = \emptyset$. Let now $\tilde{\boldsymbol{g}}_{\mathcal{A}} : \mathcal{E}_{\widetilde{\mathrm{pa}}(\mathcal{A})} \to \mathcal{X}_{\mathcal{A}}$ be a measurable solution function for $\tilde{\mathcal{M}}$ w.r.t. $\mathcal{A}$, i.e., we have for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $\boldsymbol{x} \in \mathcal{X}$

$$\boldsymbol{x}_{\mathcal{A}} = \tilde{\boldsymbol{f}}_{\mathcal{A}}(\boldsymbol{x}, \boldsymbol{e}) \quad \Longleftrightarrow \quad \boldsymbol{x}_{\mathcal{A}} = \tilde{\boldsymbol{g}}_{\mathcal{A}}(\boldsymbol{e}_{\widetilde{\mathrm{pa}}(\mathcal{A})}),$$

where $\tilde{\boldsymbol{f}}$ is the ausal mechanism of $\tilde{\mathcal{M}}$. Because $i \notin \mathcal{A}$ and $j \in \mathcal{A}$, it follows that for the intervened model $(\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)})_{\mathrm{do}(\{i\}, \xi_i)}$ the marginal solution $X_j$ is also a marginal solution of $(\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)})_{\mathrm{do}(\{i\}, \tilde{\xi}_i)}$ and vice versa, which is in contradiction with the assumption. $\qquad\square$

PROOF OF PROPOSITION 7.3. Let's define $\tilde{\mathcal{M}} := \mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$, $\widetilde{\mathrm{pa}} := \mathrm{pa}_{\mathcal{G}^a(\tilde{\mathcal{M}})}$, $\mathcal{A}_i := \mathrm{an}_{\mathcal{G}(\tilde{\mathcal{M}})}(i)$ and $\mathcal{A}_j^{\setminus i} := \mathrm{an}_{\mathcal{G}(\tilde{\mathcal{M}})_{\setminus i}}(j)$. Suppose that there does not exist a bidirected edge $i \leftrightarrow j$

in the latent projection $\mathrm{marg}(\mathcal{I} \setminus \mathcal{O})(\mathcal{G}(\mathcal{M}))$. Because $i \leftrightarrow j \notin \mathrm{marg}(\mathcal{I} \setminus \mathcal{O})(\mathcal{G}(\tilde{\mathcal{M}}))$, where here $\tilde{\mathcal{M}}$ is the intervened model $\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}$, we have that $\mathrm{an}_{\mathcal{G}^a(\tilde{\mathcal{M}})_{\setminus j}}(i) \cap \mathrm{an}_{\mathcal{G}^a(\tilde{\mathcal{M}})_{\setminus i}}(j) \cap \mathcal{J} = \emptyset$. From $j \notin \mathrm{an}_{\mathcal{G}(\tilde{\mathcal{M}})}(i)$ it follows that $\mathrm{an}_{\mathcal{G}(\tilde{\mathcal{M}})_{\setminus j}}(i) = \mathrm{an}_{\mathcal{G}(\tilde{\mathcal{M}})}(i)$, and hence $\mathrm{an}_{\mathcal{G}^a(\tilde{\mathcal{M}})}(i) \cap \mathrm{an}_{\mathcal{G}^a(\tilde{\mathcal{M}})_{\setminus i}}(j) \cap \mathcal{J} = \emptyset$. Observe that $\widetilde{\mathrm{pa}}(\mathcal{A}_i) \subseteq \mathrm{an}_{\mathcal{G}^a(\tilde{\mathcal{M}})}(i)$ and $\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i}) \subseteq \mathrm{an}_{\mathcal{G}^a(\tilde{\mathcal{M}})_{\setminus i}}(j) \cup \{i\}$, and thus $\widetilde{\mathrm{pa}}(\mathcal{A}_i) \cap \widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i}) \cap \mathcal{J} = \emptyset$. Let $\boldsymbol{g}_{\mathcal{A}_i} : \boldsymbol{\mathcal{E}}_{\widetilde{\mathrm{pa}}(\mathcal{A}_i)} \to \boldsymbol{\mathcal{X}}_{\mathcal{A}_i}$ be a measurable solution function for $\tilde{\mathcal{M}}$ w.r.t. $\mathcal{A}_i$, i.e., we have for $\mathbb{P}_{\mathcal{E}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$

$$\boldsymbol{x}_{\mathcal{A}_i} = \tilde{\boldsymbol{f}}_{\mathcal{A}_i}(\boldsymbol{x}, \boldsymbol{e}) \quad \Longleftrightarrow \quad \boldsymbol{x}_{\mathcal{A}_i} = \boldsymbol{g}_{\mathcal{A}_i}(\boldsymbol{e}_{\widetilde{\mathrm{pa}}(\mathcal{A}_i)}),$$

where $\tilde{\boldsymbol{f}}$ is the intervened causal mechanism of $\tilde{\mathcal{M}}$. Because $\widetilde{\mathrm{pa}}(\mathcal{A}_i) \cap \widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i}) \cap \mathcal{J} = \emptyset$ and $i \in \mathcal{A}_i$, we have that $X_i \perp\!\!\!\perp \boldsymbol{E}_{\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i})}$ for every solution $(\boldsymbol{X}, \boldsymbol{E})$ of $\tilde{\mathcal{M}}$.

Assume for the moment that $i \in \widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i}) \setminus \mathcal{A}_j^{\setminus i}$, then $(\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i}) \setminus \mathcal{A}_j^{\setminus i}) \cap \mathcal{I} = \{i\}$. Let $\boldsymbol{g}_{\mathcal{A}_j^{\setminus i}} : \mathcal{X}_i \times \boldsymbol{\mathcal{E}}_{\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i})} \to \boldsymbol{\mathcal{X}}_{\mathcal{A}_j^{\setminus i}}$ be a measurable solution function for $\tilde{\mathcal{M}}$ w.r.t. $\mathcal{A}_j^{\setminus i}$, i.e., we have for $\mathbb{P}_{\mathcal{E}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$ and for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$

$$\boldsymbol{x}_{\mathcal{A}_j^{\setminus i}} = \tilde{\boldsymbol{f}}_{\mathcal{A}_j^{\setminus i}}(\boldsymbol{x}, \boldsymbol{e}) \Longleftrightarrow \boldsymbol{x}_{\mathcal{A}_j^{\setminus i}} = \boldsymbol{g}_{\mathcal{A}_j^{\setminus i}}(x_i, \boldsymbol{e}_{\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i})}).$$

For every measurable set $\mathcal{B}_j \subseteq \mathcal{X}_j$ there exists a version of the regular conditional probability $\mathbb{P}_{\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}}(X_j \in \mathcal{B} \mid X_i = \xi_i)$ such that for every value $\xi_i \in \mathcal{X}_i$ it satisfies

$$\begin{aligned}
\mathbb{P}_{\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}}(X_j \in \mathcal{B}_j \mid X_i = \xi_i) &= \mathbb{P}_{\tilde{\mathcal{M}}}(X_j \in \mathcal{B}_j \mid X_i = \xi_i) \\
&= \mathbb{P}_{\tilde{\mathcal{M}}}\big((\boldsymbol{g}_{\mathcal{A}_j^{\setminus i}})_j(X_i, \boldsymbol{E}_{\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i})}) \in \mathcal{B}_j \mid X_i = \xi_i\big) \\
&= \mathbb{P}_{\tilde{\mathcal{M}}}\big((\boldsymbol{g}_{\mathcal{A}_j^{\setminus i}})_j(\xi_i, \boldsymbol{E}_{\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i})}) \in \mathcal{B}_j \mid X_i = \xi_i\big) \\
&= \mathbb{P}_{\tilde{\mathcal{M}}}\big((\boldsymbol{g}_{\mathcal{A}_j^{\setminus i}})_j(\xi_i, \boldsymbol{E}_{\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i})}) \in \mathcal{B}_j\big) \\
&= \mathbb{P}_{\tilde{\mathcal{M}}_{\mathrm{do}(\{i\}, \xi_i)}}\big((\boldsymbol{g}_{\mathcal{A}_j^{\setminus i}})_j(X_i, \boldsymbol{E}_{\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i})}) \in \mathcal{B}_j\big) \\
&= \mathbb{P}_{\tilde{\mathcal{M}}_{\mathrm{do}(\{i\}, \xi_i)}}\big(X_j \in \mathcal{B}_j\big) \\
&= \mathbb{P}_{\left(\mathcal{M}_{\mathrm{do}(I, \boldsymbol{\xi}_I)}\right)_{\mathrm{do}(\{i\}, \xi_i)}}\big(X_j \in \mathcal{B}_j\big),
\end{aligned}$$

where we used $X_i \perp\!\!\!\perp \boldsymbol{E}_{\widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i})}$ in the fourth equality.

If we assume $i \notin \widetilde{\mathrm{pa}}(\mathcal{A}_j^{\setminus i}) \setminus \mathcal{A}_j^{\setminus i}$ instead of $i \in \mathrm{pa}(\mathcal{A}_j^{\setminus i}) \setminus \mathcal{A}_j^{\setminus i}$, then we similarly arrive at the same conclusion. $\qquad \square$

*Section 8*

PROOF OF PROPOSITION 8.2. We first show that the class of simple SCMs is closed under marginalization. Take two disjoint subsets $\mathcal{L}_1$ and $\mathcal{L}_2$ in $\mathcal{I}$. Then, it suffices to show that $\mathcal{M}_{\mathrm{marg}(\mathcal{L}_1)}$ is uniquely solvable w.r.t. $\mathcal{L}_2$. This follows directly from Proposition 5.4.

To show that the class of simple SCMs is closed under perfect intervention. Let $\mathcal{M}$ be a simple SCM, $\mathcal{O} \subseteq \mathcal{I}$, $I \subseteq \mathcal{I}$ and $\boldsymbol{\xi}_I \in \boldsymbol{\mathcal{X}}_I$. Define $\mathcal{O}_1 := \mathcal{O} \cap I$ and $\mathcal{O}_2 := \mathcal{O} \setminus I$, then $\mathcal{O} = \mathcal{O}_1 \cup \mathcal{O}_2$. Note that $\mathrm{pa}(\mathcal{O}_2) \setminus \mathcal{O}_2 = (\mathrm{pa}(\mathcal{O}_2) \setminus (\mathcal{O}_2 \cup I)) \cup (\mathrm{pa}(\mathcal{O}_2) \cap I)$ and $\mathrm{pa}(\mathcal{O}_2) \setminus (\mathcal{O}_2 \cup I) \subseteq \mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}$. Let $\boldsymbol{g}_{\mathcal{O}_2} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{O}_2) \setminus \mathcal{O}_2} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{O}_2)} \to \boldsymbol{\mathcal{X}}_{\mathcal{O}_2}$ be a measurable solution function for $\mathcal{M}$ w.r.t. $\mathcal{O}_2$. The mapping $\tilde{\boldsymbol{g}}_{\mathcal{O}} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{O})} \to \boldsymbol{\mathcal{X}}_{\mathcal{O}}$ defined by

$$\begin{cases} (\tilde{\boldsymbol{g}}_{\mathcal{O}})_{\mathcal{O}_1}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}) := \boldsymbol{\xi}_{\mathcal{O}_1} \\ (\tilde{\boldsymbol{g}}_{\mathcal{O}})_{\mathcal{O}_2}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}) \setminus \mathcal{O}}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O})}) := \boldsymbol{g}_{\mathcal{O}_2}(\boldsymbol{x}_{\mathrm{pa}(\mathcal{O}_2) \setminus (\mathcal{O}_2 \cup I)}, \boldsymbol{\xi}_{\mathrm{pa}(\mathcal{O}_2) \cap I}, \boldsymbol{e}_{\mathrm{pa}(\mathcal{O}_2)}) \end{cases}$$

is a measurable solution function for $\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}$ w.r.t. $\mathcal{O}$, and it is clear that $\mathcal{M}_{\mathrm{do}(I,\boldsymbol{\xi}_I)}$ is uniquely solvable w.r.t. $\mathcal{O}$.

Next, we show that the class of simple SCMs is closed under the twin operation. Let $\tilde{\mathcal{O}} \subseteq \mathcal{I} \cup \mathcal{I}'$. Take $\mathcal{O}_1 = \tilde{\mathcal{O}} \cap \mathcal{I}$, $\mathcal{O}_2' = \tilde{\mathcal{O}} \cap \mathcal{I}'$ and $\mathcal{O}_2$ the original copy of $\mathcal{O}_2'$ in $\mathcal{I}$. Let $\boldsymbol{g}_{\mathcal{O}_1} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{O}_1)\setminus\mathcal{O}_1} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{O}_1)} \to \boldsymbol{\mathcal{X}}_{\mathcal{O}_1}$ and $\boldsymbol{g}_{\mathcal{O}_2} : \boldsymbol{\mathcal{X}}_{\mathrm{pa}(\mathcal{O}_2)\setminus\mathcal{O}_2} \times \boldsymbol{\mathcal{E}}_{\mathrm{pa}(\mathcal{O}_2)} \to \boldsymbol{\mathcal{X}}_{\mathcal{O}_2}$ be measurable solution functions for $\mathcal{M}$ w.r.t. $\mathcal{O}_1$ and $\mathcal{O}_2$ respectively. Define now the mapping $\boldsymbol{h}_{\tilde{\mathcal{O}}} : \boldsymbol{\mathcal{X}}_{\widetilde{\mathrm{pa}}(\tilde{\mathcal{O}})\setminus\tilde{\mathcal{O}}} \times \boldsymbol{\mathcal{E}}_{\widetilde{\mathrm{pa}}(\tilde{\mathcal{O}})} \to \boldsymbol{\mathcal{X}}_{\tilde{\mathcal{O}}}$ by

$$(\boldsymbol{h}_{\tilde{\mathcal{O}}})_{\tilde{\mathcal{O}}\cap\mathcal{I}}(\boldsymbol{x}_{\widetilde{\mathrm{pa}}(\tilde{\mathcal{O}})\setminus\tilde{\mathcal{O}}}, \boldsymbol{e}_{\widetilde{\mathrm{pa}}(\tilde{\mathcal{O}})}) := \boldsymbol{g}_{\mathcal{O}_1}(\boldsymbol{x}_{\widetilde{\mathrm{pa}}(\mathcal{O}_1)\setminus\mathcal{O}_1}, \boldsymbol{e}_{\widetilde{\mathrm{pa}}(\mathcal{O}_1)})$$

$$(\boldsymbol{h}_{\tilde{\mathcal{O}}})_{\tilde{\mathcal{O}}\cap\mathcal{I}'}(\boldsymbol{x}_{\widetilde{\mathrm{pa}}(\tilde{\mathcal{O}})\setminus\tilde{\mathcal{O}}}, \boldsymbol{e}_{\widetilde{\mathrm{pa}}(\tilde{\mathcal{O}})}) := \boldsymbol{g}_{\mathcal{O}_2}(\boldsymbol{x}_{\widetilde{\mathrm{pa}}(\mathcal{O}_2')\setminus\mathcal{O}_2'}, \boldsymbol{e}_{\widetilde{\mathrm{pa}}(\mathcal{O}_2')}),$$

where we define $\widetilde{\mathrm{pa}} := \mathrm{pa}_{\mathcal{G}^a(\mathcal{M}^{\mathrm{twin}})}$ as the parents w.r.t. the twin graph $\mathcal{G}^a(\mathcal{M}^{\mathrm{twin}})$. Then by construction this mapping $\boldsymbol{h}_{\tilde{\mathcal{O}}}$ is a measurable solution function for $\mathcal{M}^{\mathrm{twin}}$ w.r.t. $\tilde{\mathcal{O}}$, and it is clear that $\mathcal{M}^{\mathrm{twin}}$ is uniquely solvable w.r.t. $\tilde{\mathcal{O}}$.

Lastly, it follows that the observational and all the intervened models of $\mathcal{M}$ and $\mathcal{M}^{\mathrm{twin}}$ are uniquely solvable. From Theorem 3.8 we conclude that $\mathcal{M}$ induces unique observational, interventional and counterfactual distributions. $\qquad\square$

PROOF OF COROLLARY 8.3. This follows from Corollary A.22. $\qquad\square$

## APPENDIX F: MEASURABLE SELECTION THEOREMS

In this appendix we derive some lemmas and state two measurable selection theorems that are used in several proofs in Appendix E. First we introduce the measure theoretic notation and terminology needed to understand the results (see [30] for more details).

DEFINITION F.1 (Standard measurable space). *A measurable space $(\boldsymbol{\mathcal{X}}, \boldsymbol{\Sigma})$ is a* standard measurable space *if it is isomorphic to $(\boldsymbol{\mathcal{Y}}, \mathcal{B}(\boldsymbol{\mathcal{Y}}))$, where $\boldsymbol{\mathcal{Y}}$ is a Polish space, i.e., a separable completely metrizable space,[22] and $\mathcal{B}(\boldsymbol{\mathcal{Y}})$ are the Borel subsets of $\boldsymbol{\mathcal{Y}}$, i.e., the $\sigma$-algebra generated by the open sets in $\boldsymbol{\mathcal{Y}}$. A measure space $(\boldsymbol{\mathcal{X}}, \boldsymbol{\Sigma}, \boldsymbol{\mu})$ is a* standard probability space *if $(\boldsymbol{\mathcal{X}}, \boldsymbol{\Sigma})$ is a standard measurable space and $\boldsymbol{\mu}$ is a probability measure.*

Examples of standard measurable spaces are the open and closed subsets of $\mathbb{R}^d$, and the finite sets with the usual complete metric. If we say that $\boldsymbol{\mathcal{X}}$ is a standard measurable space, then we implicitly assume that there exists a $\sigma$-algebra $\boldsymbol{\Sigma}$ such that $(\boldsymbol{\mathcal{X}}, \boldsymbol{\Sigma})$ is a standard measurable space. Similarly, if we say that $\boldsymbol{\mathcal{X}}$ is a standard probability space with probability measure $\mathbb{P}_{\boldsymbol{\mathcal{X}}}$, then we implicitly assume that there exists a $\sigma$-algebra $\boldsymbol{\Sigma}$ such that $(\boldsymbol{\mathcal{X}}, \boldsymbol{\Sigma}, \mathbb{P}_{\boldsymbol{\mathcal{X}}})$ is a standard probability space.

DEFINITION F.2 (Analytic set). *Let $\boldsymbol{\mathcal{X}}$ be a Polish space. A set $\boldsymbol{\mathcal{A}} \subseteq \boldsymbol{\mathcal{X}}$ is called* analytic *if there exist a Polish space $\boldsymbol{\mathcal{Y}}$ and a continuous mapping $\boldsymbol{f} : \boldsymbol{\mathcal{Y}} \to \boldsymbol{\mathcal{X}}$ with $\boldsymbol{f}(\boldsymbol{\mathcal{Y}}) = \boldsymbol{\mathcal{A}}$.*

---

[22] A *metrizable space* is a topological space $\boldsymbol{\mathcal{X}}$ for which there exists a metric $d$ such that $(\boldsymbol{\mathcal{X}}, d)$ is a metric space and induces the topology on $\boldsymbol{\mathcal{X}}$. For a metric space $(\boldsymbol{\mathcal{X}}, d)$, a *Cauchy sequence* is a sequence $(x_n)_{n\in\mathbb{N}}$ of elements of $\boldsymbol{\mathcal{X}}$ such that for every $\epsilon > 0$ there exists an $N \in \mathbb{N}$ such that for all natural numbers $p, q > N$ we have $d(x_n, x_m) < \epsilon$. We call $(\boldsymbol{\mathcal{X}}, d)$ *complete* if every Cauchy sequence has a limit in $\boldsymbol{\mathcal{X}}$. A *completely metrizable space* is a topological space $\boldsymbol{\mathcal{X}}$ for which there exists a metric $d$ such that $(\boldsymbol{\mathcal{X}}, d)$ is a complete metric space that induces the topology on $\boldsymbol{\mathcal{X}}$. A topological space $\boldsymbol{\mathcal{X}}$ is called *separable* if it contains a countable dense subset, i.e., there exists a sequence $(x_n)_{n\in\mathbb{N}}$ of elements in $\boldsymbol{\mathcal{X}}$ such that every non-empty open subset of $\boldsymbol{\mathcal{X}}$ contains at least one element of the sequence. A separable completely metrizable space is called a *Polish space* (see [9] and [30] for more details).

LEMMA F.3.  *Let $\mathcal{X}$ and $\mathcal{Y}$ be standard measurable spaces and $f : \mathcal{X} \to \mathcal{Y}$ a measurable mapping. Then*

1. *every measurable set $\mathcal{A} \subseteq \mathcal{X}$ is analytic;*
2. *if the subsets $\mathcal{A} \subseteq \mathcal{X}$ and $\tilde{\mathcal{A}} \subseteq \mathcal{Y}$ are analytic, then the sets $f(\mathcal{A})$ and $f^{-1}(\tilde{\mathcal{A}})$ are analytic.*

PROOF. From Proposition 13.7 in [30] it follows that every measurable set $\mathcal{A} \subseteq \mathcal{X}$ is analytic. From Proposition 14.4.(ii) in [30] it follows that the image and the preimage of an analytic set is an analytic set. $\qquad\square$

DEFINITION F.4 ($\mu$-measurability).  *Let $(\mathcal{X}, \Sigma, \mu)$ be a measure space. A set $\mathcal{E} \subseteq \mathcal{X}$ is called a $\mu$-null set if there exists a $\mathcal{A} \in \Sigma$ with $\mathcal{E} \subseteq \mathcal{A}$ and $\mu(\mathcal{A}) = 0$. We denote the class of $\mu$-null sets by $\mathcal{N}$, and we denote the $\sigma$-algebra generated by $\Sigma \cup \mathcal{N}$ by $\bar{\Sigma}$, and its members are called the $\mu$-measurable sets. Note that each member of $\bar{\Sigma}$ is of the form $\mathcal{A} \cup \mathcal{E}$ with $\mathcal{A} \in \Sigma$ and $\mathcal{E} \in \mathcal{N}$. The measure $\mu$ is extended to a measure $\bar{\mu}$ on $\bar{\Sigma}$, by $\bar{\mu}(\mathcal{A} \cup \mathcal{E}) = \mu(\mathcal{A})$ for every $\mathcal{A} \in \Sigma$ and $\mathcal{E} \in \mathcal{N}$, and is called its completion. A mapping $f : \mathcal{X} \to \mathcal{Y}$ between measurable spaces is called $\mu$-measurable if the inverse image $f^{-1}(\mathcal{C})$ of every measurable set $\mathcal{C} \subseteq \mathcal{Y}$ is $\mu$-measurable.*

DEFINITION F.5 (Universal measurability).  *Let $(\mathcal{X}, \Sigma)$ be a standard measurable space. A set $\mathcal{A} \subseteq \mathcal{X}$ is called universally measurable if it is $\mu$-measurable for every $\sigma$-finite measure[23] $\mu$ on $\mathcal{X}$ (i.e., in particular every probability measure). A mapping $f : \mathcal{X} \to \mathcal{Y}$ between standard measurable spaces is universally measurable if it is $\mu$-measurable for every $\sigma$-finite measure $\mu$.*

LEMMA F.6.  *Let $\mathcal{E}$ be a standard probability space with probability measure $\mathbb{P}_{\mathcal{E}}$ and $\mathcal{A} \subseteq \mathcal{E}$ an analytic set. Then $\mathcal{A}$ is $\mathbb{P}_{\mathcal{E}}$-measurable and there exist measurable sets $\mathcal{S}, \mathcal{T} \subseteq \mathcal{E}$ such that $\mathcal{S} \subseteq \mathcal{A} \subseteq \mathcal{T}$ and $\mathbb{P}_{\mathcal{E}}(\mathcal{S}) = \bar{\mathbb{P}}_{\mathcal{E}}(\mathcal{A}) = \mathbb{P}_{\mathcal{E}}(\mathcal{T})$, where $\bar{\mathbb{P}}_{\mathcal{E}}$ is the completion of $\mathbb{P}_{\mathcal{E}}$.*

PROOF. Let $\mathcal{A} \subseteq \mathcal{E}$ be an analytic set. Since every analytic set in a standard measurable space is a universally measurable set (see Theorem 21.10 in [30]), we know that $\mathcal{A}$ is a universally measurable set, and hence it is in particular a $\mathbb{P}_{\mathcal{E}}$-measurable set. Thus, there exist a measurable set $\mathcal{S} \subseteq \mathcal{E}$ and a $\mathbb{P}_{\mathcal{E}}$-null set $\mathcal{C} \subseteq \mathcal{E}$ such that $\mathcal{A} = \mathcal{S} \cup \mathcal{C}$ and $\bar{\mathbb{P}}_{\mathcal{E}}(\mathcal{A}) = \mathbb{P}_{\mathcal{E}}(\mathcal{S})$, where $\bar{\mathbb{P}}_{\mathcal{E}}$ is the completion of $\mathbb{P}_{\mathcal{E}}$. Moreover, there exists a measurable set $\tilde{\mathcal{C}} \subseteq \mathcal{E}$ such that $\mathcal{C} \subseteq \tilde{\mathcal{C}}$ and $\mathbb{P}_{\mathcal{E}}(\tilde{\mathcal{C}}) = 0$. Let $\mathcal{T} := \mathcal{S} \cup \tilde{\mathcal{C}}$, then $\mathcal{A} \subseteq \mathcal{T}$ and $\mathbb{P}_{\mathcal{E}}(\mathcal{T}) = \mathbb{P}_{\mathcal{E}}(\mathcal{S})$. $\qquad\square$

LEMMA F.7.  *Let $f : \mathcal{X} \to \mathcal{Y}$ be a $\mu$-measurable mapping. If $\mathcal{Y}$ is countably generated, then there exists a measurable mapping $g : \mathcal{X} \to \mathcal{Y}$ such that $f(x) = g(x)$ holds $\mu$-a.e..*

PROOF. Let the $\sigma$-algebra of $\mathcal{Y}$ be generated by the countable generating set $\{\mathcal{C}_n\}_{n \in \mathbb{N}}$. The $\mu$-measurable set $f^{-1}(\mathcal{C}_n) = \mathcal{A}_n \cup \mathcal{E}_n$ for some $\mathcal{A}_n \in \Sigma$ and some $\mathcal{E}_n \in \mathcal{N}$ and hence there is some $\mathcal{E}_n \subseteq \mathcal{B}_n \in \Sigma$ such that $\mu(\mathcal{B}_n) = 0$. Let $\hat{\mathcal{B}} = \cup_{n \in \mathbb{N}} \mathcal{B}_n$, $\hat{\mathcal{A}}_n = \mathcal{A}_n \setminus \hat{\mathcal{B}}$ and $\hat{\mathcal{A}} = \cup_{n \in \mathbb{N}} \hat{\mathcal{A}}_n$, then $\mu(\hat{\mathcal{B}}) = 0$, $\hat{\mathcal{A}}$ and $\hat{\mathcal{B}}$ are disjoint and $\mathcal{X} = \hat{\mathcal{A}} \cup \hat{\mathcal{B}}$. Now define the mapping $g : \mathcal{X} \to \mathcal{Y}$ by

$$g(x) := \begin{cases} f(x) & \text{if } x \in \hat{\mathcal{A}}, \\ y_0 & \text{otherwise,} \end{cases}$$

---

[23] A measure $\mu$ on a measurable space $(\mathcal{X}, \Sigma)$ is called $\sigma$-finite if $\mathcal{X} = \cup_{n \in \mathbb{N}} \mathcal{A}_n$, with $\mathcal{A}_n \in \Sigma$, $\mu(\mathcal{A}_n) < \infty$.

where for $y_0$ we can take an arbitrary point in $\mathcal{Y}$. This mapping $g$ is measurable since for each generator $\mathcal{C}_n$ we have

$$g^{-1}(\mathcal{C}_n) = \begin{cases} \hat{\mathcal{A}}_n & \text{if } y_0 \notin \mathcal{C}_n, \\ \hat{\mathcal{A}}_n \cup \hat{\mathcal{B}} & \text{otherwise.} \end{cases}$$

is in $\Sigma$. Moreover, $f(x) = g(x)$ $\mu$-almost everywhere. $\qquad\square$

With this result at hand we can now prove the first measurable selection theorem.

THEOREM F.8 (Measurable selection theorem). *Let $\mathcal{E}$ be a standard probability space with probability measure $\mathbb{P}_{\mathcal{E}}$, $\mathcal{X}$ a standard measurable space and $\mathcal{S} \subseteq \mathcal{E} \times \mathcal{X}$ a measurable set such that $\mathcal{E} \setminus pr_{\mathcal{E}}(\mathcal{S})$ is a $\mathbb{P}_{\mathcal{E}}$-null set, where $pr_{\mathcal{E}} : \mathcal{E} \times \mathcal{X} \to \mathcal{E}$ is the projection mapping on $\mathcal{E}$. Then there exists a measurable mapping $g : \mathcal{E} \to \mathcal{X}$ such that $(e, g(e)) \in \mathcal{S}$ for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$.*

PROOF. Take the subset $\hat{\mathcal{E}} := \mathcal{E} \setminus \mathcal{B}$, for some measurable set $\mathcal{B} \supseteq \mathcal{E} \setminus pr_{\mathcal{E}}(\mathcal{S})$ and $\mathbb{P}_{\mathcal{E}}(\mathcal{B}) = 0$, and note that $\hat{\mathcal{E}}$ is a standard measurable space (see Corollary 13.4 in [30]) and $\hat{\mathcal{E}} \subseteq pr_{\mathcal{E}}(\mathcal{S})$. Let $\hat{\mathcal{S}} = \mathcal{S} \cap (\hat{\mathcal{E}} \times \mathcal{X})$. Because the set $\hat{\mathcal{S}}$ is measurable, it is in particular analytic (see Lemma F.3). It follows by the Jankov-von Neumann Theorem (see Theorem 18.8 or 29.9 in [30]) that $\hat{\mathcal{S}}$ has a universally measurable uniformizing function, that is, there exists a universally measurable mapping $\hat{g} : \hat{\mathcal{E}} \to \mathcal{X}$ such that for all $e \in \hat{\mathcal{E}}$, $(e, \hat{g}(e)) \in \hat{\mathcal{S}}$. Hence, in particular, it is $\mathbb{P}_{\mathcal{E}}\big|_{\hat{\mathcal{E}}}$-measurable, where $\mathbb{P}_{\mathcal{E}}\big|_{\hat{\mathcal{E}}}$ is the restriction of $\mathbb{P}_{\mathcal{E}}$ to $\hat{\mathcal{E}}$.

Now define the mapping $g^* : \mathcal{E} \to \mathcal{X}$ by

$$g^*(e) := \begin{cases} \hat{g}(e) & \text{if } e \in \hat{\mathcal{E}} \\ x_0 & \text{otherwise,} \end{cases}$$

where for $x_0$ we can take an arbitrary point in $\mathcal{X}$. Then this mapping $g^*$ is $\mathbb{P}_{\mathcal{E}}$-measurable. To see this, take any measurable set $\mathcal{C} \subseteq \mathcal{X}$, then

$$g^{*-1}(\mathcal{C}) = \begin{cases} \hat{g}^{-1}(\mathcal{C}) & \text{if } x_0 \notin \mathcal{C} \\ \hat{g}^{-1}(\mathcal{C}) \cup \mathcal{B} & \text{otherwise.} \end{cases}$$

Because $\hat{g}^{-1}(\mathcal{C})$ is $\mathbb{P}_{\mathcal{E}}\big|_{\hat{\mathcal{E}}}$-measurable it is also $\mathbb{P}_{\mathcal{E}}$-measurable and thus $g^{*-1}(\mathcal{C})$ is $\mathbb{P}_{\mathcal{E}}$-measurable.

By Lemma F.7 and the fact that standard measurable spaces are countably generated (see Proposition 12.1 in [30]), we prove the existence of a measurable mapping $g : \mathcal{E} \to \mathcal{X}$ such that $g^* = g$ $\mathbb{P}_{\mathcal{E}}$-a.e. and thus it satisfies $(e, g(e)) \in \mathcal{S}$ for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$. $\qquad\square$

This theorem rests on the assumption that the standard measurable space $\mathcal{E}$ has a probability measure $\mathbb{P}_{\mathcal{E}}$. If this space becomes the product space $\mathcal{Y} \times \mathcal{E}$, for some standard measurable space $\mathcal{Y}$ where only the space $\mathcal{E}$ has a probability measure, then in general this theorem does not hold anymore. However, if we assume in addition that the fibers of $\mathcal{S}$ in $\mathcal{Y}$ are $\sigma$-compact for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$, then we can prove a second measurable selection theorem. A topological space is $\sigma$-compact if it is the union of countably many compact subspaces. For example, all countable discrete spaces, every interval of the real line, and moreover all the Euclidean spaces are $\sigma$-compact spaces.

THEOREM F.9 (Second measurable selection theorem). *Let $\mathcal{E}$ be a standard probability space with probability measure $\mathbb{P}_{\mathcal{E}}$, $\mathcal{X}$ and $\mathcal{Y}$ standard measurable spaces and $\mathcal{S} \subseteq \mathcal{X} \times \mathcal{E} \times \mathcal{Y}$ a measurable set such that $\mathcal{E} \setminus \mathcal{K}_\sigma$ is a $\mathbb{P}_{\mathcal{E}}$-null set, where*

$$\mathcal{K}_\sigma := \{e \in \mathcal{E} : \forall x \in \mathcal{X}(\mathcal{S}_{(x,e)} \text{ is non-empty and } \sigma\text{-compact})\},$$

*with $\mathcal{S}_{(x,e)}$ denoting the fiber over $(x, e)$, that is*

$$\mathcal{S}_{(x,e)} := \{y \in \mathcal{Y} : (x, e, y) \in \mathcal{S}\}.$$

*Then there exists a measurable mapping $g : \mathcal{X} \times \mathcal{E} \to \mathcal{Y}$ such that for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$ we have $(x, e, g(x, e)) \in \mathcal{S}$.*

PROOF. Take the subset $\hat{\mathcal{E}} := \mathcal{E} \setminus \mathcal{B}$, for some measurable set $\mathcal{B} \supseteq \mathcal{E} \setminus \mathcal{K}_\sigma$ and $\mathbb{P}_{\mathcal{E}}(\mathcal{B}) = 0$. Note that $\hat{\mathcal{E}}$ is a standard measurable space, $\hat{\mathcal{E}} \subseteq \mathcal{K}_\sigma$ and $\hat{\mathcal{S}} = \mathcal{S} \cap (\mathcal{X} \times \hat{\mathcal{E}} \times \mathcal{Y})$ is measurable. By assumption, for each $(x, e) \in \mathcal{X} \times \hat{\mathcal{E}}$ the fiber $\hat{\mathcal{S}}_{(x,e)}$ is non-empty and $\sigma$-compact and hence by applying the Theorem of Arsenin-Kunugui (see Theorem 35.46 in [30]) it follows that the set $\hat{\mathcal{S}}$ has a measurable uniformizing function, that is, there exists a measurable mapping $\hat{g} : \mathcal{X} \times \hat{\mathcal{E}} \to \mathcal{Y}$ such that for all $(x, e) \in \mathcal{X} \times \hat{\mathcal{E}}$, $(x, e, \hat{g}(x, e)) \in \hat{\mathcal{S}}$. Now define the mapping $g : \mathcal{X} \times \mathcal{E} \to \mathcal{Y}$ by

$$g(x, e) := \begin{cases} \hat{g}(x, e) & \text{if } e \in \hat{\mathcal{E}} \\ y_0 & \text{otherwise,} \end{cases}$$

where for $y_0$ we can take an arbitrary point in $\mathcal{Y}$. This mapping $g$ inherits the measurability from $\hat{g}$ and it satisfies for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$ and for all $x \in \mathcal{X}$ that $(x, e, g(x, e)) \in \mathcal{S}$. $\qquad\square$

The next two lemmas provide some useful properties for the "for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$" quantifier.

LEMMA F.10. *Let $\phi : \mathcal{E} \to \tilde{\mathcal{E}}$ be a measurable map between two standard measurable spaces. Let $\mathbb{P}_{\mathcal{E}}$ be a probability measure on $\mathcal{E}$ and let $\mathbb{P}_{\tilde{\mathcal{E}}} = \mathbb{P}_{\mathcal{E}} \circ \phi^{-1}$ be its push-forward under $\phi$. Let $\tilde{P} : \tilde{\mathcal{E}} \to \{0, 1\}$ be a property, i.e., a (measurable) boolean-valued function on $\tilde{\mathcal{E}}$. Then the property $P = \tilde{P} \circ \phi$ on $\mathcal{E}$ holds $\mathbb{P}_{\mathcal{E}}$-a.e. if and only if the property $\tilde{P}$ holds $\mathbb{P}_{\tilde{\mathcal{E}}}$-a.e..*

PROOF. Assume the property $P = \tilde{P} \circ \phi$ holds $\mathbb{P}_{\mathcal{E}}$-a.e., then $\mathcal{C} = \{e \in \mathcal{E} : P(e) = 1\}$ contains a measurable set $\mathcal{C}^*$ with $\mathbb{P}_{\mathcal{E}}$-measure 1, i.e., $\mathcal{C}^* \subseteq \mathcal{C}$ and $\mathbb{P}_{\mathcal{E}}(\mathcal{C}^*) = 1$. By Lemma F.3, $\phi(\mathcal{C}^*)$ is analytic. By Lemma F.6, there exist measurable sets $\mathcal{A}, \mathcal{B}$ such that $\mathcal{A} \subseteq \phi(\mathcal{C}^*) \subseteq \mathcal{B}$ and $\mathbb{P}_{\tilde{\mathcal{E}}}(\mathcal{A}) = \mathbb{P}_{\tilde{\mathcal{E}}}(\mathcal{B})$. Because $\phi$ is measurable, $\phi^{-1}(\mathcal{A})$ and $\phi^{-1}(\mathcal{B})$ are both measurable. Also, $\phi^{-1}(\mathcal{A}) \subseteq \phi^{-1}(\phi(\mathcal{C}^*)) \subseteq \phi^{-1}(\mathcal{B})$. As $\mathcal{C}^* \subseteq \phi^{-1}(\phi(\mathcal{C}^*))$, we must have that $\mathbb{P}_{\mathcal{E}}(\phi^{-1}(\mathcal{B})) \geq \mathbb{P}_{\mathcal{E}}(\mathcal{C}^*) = 1$. Hence $\mathbb{P}_{\tilde{\mathcal{E}}}(\mathcal{A}) = \mathbb{P}_{\tilde{\mathcal{E}}}(\mathcal{B}) = 1$. Note that as $\mathcal{C}^* \subseteq \mathcal{C}$, $\mathcal{A} \subseteq \phi(\mathcal{C}^*) \subseteq \phi(\mathcal{C}) \subseteq \{\tilde{e} \in \tilde{\mathcal{E}} : \tilde{P}(\tilde{e}) = 1\}$. Hence the set $\tilde{\mathcal{C}} := \{\tilde{e} \in \tilde{\mathcal{E}} : \tilde{P}(\tilde{e}) = 1\}$ contains a measurable set of $\mathbb{P}_{\tilde{\mathcal{E}}}$-measure 1, in other words, $\tilde{P}$ holds $\mathbb{P}_{\tilde{\mathcal{E}}}$-a.s..

The converse is easier to prove. Suppose $\tilde{\mathcal{C}} = \{\tilde{e} \in \tilde{\mathcal{E}} : \tilde{P}(\tilde{e}) = 1\}$ contains a measurable set $\tilde{\mathcal{C}}^*$ with $\mathbb{P}_{\tilde{\mathcal{E}}}$-measure 1, i.e., $\tilde{\mathcal{C}}^* \subseteq \tilde{\mathcal{C}}$ and $\mathbb{P}_{\tilde{\mathcal{E}}}(\tilde{\mathcal{C}}^*) = 1$. Because $\phi$ is measurable, the set $\phi^{-1}(\tilde{\mathcal{C}}^*)$ is measurable and $\mathbb{P}_{\mathcal{E}}(\phi^{-1}(\tilde{\mathcal{C}}^*)) = 1$, and furthermore, $\phi^{-1}(\tilde{\mathcal{C}}^*) \subseteq \phi^{-1}(\tilde{\mathcal{C}}) = \mathcal{C}$. $\qquad\square$

LEMMA F.11 (Some properties for the for-almost-every quantifier). *Let $\boldsymbol{\mathcal{X}} = \mathcal{X} \times \tilde{\mathcal{X}}$ and $\boldsymbol{\mathcal{E}} = \mathcal{E} \times \tilde{\mathcal{E}}$ be products of non-empty standard measurable spaces and $\mathbb{P}_{\boldsymbol{\mathcal{E}}} = \mathbb{P}_{\mathcal{E}} \times \mathbb{P}_{\tilde{\mathcal{E}}}$ be the product measure of probability measures $\mathbb{P}_{\mathcal{E}}$ and $\mathbb{P}_{\tilde{\mathcal{E}}}$ on $\mathcal{E}$ and $\tilde{\mathcal{E}}$ respectively. Denote by "$\forall\!\!\!\!\forall e$" the quantifier "for $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-almost every $\boldsymbol{e} \in \boldsymbol{\mathcal{E}}$" and by "$\forall \boldsymbol{x}$" the quantifier "for all $\boldsymbol{x} \in \boldsymbol{\mathcal{X}}$", and similarly for their components, e.g. "$\forall\!\!\!\!\forall e$" for "for $\mathbb{P}_{\mathcal{E}}$-almost every $e \in \mathcal{E}$" and "$\forall x$" for "for all $x \in \mathcal{X}$". Then we have the following properties:*

1. $\forall\!\!\!\!\forall e : P(e) \implies \exists e : P(e)$ *(similarly to $\forall x : P(x) \implies \exists x : P(x)$)*;
2. $\forall\!\!\!\!\forall e : P(e) \iff \forall\!\!\!\!\forall \boldsymbol{e} : P(e)$ *(similarly to $\forall x : P(x) \iff \forall \boldsymbol{x} : P(x)$)*;
3. $\exists x \forall\!\!\!\!\forall e : P(x,e) \implies \forall\!\!\!\!\forall e \exists x : P(x,e)$ *(similarly to $\exists x \forall e : P(x,e) \implies \forall e \exists x : P(x,e)$)*;
4. $\forall\!\!\!\!\forall e \forall x : P(x,e) \implies \forall x \forall\!\!\!\!\forall e : P(x,e)$ *(similarly to $\forall e \forall x : P(x,e) \implies \forall x \forall e : P(x,e)$)*;
5. $\forall\!\!\!\!\forall \boldsymbol{e} : P(\boldsymbol{e}) \implies \exists \tilde{e} \forall\!\!\!\!\forall e : P(\boldsymbol{e})$ *(similarly to $\forall \boldsymbol{x} : P(\boldsymbol{x}) \implies \exists \tilde{x} \forall x : P(\boldsymbol{x})$)*;
6. $\forall\!\!\!\!\forall e \forall x : P(x,e) \iff \forall\!\!\!\!\forall \boldsymbol{e} \forall \boldsymbol{x} : P(x,e)$;
7. $\forall\!\!\!\!\forall \boldsymbol{e} \forall \boldsymbol{x} : P(\boldsymbol{x}, \boldsymbol{e}) \implies \exists \tilde{e} \exists \tilde{x} \forall\!\!\!\!\forall e \forall x : P(\boldsymbol{x}, \boldsymbol{e})$,

*where $P$ denotes a property, i.e., a measurable boolean-valued function, on the corresponding measurable spaces and we write $\boldsymbol{e}$ and $\boldsymbol{x}$ for $(e, \tilde{e})$ and $(x, \tilde{x})$ respectively.*

PROOF. We only prove the statements that may not be immediately obvious.

Property 2. Let $pr_{\mathcal{E}} : \boldsymbol{\mathcal{E}} \to \mathcal{E}$ be the projection mapping on $\mathcal{E}$. Then by Lemma F.10 we have

$$\forall\!\!\!\!\forall e : P(e) \iff \forall\!\!\!\!\forall \boldsymbol{e} : P \circ pr_{\mathcal{E}}(\boldsymbol{e}) \iff \forall\!\!\!\!\forall \boldsymbol{e} : P(e).$$

Property 4: We have

$$\forall\!\!\!\!\forall e \forall x : P(x,e)$$
$$\implies \exists \, \mathbb{P}_{\mathcal{E}}\text{-null set } N \, \forall e \in \mathcal{E} \setminus N \, \forall x : P(x,e)$$
$$\implies \exists \, \mathbb{P}_{\mathcal{E}}\text{-null set } N \, \forall x \, \forall e \in \mathcal{E} \setminus N : P(x,e)$$
$$\implies \forall x \, \exists \, \mathbb{P}_{\mathcal{E}}\text{-null set } N \, \forall e \in \mathcal{E} \setminus N : P(x,e)$$
$$\implies \forall x \forall\!\!\!\!\forall e : P(x,e).$$

Property 5: Let $\boldsymbol{N}$ be a measurable $\mathbb{P}_{\boldsymbol{\mathcal{E}}}$-null set such that $P(\boldsymbol{e})$ holds for all $\boldsymbol{e} \in \boldsymbol{\mathcal{E}} \setminus \boldsymbol{N}$. Define for $\tilde{e} \in \tilde{\mathcal{E}}$ the set $N_{\tilde{e}} := \{e \in \mathcal{E} : (e, \tilde{e}) \in \boldsymbol{N}\}$. Note that the sets $N_{\tilde{e}}$ are measurable. From Fubini's theorem it follows that for $\mathbb{P}_{\tilde{\mathcal{E}}}$-almost every $\tilde{e} \in \tilde{\mathcal{E}}$ we have $\mathbb{P}_{\mathcal{E}}(N_{\tilde{e}}) = 0$. That is, there exists a measurable $\mathbb{P}_{\tilde{\mathcal{E}}}$-null set $\tilde{N}$ such that $\mathbb{P}_{\mathcal{E}}(N_{\tilde{e}}) = 0$ for all $\tilde{e} \in \tilde{\mathcal{E}} \setminus \tilde{N}$. Hence, there exists $\tilde{e} \in \tilde{\mathcal{E}} \setminus \tilde{N}$ such that $\mathbb{P}_{\mathcal{E}}(N_{\tilde{e}}) = 0$; for all $e \in \mathcal{E} \setminus N_{\tilde{e}}$, $P(\boldsymbol{e})$ then holds. This means $\exists \tilde{e} \forall\!\!\!\!\forall e : P(\boldsymbol{e})$.

Property 7: We have

$$\forall\!\!\!\!\forall \boldsymbol{e} \forall \boldsymbol{x} : P(\boldsymbol{x}, \boldsymbol{e}) \implies \exists \tilde{e} \forall\!\!\!\!\forall e \forall \boldsymbol{x} : P(\boldsymbol{x}, \boldsymbol{e}) \implies \exists \tilde{e} \forall\!\!\!\!\forall e \forall \tilde{x} \forall x : P(\boldsymbol{x}, \boldsymbol{e})$$
$$\implies \exists \tilde{e} \forall \tilde{x} \forall\!\!\!\!\forall e \forall x : P(\boldsymbol{x}, \boldsymbol{e}) \implies \exists \tilde{e} \exists \tilde{x} \forall\!\!\!\!\forall e \forall x : P(\boldsymbol{x}, \boldsymbol{e}),$$

where in the first equivalence we used Property 5, in the third equivalence we used Property 4 and in the last equivalence we used Property 1. □

## REFERENCES

[1] BALKE, A. and PEARL, J. (1994). Probabilistic Evaluation of Counterfactual Queries. In *Proceedings of the Twelfth National Conference on Artificial Intelligence (AAAI-94)* **1** 230–237. AAAI Press.

[2] BECKERS, S. and HALPERN, J. Y. (2019). Abstracting Causal Models. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)* **33** 2678–2685. AAAI Press.

[3] BLOM, T., BONGERS, S. and MOOIJ, J. M. (2019). Beyond Structural Causal Models: Causal Constraints Models. In *Proceedings of the 35th Conference on Uncertainty in Artificial Intelligence (UAI-19)* (R. P. ADAMS and V. GOGATE, eds.). AUAI Press.

[4] BLOM, T., VAN DIEPEN, M. M. and MOOIJ, J. M. (2020). Conditional Independences and Causal Relations implied by Sets of Equations. *arXiv.org preprint* arXiv:2007.07183 [cs.AI].

[5] BOLLEN, K. A. (1989). *Structural Equations with Latent Variables*. John Wiley & Sons, New York, USA.

[6] BONGERS, S. and MOOIJ, J. M. (2018). From Random Differential Equations to Structural Causal Models: the stochastic case. *arXiv.org preprint* arXiv:1803.08784 [cs.AI].

[7] BÜHLMANN, P., PETERS, J. and ERNEST, J. (2014). CAM: Causal Additive Models, high-dimensional order search and penalized regression. *The Annals of Statistics* **42** 2526–2556.

[8] BYRNE, R. M. J. (2007). *The Rational Imagination: How People Create Alternatives to Reality. A Bradford Book*. MIT Press, Cambridge, MA.

[9] COHN, D. L. (2013). *Measure Theory*, 2nd ed. Birkhäuser, Boston, USA.

[10] COOPER, G. F. (1997). A Simple Constraint-Based Algorithm for Efficiently Mining Observational Databases for Causal Relationships. *Data Mining and Knowledge Discovery* **1** 203–224.

[11] DAWID, A. P. (2002). Influence Diagrams for Causal Modelling and Inference. *International Statistical Review* **70** 161–189.

[12] DUNCAN, O. D. (1975). *Introduction to Structural Equation Models*. Academic Press, New York.

[13] EATON, D. and MURPHY, K. (2007). Exact Bayesian structure learning from uncertain interventions. In *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics* (M. MEILA and X. SHEN, eds.). *Proceedings of Machine Learning Research* **2** 107–114.

[14] EBERHARDT, F., HOYER, P. and SCHEINES, R. (2010). Combining Experiments to Discover Linear Cyclic Models with Latent Variables. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics* (Y. W. TEH and M. TITTERINGTON, eds.). *Proceedings of Machine Learning Research* **9** 185–192.

[15] EVANS, R. J. (2016). Graphs for Margins of Bayesian Networks. *Scandinavian Journal of Statistics* **43** 625–648.

[16] EVANS, R. J. (2018). Margins of discrete Bayesian networks. *The Annals of Statistics* **46** 2623–2656.

[17] FISHER, F. M. (1970). A Correspondence Principle For Simultaneous Equation Models. *Econometrica* **38** 73–92.

[18] FORRÉ, P. and MOOIJ, J. M. (2017). Markov Properties for Graphical Models with Cycles and Latent Variables. *arXiv.org preprint* arXiv:1710.08775 [math.ST].

[19] FORRÉ, P. and MOOIJ, J. M. (2018). Constraint-based Causal Discovery for Non-Linear Structural Causal Models with Cycles and Latent Confounders. In *Proceedings of the 34th Conference on Uncertainty in Artificial Intelligence (UAI-18)* (A. GLOBERSON and R. SILVA, eds.). AUAI Press.

[20] FORRÉ, P. and MOOIJ, J. M. (2019). Causal Calculus in the Presence of Cycles, Latent Confounders and Selection Bias. In *Proceedings of the 35th Conference on Uncertainty in Artificial Intelligence (UAI-19)* (R. P. ADAMS and V. GOGATE, eds.). AUAI Press.

[21] FOYGEL, R., DRAISMA, J. and DRTON, M. (2012). Half-trek Criterion for Generic Identifiability of Linear Structural Equation Models. *The Annals of Statistics* **40** 1682–1713.

[22] GEIGER, D. (1990). Graphoids: A Qualitative Framework for Probabilistic Inference Technical Report No. R-142, Computer Science Department, University of California, Los Angeles, USA.

[23] GOLDBERGER, A. S. and DUNCAN, O. D. (1973). *Structural Equation Models in the Social Sciences*. Seminar Press, New York.

[24] GOLUB, G. and KAHAN, W. (1965). Calculating the Singular Values and Pseudo-Inverse of a Matrix. *Journal of the Society for Industrial and Applied Mathematics: Series B, Numerical Analysis* **2** 205–224.

[25] HAAVELMO, T. (1943). The Statistical Implications of a System of Simultaneous Equations. *Econometrica* **11** 1–12.

[26] HALPERN, J. (1998). Axiomatizing Causal Reasoning. In *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence (UAI-98)* (G. COOPER and S. MORAL, eds.) 202–210. Morgan Kaufmann, San Francisco, CA, USA.

[27] HYTTINEN, A., EBERHARDT, F. and HOYER, P. O. (2012). Learning Linear Cyclic Causal Models with Latent Variables. *Journal of Machine Learning Research* **13** 3387–3439.

[28] HYTTINEN, A., HOYER, P. O., EBERHARDT, F. and JÄRVISALO, M. (2013). Discovering Cyclic Causal Models with Latent Variables: A General SAT-based Procedure. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence (UAI-13)* (A. NICHOLSON and P. SMYTH, eds.) 301–310. AUAI Press, Corvallis, Oregon, USA.

[29] IWASAKI, Y. and SIMON, H. A. (1994). Causality and model abstraction. *Artificial Intelligence* **67** 143–194.

[30] KECHRIS, A. S. (1995). *Classical Descriptive Set Theory. Graduate Texts in Mathematics* **156**. Springer-Verlag, New York, USA.

[31] KOSTER, J. T. A. (1996). Markov Properties of Nonrecursive Causal Models. *The Annals of Statistics* **24** 2148–2177.

[32] KOSTER, J. T. A. (1999). On the Validity of the Markov Interpretation of Path Diagrams of Gaussian Structural Equations Systems with Correlated Errors. *Scandinavian Journal of Statistics* **26** 413–431.

[33] LACERDA, G., SPIRTES, P. L., RAMSEY, J. and HOYER, P. O. (2008). Discovering cyclic causal models by independent components analysis. In *Proceedings of the Twenty-Fourth Conference on Uncertainty in Artificial Intelligence (UAI-08)* (D. MCALLESTER and P. MYLLYMAKI, eds.) 366–374. AUAI Press, Corvallis, Oregon, USA.

[34] LAURITZEN, S. L. (1996). *Graphical Models. Oxford Statistical Science Series* **17**. Clarendon Press, Oxford.

[35] LAURITZEN, S. L., DAWID, A. P., LARSEN, B. N. and LEIMER, H. G. (1990). Independence Properties of Directed Markov Fields. *Networks* **20** 491–505.

[36] LEWIS, D. K. (1979). Counterfactual Dependence and Time's Arrow. *Noûs* **13** 455–476.

[37] MAATHUIS, M. H., COLOMBO, D., KALISCH, M. and BÜHLMANN, P. (2009). Estimating High-Dimensional Intervention Effects from Observational Data. *The Annals of Statistics* **37** 3133–3164.

[38] MANI, S. (2006). A Bayesian Local Causal Discovery Framework, PhD thesis, University of Pittsburg.

[39] MASON, S. J. (1953). Feedback Theory - Some Properties of Signal Flow Graphs. In *Proceedings of the IRE* **41** 1144-1156. IEEE.

[40] MASON, S. J. (1956). Feedback Theory - Further Properties of Signal Flow Graphs. In *Proceedings of the IRE* **44** 920–926. IEEE.

[41] MEEK, C. (1995). Strong Completeness and Faithfulness in Bayesian Networks. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence (UAI-95)* (P. BESNARD and S. HANKS, eds.) 411–418. Morgan Kaufmann, San Francisco, CA, USA.

[42] MOGENSEN, S. W. and HANSEN, N. R. (2020). Markov equivalence of marginalized local independence graphs. *Ann. Statist.* **48** 539–559.

[43] MOGENSEN, S. W., MALINSKY, D. and HANSEN, N. R. (2018). Causal Learning for Partially Observed Stochastic Dynamical Systems. In *Proceedings of the Thirty-Fourth conference on Uncertainty in Artificial Intelligence (UAI-18)* (A. GLOBERSON and R. SILVA, eds.). AUAI Press.

[44] MOOIJ, J. M. and CLAASSEN, T. (2020). Constraint-Based Causal Discovery using Partial Ancestral Graphs in the presence of Cycles. In *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI-20)* (J. PETERS and D. SONTAG, eds.) **124** 1159–1168. PMLR.

[45] MOOIJ, J. M. and HESKES, T. (2013). Cyclic Causal Discovery from Continuous Equilibrium Data. In *Proceedings of the 29th Conference on Uncertainty in Artificial Intelligence (UAI-13)* (A. NICHOLSON and P. SMYTH, eds.) 431–439. AUAI Press, Corvallis, Oregon, USA.

[46] MOOIJ, J. M., JANZING, D. and SCHÖLKOPF, B. (2013). From Ordinary Differential Equations to Structural Causal Models: the deterministic case. In *Proceedings of the 29th Conference on Uncertainty in Artificial Intelligence (UAI-13)* (A. NICHOLSON and P. SMYTH, eds.) 440–448. AUAI Press.

[47] MOOIJ, J. M., MAGLIACANE, S. and CLAASSEN, T. (2020). Joint Causal Inference from Multiple Contexts. *Journal of Machine Learning Research* **21** 1–108.

[48] MOOIJ, J. M., PETERS, J., JANZING, D., ZSCHEISCHLER, J. and SCHÖLKOPF, B. (2016). Distinguishing Cause from Effect using Observational Data: Methods and Benchmarks. *Journal of Machine Learning Research* **17** 1–102.

[49] NEAL, R. M. (2000). On Deducing Conditional Independence from $d$-Separation in Causal Graphs with Feedback. *Journal of Artificial Intelligence Research* **12** 87–91.

[50] PEARL, J. (1985). A Constraint Propagation Approach to Probabilistic Reasoning. In *Proceedings of the First Conference on Uncertainty in Artificial Intelligence (UAI-85)* (L. KANAL and J. LEMMER, eds.) 31–42. AUAI Press, Corvallis, Oregon, USA.

[51] PEARL, J. (2009). *Causality: Models, Reasoning, and Inference*, 2nd ed. Cambridge University Press, New York, USA.

[52] PEARL, J. and DECHTER, R. (1996). Identifying Independence in Causal Graphs with Feedback. In *Proceedings of the Twelfth Conference on Uncertainty in Artificial Intelligence (UAI-96)* (E. HORVITZ and F. JENSEN, eds.) 420–426. Morgan Kaufmann, San Francisco, CA, USA.

[53] PEARL, J. and MACKENZIE, D. (2018). *The Book of Why: The New Science of Cause and Effect*, 1st ed. Basic Books, New York, USA.

[54] PENROSE, R. (1955). A generalized inverse for matrices. *Mathematical Proceedings of the Cambridge Philosophical Society* **51** 406–413.

[55] PETERS, J., JANZING, D. and SCHÖLKOPF, B. (2017). *Elements of Causal Inference: Foundations and Learning Algorithms*. MIT Press, Cambridge, MA, USA.

[56] PETERS, J., MOOIJ, J. M., JANZING, D. and SCHÖLKOPF, B. (2014). Causal Discovery with Continuous Additive Noise Models. *Journal of Machine Learning Research* **15** 2009–2053.

[57] PFISTER, N., BAUER, S. and PETERS, J. (2019). Learning Stable and Predictive Structures in Kinetic Systems. *Proceedings of the National Academy of Sciences* **116** 25405–25411.

[58] RICHARDSON, T. S. (1996). A Discovery Algorithm for Directed Cyclic Graphs. In *Proceedings of the Twelfth Conference on Uncertainty in Artificial Intelligence (UAI-96)* (E. HORVITZ and F. JENSEN, eds.) 454–461. Morgan Kaufmann, San Francisco, CA, USA.

[59] RICHARDSON, T. S. (1996). Discovering Cyclic Causal Structure Technical Report No. CMU-PHIL-68, Carnegie Mellon University.

[60] RICHARDSON, T. (2003). Markov Properties for Acyclic Directed Mixed Graphs. *Scandinavian Journal of Statistics* **30** 145–157.

[61] RICHARDSON, T. S. and SPIRTES, P. (1999). Automated Discovery of Linear Feedback Models. In *Computation, Causation, and Discovery* (C. Glymour and G. F. Cooper, eds.) 253-–304. MIT Press.

[62] RICHARDSON, T. S. and SPIRTES, P. (2002). Ancestral Graph Markov Models. *The Annals of Statistics* **30** 962–1030.

[63] RICHARDSON, T. S. (1996). Models of Feedback: Interpretation and Discovery, PhD thesis, Carnegie Mellon University.

[64] RICHARDSON, T. S. and ROBINS, J. (2013). Single World Intervention Graphs (SWIGs): A Unification of the Counterfactual and Graphical Approaches to Causality Technical Report No. 128, Center for Statistics and the Social Sciences.

[65] RICHARDSON, T. S. and ROBINS, J. M. (2014). ACE Bounds; SEMs with Equilibrium Conditions. *Statistical Science* **29** 363-366.

[66] ROESE, N. J. (1997). Counterfactual Thinking. *Psychological Bulletin* **121** 133–148.

[67] RUBENSTEIN, P. K., WEICHWALD, S., BONGERS, S., MOOIJ, J. M., JANZING, D., GROSSE-WENTRUP, M. and SCHÖLKOPF, B. (2017). Causal Consistency of Structural Equation Models. In *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence (UAI-17)* (G. ELIDAN and K. KERSTING, eds.). AUAI Press.

[68] RUBIN, D. B. (1974). Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies. *Journal of Educational Psychology* **66** 688–701.

[69] SHPITSER, I. and PEARL, J. (2008). Complete Identification Methods for the Causal Hierarchy. *Journal of Machine Learning Research* **9** 1941–1979.

[70] SPIRTES, P. (1993). Directed Cyclic Graphs, Conditional Independence, and Non-recursive Linear Structural Equation Models Technical Report No. CMU-PHIL-35, Carnegie Mellon University.

[71] SPIRTES, P. (1994). Conditional Independence in Directed Cyclic Graphical Models for Feedback Technical Report No. CMU-PHIL-54, Carnegie Mellon University.

[72] SPIRTES, P. (1995). Directed Cyclic Graphical Representations of Feedback Models. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence (UAI-95)* (P. BESNARD and S. HANKS, eds.) 499–506. Morgan Kaufmann, San Francisco, CA, USA.

[73] SPIRTES, P., GLYMOUR, C. and SCHEINES, R. (2000). *Causation, Prediction, and Search*, 2nd ed. *Adaptive Computation and Machine Learning*. MIT Press, Cambridge, Massachusetts.

[74] SPIRTES, P., MEEK, C. and RICHARDSON, T. S. (1999). An Algorithm for Causal Inference in the Presence of Latent Variables and Selection Bias. In *Computation, Causation and Discovery* (C. Glymour and G. F. Cooper, eds.) 6, 211-252. The MIT Press.

[75] SPIRTES, P., RICHARDSON, T., MEEK, C., SCHEINES, R. and GLYMOUR, C. (1998). Using Path Diagrams as a Structural Equation Modelling Tool. *Sociological Methods & Research* **27** 182–225.

[76] TIAN, J. (2002). Studies in Causal Reasoning and Learning Technical Report No. R-309, Cognitive Systems Laboratory, University of California, Los Angeles, USA.

[77] TIAN, J. and PEARL, J. (2001). Causal Discovery from Changes. In *Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence (UAI-01)* (J. BREESE and D. KOLLER, eds.) 512–521. Morgan Kaufmann, San Francisco, CA, USA.

[78] VERMA, T. S. (1993). Graphical Aspects of Causal Models Technical Report No. R-191, Computer Science Department, University of California, Los Angeles, USA.

[79] WRIGHT, S. (1921). Correlation and Causation. *Journal of Agricultural Research* **20** 557–585.

[80] ZHANG, J. (2008). On the Completeness of Orientation Rules for Causal Discovery in the Presence of Latent Confounders and Selection Bias. *Artificial Intelligence* **172** 1873–1896.