

# Causal Modeling With Infinitely Many Variables

Spencer Peters and Joseph Y. Halpern

Cornell University

Ithaca, NY 14850

sp2473@cornell.edu, halpern@cs.cornell.edu

## Abstract

Structural-equations models (SEMs) are perhaps the most commonly used framework for modeling causality. However, as we show, naively extending this framework to infinitely many variables, which is necessary, for example, to model dynamical systems, runs into several problems. We introduce *GSEMs* (*generalized SEMs*), a flexible generalization of SEMs that directly specify the results of interventions, in which (1) systems of differential equations can be represented in a natural and intuitive manner, (2) certain natural situations, which cannot be represented by SEMs at all, can be represented easily, (3) the definition of actual causality in SEMs carries over essentially without change.

## 1 Introduction

For scientists trying to understand causal relationships, and policymakers grappling with the consequences of their decisions, the structure of causality itself is of great importance. One influential paradigm for formalizing causality, *structural-equations models* (SEMs), describes causal relationships as a collection of *structural equations*. Actions taken by a scientist or policymaker (interventions) manifest as modifications to the structural equations; for example, if  $X$  represents the rent price of an apartment, imposing rent control amounts to replacing the equation for  $X$  with  $X = r$ , producing a new SEM.

SEMs are well studied, and there are standard techniques for reasoning about the outcomes of given interventions. For example, SEMs without cyclic dependencies have a unique outcome for any given intervention, which can be obtained simply by solving the equations in any order consistent with the dependencies between equations. However, this property does not hold when there are infinitely many variables. Consider the simple SEM with binary variables  $X_1, X_2, \dots$ , where the equations are  $X_1 = X_2, X_2 = X_3, \dots$ <sup>1</sup>. Intuitively, this says that  $X_1$  gets the value that  $X_2$  has,  $X_2$  gets the value that  $X_3$  has, and so on. There are clearly no acyclic

dependencies here; nevertheless, these equations have two solutions:  $X_i = 0$  for all  $i$ , and  $X_i = 1$  for all  $i$ . The problem here is that the dependency relation above is not well founded, since from  $X_1$  an infinite chain of dependencies  $X_2, X_3, \dots$  can be traced. Indeed, it is easy to establish can be traced:  $X_1$  depends on  $X_2$ , which depends on  $X_3$ , and so on. Indeed, it is easy to establish using that if the dependency relation is well-founded, then there will also be a unique outcome.

The problem of ill-founded dependencies is unavoidable when working with dynamical systems. Consider variables  $X_t$  representing the state of a dynamical system at time  $t$ , where  $t$  ranges over an interval of real numbers. In general, we expect that changing the value of  $X_s$  will affect the value of  $X_t$  for all  $t > s$ . That is, the dependency relation can be identified with the less-than relation on the real numbers, which is not well founded.

To capture dynamical systems, we propose a more flexible class of models that we call *generalized structural-equations models* (GSEMs). Given a SEM and an intervention, we can produce a new SEM that represents the result of the intervention by modifying the relevant equations. GSEMs represent the same input-output relationships as SEMs—a GSEM and an allowed intervention maps to a new GSEM, while a GSEM and a context together determine a set of possible assignments to the endogenous variables—without committing to a specific mechanism for producing the assignments. It is easy to show that GSEMs generalize SEMs (see Theorem 3.1).

The way that GSEMs are defined makes it easy to lift definitions of causal notions from SEMs to GSEMs. Any definition depending only on inputs and outputs (interventions and their outcomes) can be immediately applied to GSEMs. In particular, the notion of actual causality (whether event  $X$  caused event  $Y$  in some concrete situation) given by Halpern and Pearl [2005] can be applied to GSEMs almost without modification.

What is particularly significant for our purposes is that many standard formalisms for representing causality in infinitary settings can be captured with GSEMs, including dynamical systems, *hybrid automata* [Alur *et al.*, 1992] (a popular formalism for describing mixed discrete-continuous systems), and the *rule-based models* commonly used in molecular biology and organic chemistry (see [Laurent *et al.*, 2018] and the references therein). In particular, dynamical systems can be represented in a direct and natural way using GSEMs.

<sup>1</sup>SEMs are typically defined to have only finitely many variables; we relax this restriction for the purposes of illustration.

The definition (see Section 4) is nothing more than the textbook definition of the solution of a system of differential equations.

But even if we restrict to settings with only finitely many variables, GSEMs are more expressive than SEMs.

This flexibility can be critical if, for example, we want to model quantum systems. Consider, for example, a SEM describing a spin qubit using three ternary variables  $X, Y, Z$  that represent the value of the spin in the  $x, y, z$  directions, respectively. The values  $+1$  and  $-1$  represent “spin up” and “spin down”; the value  $0$  is interpreted as “not well defined”. A SEM with these variables has to specify what happens after the intervention  $X \leftarrow 1, Y \leftarrow 1$ , and must have  $X = 1$  and  $Y = 1$  after this intervention. But since the observables  $X$  and  $Y$  do not commute, quantum mechanics says that if the qubit is spin-up in the  $X$ -direction, then the spin in the  $Y$  direction is not well defined. That is, if  $X = 1$  then  $Y = 0$ . Hence, this system cannot be modeled by a SEM. However, a generalized SEM is easily constructed that corresponds to precisely this situation. We might hope to solve this problem by limiting the set of allowed interventions in a SEM, as is done, for example, by Rubinstein et al. [2017]. While that will work in this case, it does not solve the problem in general, as we show in Example 3.6. These problems can also be handled by *causal constraints models* (CCMs) [Blom et al., 2019], which were introduced to get around some of the restrictions of SEMs when modeling equilibrium solutions of dynamical systems. GSEMs and CCMs are actually equally expressive (Theorem 5.1). However, because GSEMs are a more straightforward generalization of SEMs than CCMs, they are arguably easier to use for those familiar with SEMs, and make it easy to carry over standard notions like the definition of actual causality [2005; 2016].

## 2 SEMs: a review

Formally, a *structural-equations model*  $M$  is a pair  $(\mathcal{S}, \mathcal{F})$ , where  $\mathcal{S}$  is a *signature*, which explicitly lists the endogenous and exogenous variables and characterizes their possible values, and  $\mathcal{F}$  defines a set of *modifiable structural equations*, relating the values of the variables. We extend the signature to include a set of *allowed interventions*, as was done in earlier work [Beckers and Halpern, 2019; Rubenstein et al., 2017]. Intuitively, allowed interventions are the ones that are feasible or meaningful. A signature  $\mathcal{S}$  is a tuple  $(\mathcal{U}, \mathcal{V}, \mathcal{R}, \mathcal{I})$ .  $\mathcal{U}$  is a set of exogenous variables,  $\mathcal{V}$  is a set of endogenous variables, and  $\mathcal{R}$  associates with every variable  $Y \in \mathcal{U} \cup \mathcal{V}$  a nonempty, finite set  $\mathcal{R}(Y)$  of possible values for  $Y$  (i.e., the set of values over which  $Y$  ranges). We assume (as is typical for SEMs) that  $\mathcal{U}$  and  $\mathcal{V}$  are finite sets, and adopt the convention that for  $\vec{Y} \subseteq \mathcal{U} \cup \mathcal{V}$ ,  $\mathcal{R}(\vec{Y})$  denotes the product of the ranges of the variables appearing in  $\vec{Y}$ ; that is,  $\mathcal{R}(\vec{Y}) := \times_{Y \in \vec{Y}} \mathcal{R}(Y)$ . Finally, an intervention  $I \in \mathcal{I}$  is a set of pairs  $(X, x)$ , where  $X \in \mathcal{V}$  and  $x \in \mathcal{R}(X)$ . For each  $X \in \mathcal{V}$ , there is at most one  $x \in \mathcal{R}(X)$  with  $(X, x) \in I$ . We abbreviate an intervention  $I$  by  $\vec{X} \leftarrow \vec{x}$ , where  $\vec{X} \subseteq \mathcal{V}$  and, unless  $\vec{X}$  is empty,  $\vec{x} \in \mathcal{R}(\vec{X})$ . Although this notation makes most sense if  $\vec{X}$  is nonempty, we allow  $\vec{X}$  to be empty (which

amounts to not intervening at all). If  $I$  consists of exactly one pair  $(Y, y)$ , we abbreviate  $I$  as  $Y \leftarrow y$ .

$\mathcal{F}$  associates with each endogenous variable  $X \in \mathcal{V}$  a function denoted  $F_X$  such that  $F_X : \mathcal{R}(\mathcal{U} \cup \mathcal{V} - \{X\}) \rightarrow \mathcal{R}(X)$ . This mathematical notation just makes precise the fact that  $F_X$  determines the value of  $X$ , given the values of all the other variables in  $\mathcal{U} \cup \mathcal{V}$ . If there is one exogenous variable  $U$  and three endogenous variables,  $X, Y$ , and  $Z$ , then  $F_X$  defines the values of  $X$  in terms of the values of  $Y, Z$ , and  $U$ . For example, we might have  $F_X(u, y, z) = u + y$ , which is usually written as  $X = U + Y$ . Thus, if  $Y = 3$  and  $U = 2$ , then  $X = 5$ , regardless of how  $Z$  is set.

The structural equations define what happens in the presence of external interventions. Setting the value of some variable  $X$  to  $x$  in a SEM  $M = (\mathcal{S}, \mathcal{F})$  results in a new SEM, denoted  $M_{\vec{X} \leftarrow \vec{x}}$ , which is identical to  $M$ , except that the equation for  $X$  in  $\mathcal{F}$  is replaced by  $X = x$ . Interventions on subsets  $\vec{X}$  of  $\mathcal{V}$  are defined similarly. Notice that  $M_{\vec{X} \leftarrow \vec{x}}$  is always well defined, even if  $(\vec{X} \leftarrow \vec{x}) \notin \mathcal{I}$ . In earlier work, the reason that the model included allowed interventions was that, for example, relationships between two models were required to hold only for allowed interventions (i.e., the interventions that were meaningful). As we shall see, here, the fact that we do not have to specify what happens for certain interventions has a more significant impact.

Given a context  $\mathbf{u} \in \mathcal{R}(\mathcal{U})$ , the *outcomes* of a SEM  $M$  under intervention  $\vec{X} \leftarrow \vec{x}$  are all assignments of values  $\mathbf{v} \in \mathcal{R}(\mathcal{V})$  such that the assignments  $\mathbf{u}$  and  $\mathbf{v}$  together satisfy the structural equations of  $M_{\vec{X} \leftarrow \vec{x}}$ . This set of outcomes is denoted  $M(\mathbf{u}, \vec{X} \leftarrow \vec{x})$ . Given an outcome  $\mathbf{v}$ , we denote by  $\mathbf{v}[X]$  and  $\mathbf{v}[\vec{X}]$  the value that  $\mathbf{v}$  assigns to  $X$  and the restriction of  $\mathbf{v}$  to  $\mathcal{R}(\vec{X})$  respectively. We also use this notation for interventions; for example,  $\vec{y}[X]$  is the value that intervention  $\vec{Y} \leftarrow \vec{y}$  assigns to variable  $X \in \vec{Y}$ .

As discussed in the introduction, an important special case of SEMs are acyclic (or recursive) SEMs. Formally, an acyclic SEM is one for which, for every context  $\mathbf{u} \in \mathcal{R}(\mathcal{U})$ , there is some total ordering  $\prec_{\mathbf{u}}$  of the endogenous variables (the ones in  $\mathcal{V}$ ) such that if  $X \prec_{\mathbf{u}} Y$ , then  $X$  is independent of  $Y$ , that is,  $F_X(\mathbf{u}, \dots, y, \dots) = F_X(\mathbf{u}, \dots, y', \dots)$  for all  $y, y' \in \mathcal{R}(Y)$ . Intuitively, if a theory is acyclic, there is no feedback. Acyclic models always have unique outcomes; this is a consequence of assuming that  $\mathcal{V}$  is finite.

In order to talk about SEMs and the information they represent more precisely, we use the formal language  $\mathcal{L}(\mathcal{S})$  for SEMs having signature  $\mathcal{S}$ , introduced by Halpern [2000]; see also [Galles and Pearl, 1998]. An informal description of this language follows; for more details, see [Halpern, 2000]. We restrict the language used by Halpern [2000] to formulas that mention only allowed interventions. Fix a signature  $\mathcal{S} = (\mathcal{U}, \mathcal{V}, \mathcal{R}, \mathcal{I})$ . Given an assignment  $\mathbf{v} \in \mathcal{R}(\mathcal{V})$ , the *primitive event*  $X = x$  is true of  $\mathbf{v}$ , written  $\mathbf{v} \models (X = x)$ , if  $\mathbf{v}[X] = x$ ; otherwise it is false. We extend this definition to *events*  $\varphi$ , which are Boolean combinations of primitive events, in the obvious way. Given a SEM  $M$  with signature  $\mathcal{S}$  and an allowed intervention  $\vec{Y} \leftarrow \vec{y} \in \mathcal{I}$ , the *atomic causal formula*  $[\vec{Y} \leftarrow \vec{y}]\varphi$  is true in context  $\mathbf{u}$ , written

$(M, \mathbf{u}) \models [\vec{Y} \leftarrow \vec{y}] \varphi$  if, for all outcomes  $\mathbf{v} \in M(\mathbf{u}, \vec{Y} \leftarrow \vec{y})$ , we have  $\mathbf{v} \models \varphi$ . Again, we extend this definition to *causal formulas*, which are Boolean combinations of atomic formulas, in the obvious way. The language  $\mathcal{L}(\mathcal{S})$  consists of all causal formulas (over  $\mathcal{S}$ ). Using these formulas, we can also talk about properties that only *some* of the outcomes  $\mathbf{v} \in M(\mathbf{u}, \vec{Y} \leftarrow \vec{y})$  have. For an event  $\varphi$ , define  $\langle \vec{Y} \leftarrow \vec{y} \rangle \varphi$  as  $\neg[\vec{Y} \leftarrow \vec{y}](\neg\varphi)$ . This formula is true exactly when  $\varphi$  is true of at least one outcome  $\mathbf{v} \in M(\mathbf{u}, \vec{Y} \leftarrow \vec{y})$ .

The language of causal formulas completely characterizes the outcomes of a causal model with finite outcome sets, in the following precise sense. (For the purposes of this paper, a *causal model* is either a SEM, a GSEM, or a CCM.)

**Theorem 2.1:** *If  $M$  and  $M'$  are causal models over the same signature  $\mathcal{S}$  that, given a context and intervention, return a finite set of outcomes, then  $M$  and  $M'$  have the same outcomes (that is, for all  $\mathbf{u} \in \mathcal{R}(\mathcal{U})$  and  $I \in \mathcal{I}$ ,  $M(\mathbf{u}, I) = M'(\mathbf{u}, I)$ ) if and only if they satisfy the same set of causal formulas (that is, for all  $\mathbf{u} \in \mathcal{R}(\mathcal{U})$  and  $\psi \in \mathcal{L}(\mathcal{S})$ ,  $M, \mathbf{u} \models \psi \Leftrightarrow M', \mathbf{u} \models \psi$ ).*

The proof of this and all other results not in the main text can be found in the appendix.

### 3 Generalized structural-equations models

The main purpose of causal modeling is to reason about a system's behavior under intervention. A SEM can be viewed as a function that takes a context  $\mathbf{u}$  and an intervention  $\vec{Y} \leftarrow \vec{y}$  and returns a set of outcomes, namely, the set of all solutions to the structural equations after replacing the equations for the variables in  $\vec{Y}$  with  $\vec{y}$ .

Viewed in this way, generalized structural-equations models (GSEMs) are a generalization of SEMs. In a GSEM, there is a function  $\mathbf{F}$  that takes a context  $\mathbf{u}$  and an intervention  $\vec{Y} \leftarrow \vec{y}$  and returns a set of outcomes. However, the outcomes need not be determined by solving a set of suitably modified equations as they are for SEMs. This relaxation gives GSEMs the ability to concisely represent dynamical systems and other systems with infinitely many variables, and the flexibility to handle situations involving finitely many variables that cannot be modeled by SEMs.

Because GSEMs don't have the structure that SEMs have by virtue of being defined in terms of structural equations, we may want to rule out certain unintuitive possibilities. In particular, we require that after intervening to set  $\vec{Y} \leftarrow \vec{y}$ , all outcomes satisfy  $\vec{Y} = \vec{y}$ .

#### 3.1 GSEMs and SEMs

Formally, a *generalized structural-equations model* (GSEM)  $M$  is a pair  $(\mathcal{S}, \mathbf{F})$ , where  $\mathcal{S}$  is a signature, and  $\mathbf{F}$  is a mapping from contexts and interventions to sets of outcomes. More precisely, a signature  $\mathcal{S}$  is a quadruple  $(\mathcal{U}, \mathcal{V}, \mathcal{R}, \mathcal{I})$  where, as before,  $\mathcal{U}$  is a set of exogenous variables,  $\mathcal{V}$  is a set of endogenous variables, and  $\mathcal{R}$  associates with every variable  $Y$  in  $\mathcal{U} \cup \mathcal{V}$  a nonempty, finite set  $\mathcal{R}(Y)$  of possible values for  $Y$ ; we extend  $\mathcal{R}$  to subsets of  $\mathcal{V}$  in the same way as before. However, we no longer require that  $\mathcal{U}$ ,  $\mathcal{V}$  or the sets  $\mathcal{R}(Y)$  for  $Y \in \mathcal{U} \cup \mathcal{V}$  be finite. The mapping  $\mathbf{F}$  is a

function  $\mathbf{F} : \mathcal{I} \times \mathcal{R}(\mathcal{U}) \rightarrow \mathcal{P}(\mathcal{R}(\mathcal{V}))$ , where  $\mathcal{P}$  denotes the powerset operation. That is, it maps a context  $\mathbf{u} \in \mathcal{R}(\mathcal{U})$  and an allowed intervention  $I \in \mathcal{I}$  to a set of *outcomes*  $\mathbf{F}(\mathbf{u}, I) \subseteq \mathcal{P}(\mathcal{R}(\mathcal{V}))$ . As with SEMs, we denote these outcomes by  $M(\mathbf{u}, I)$ . As stated above, we require that outcomes  $\mathbf{v} \in \mathbf{F}(\mathbf{u}, \vec{X} \leftarrow \vec{x})$  satisfy  $\mathbf{v}[\vec{X}] = \vec{x}$ . In the special case where all interventions are allowed, we take  $\mathcal{I} = \mathcal{I}_{univ}$ , the set of all interventions.

We now make precise the sense in which GSEMs generalize SEMs. Two causal models  $M$  and  $M'$  are *equivalent*, denoted  $M \equiv M'$ , if they have the same signature and they have the same outcomes, that is, if for all sets of variables  $\vec{X} \subseteq \mathcal{V}$ , all values  $\vec{x} \in \mathcal{R}(\vec{X})$  such that  $\vec{X} \leftarrow \vec{x} \in \mathcal{I}$ , and all contexts  $\mathbf{u} \in \mathcal{R}(\mathcal{U})$ , we have  $M(\mathbf{u}, \vec{X} \leftarrow \vec{x}) = M'(\mathbf{u}, \vec{X} \leftarrow \vec{x})$ .

**Theorem 3.1:** *For all SEMs  $M$ , there is a GSEM  $M'$  such that  $M \equiv M'$ .*

Just as for SEMs, the intervention  $I = \vec{Y} \leftarrow \vec{y}$  on a GSEM  $M$  induces another GSEM  $M_I$ . To define  $M_I$  precisely, we must first define the composition of interventions.

Given interventions  $\vec{X} \leftarrow \vec{x}$  and  $\vec{Y} \leftarrow \vec{y}$ , let their composition  $I = \vec{X} \leftarrow \vec{x}; \vec{Y} \leftarrow \vec{y}$  be the intervention that results by letting the intervention performed second ( $\vec{Y} \leftarrow \vec{y}$ ) override the first on variables that both interventions affect; that is,  $I = \vec{X} \cup \vec{Y} \leftarrow \vec{z}$ , where for  $Z \in \vec{X} \cup \vec{Y}$ ,

$$\vec{z}[Z] = \begin{cases} \vec{y}[Z] & \text{if } Z \in \vec{Y}, \\ \vec{x}[Z] & \text{if } Z \in \vec{X} - \vec{Y}. \end{cases}$$

Given a GSEM  $M = (\mathcal{S}, \mathbf{F}, \mathcal{I})$  and an intervention  $I \in \mathcal{I}$ , define the intervened model  $M_I$  to be  $(\mathcal{S}, \mathbf{F}', \mathcal{J})$ , where  $\mathcal{J} = \{J \in \mathcal{I}_{univ} : I; J \in \mathcal{I}\}$  and, for  $J \in \mathcal{J}$ ,  $\mathbf{F}'(\mathbf{u}, J) = \mathbf{F}(\mathbf{u}, I; J)$ . (The same relationship holds between the signatures  $\mathcal{I}$  of  $M$  and  $\mathcal{J}$  of  $M_I$  when  $M$  is a SEM.) Notice that if the set  $\mathcal{I}$  is closed under composition, that is, if for all  $I, J \in \mathcal{I}$  we have  $I; J \in \mathcal{I}$ , then  $\mathcal{J} = \{J \in \mathcal{I}_{univ} : I; J \in \mathcal{I}\} \supseteq \mathcal{I}$ , so that with  $M_I$  we have all the interventions that we had with  $M$ , and perhaps more.

The skeptical reader may wonder if the mechanism of equation modification in SEMs really is doing the same thing as the mechanism of intervention composition in GSEMs. This is indeed the case. There are two equivalent ways to see this. The first is to show that equation modification and intervention composition are the same for SEMs.

**Theorem 3.2:** *For all SEMs  $M$  and interventions  $I, J \in \mathcal{I}$  such that  $I; J \in \mathcal{I}$ , we have that  $M_I(\mathbf{u}, J) = M(\mathbf{u}, I; J)$ .*

The second is to show that interventions respect equivalences that hold between SEMs and GSEMs.

**Theorem 3.3:** *If  $M$  and  $M'$  are causal models with  $M \equiv M'$ , then for all  $I \in \mathcal{I}$ , we have that  $M_I \equiv M'_I$ .*

#### 3.2 Finite GSEMs

GSEMs clearly differ from SEMs in that the sets of endogenous and exogenous variables and the range of each individual variable can be infinite. Consider the class of GSEMs where these restrictions are retained, which we call *finite GSEMs*. How do finite GSEMs relate to SEMs? Halpern

[2000] showed that all SEMs satisfy an axiom system called  $AX^+$  (see the Appendix for more details). For example, one axiom (effectiveness) states that after setting  $X \leftarrow x$ , all outcomes have  $X = x$ :  $[\tilde{W} \leftarrow \tilde{w}; X \leftarrow x](X = x)$ . While we imposed this constraint explicitly on GSEMs (and hence this axiom is *valid* in GSEMs—it is true in all contexts of all GSEMs), in SEMs there is no need to impose it; it is a property of the way outcomes are calculated. However, there are additional axioms, for example, one that requires unique outcomes if we intervene on all but one endogenous variable, that finite GSEMs do not satisfy. If we impose these axioms on finite GSEMs, we recover SEMs.

**Theorem 3.4:** *For all finite GSEMs over a signature  $\mathcal{S}$  such that  $\mathcal{I} = \mathcal{I}_{univ}$  in which all the axioms of  $AX^+$  are valid, there is an equivalent SEM, and vice versa.*

Likewise, all acyclic SEMs satisfy an axiom system called  $AX_{rec}^+$  (also described in the Appendix), which consists of the axioms in  $AX^+$  along with two additional conditions. Imposing these axioms on finite GSEMs when all interventions are allowed recovers exactly the class of acyclic SEMs.

**Theorem 3.5:** *For all finite GSEMs over a signature  $\mathcal{S}$  such that  $\mathcal{I} = \mathcal{I}_{univ}$  and all the axioms  $AX_{rec}^+$  are valid, there is an equivalent acyclic SEM, and vice versa.*

We remark that the axiom system  $AX^+$  can be generalized so as to deal with arbitrary (not necessarily finite) GSEMs, and soundness and completeness results for GSEMs can be proved; see [Anonymous, 2021r].

Theorems 3.4 and 3.5 show that finite GSEMs satisfying  $AX^+$  and  $AX_{rec}^+$ , respectively, are equivalent to SEMs and acyclic SEMs, respectively, if *all interventions are allowed*. This equivalence breaks down once we restrict the set of interventions; GSEMs are then strictly more expressive than SEMs, as the following example shows.

**Example 3.6:** Suppose that Suzy is playing a shell game with two shells. One of the shells conceals a dollar; the other shell is empty. Suzy can choose to flip over a shell. If she does, the house flips over the other shell. If Suzy picks shell 1, which hides the dollar, she wins the dollar; otherwise she wins nothing. This example can be modeled by a GSEM  $M_{shell}$  with two binary endogenous variables  $S_1, S_2$  describing whether shell 1 is flipped over and shell 2 is flipped over, respectively, and a binary endogenous variable  $Z$  describing the change in Suzy’s winnings. (The GSEM also has a trivial exogenous variable whose range has size 1, so that there is only one context  $\mathbf{u}$ .) That defines  $\mathcal{U}, \mathcal{V}$ , and  $\mathcal{R}$ ; we set  $\mathcal{I} = \{S_1 \leftarrow 1, S_2 \leftarrow 1\}$ ; and  $\mathbf{F}$  is defined as follows:

$$\begin{aligned} \mathbf{F}(\mathbf{u}, S_1 \leftarrow 1) &= M_{shell}(\mathbf{u}, S_1 \leftarrow 1) \\ &= \{(S_1 = 1, S_2 = 1, Z = 1)\} \\ \mathbf{F}(\mathbf{u}, S_2 \leftarrow 1) &= M_{shell}(\mathbf{u}, S_2 \leftarrow 1) \\ &= \{(S_1 = 1, S_2 = 1, Z = 0)\}. \end{aligned}$$

$M_{shell}$  is clearly a valid GSEM. Furthermore, checking that  $M_{shell}$  satisfies all the axioms in  $AX^+$  is straightforward; see the Appendix (Theorem B.1) for details. However, no SEM  $M'$  with the same signature can have the outcomes  $M'(\mathbf{u}, S_1 \leftarrow 1) = \{(S_1 = 1, S_2 = 1, Z = 1)\}$  and

$M'(\mathbf{u}, S_1 \leftarrow 2) = \{(S_1 = 1, S_2 = 1, Z = 0)\}$ . This is because in a SEM, the value of  $Z$  would be specified by a structural equation  $Z = \mathcal{F}_Z(\mathcal{U}, S_1, S_2)$ . This cannot be the case here, since there are two outcomes having  $S_1 = S_2 = 1$ , but with different values of  $Z$ . ■

This example shows that the set of finite GSEMs satisfying the axioms of  $AX^+$  is more expressive than SEMs when not all interventions are allowed. The fundamental issue here is that  $Z$  is determined by the intervention (which shell Suzy picks), not the state of the shells. In SEMs, the system’s behavior cannot depend explicitly on the intervention, only on the variables altered by the intervention.

## 4 Dynamical systems

In this section, we show how GSEMs can be used to model dynamical systems characterized by a system of ordinary differential equations (ODEs). Suppose that we have a system of ODEs of the form

$$\begin{aligned} \dot{\mathcal{X}}_1 &= F_1(\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_n) \\ \dot{\mathcal{X}}_2 &= F_2(\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_n) \\ &\vdots \\ \dot{\mathcal{X}}_n &= F_n(\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_n), \end{aligned}$$

where the  $\mathcal{X}_i$  are real-valued functions of time, called *dynamical variables*, and  $\dot{\mathcal{X}}_i$  denotes the derivative of  $\mathcal{X}_i$  with respect to time. (Nearly all systems of ODEs occurring in practice can be put into this form by adding auxiliary variables [Young and Mohlenkamp, 2017].) For example,  $d^2\mathcal{X}/dt^2 = -\mathcal{X}$  becomes the pair of equations  $d\mathcal{X}/dt = \mathcal{Y}; d\mathcal{Y}/dt = -\mathcal{X}$ .) This system of ODEs, together with the initial values  $\mathcal{X}_1(0), \mathcal{X}_2(0), \dots, \mathcal{X}_n(0)$ , determines a set of solutions over the interval  $[0, T]$  for  $T > 0$  or the interval  $[0, \infty)$ .

We capture this system of ODEs using the GSEM  $M_{ODE} = (\mathcal{S}, \mathbf{F})$ . The signature  $\mathcal{S} = (\mathcal{U}, \mathcal{V}, \mathcal{R}, \mathcal{I})$  of  $M_{ODE}$  is defined as follows:

$\mathcal{V} = \{X_i^s : 1 \leq i \leq n, s \in (0, T]\}$ ,  $\mathcal{U} = \{X_1^0, \dots, X_n^0\}$ ,  $\mathcal{R}(V) = \mathbb{R}$  for  $V \in \mathcal{U} \cup \mathcal{V}$ , and  $\mathcal{I} = \mathcal{I}_{univ}$ . Here  $X_i^s$  represents the value of function  $X_i$  at time  $s$ , that is,  $X_i(s)$ . The only nontrivial part of this definition is the function  $\mathbf{F}$ . To describe  $\mathbf{F}$ , we first need some preliminary definitions. A dynamical variable  $\mathcal{X}_i$  is *intervention-free with respect to*  $I = \vec{X} \leftarrow \vec{x}$  on an interval  $(a, b)$  if for all  $t \in (a, b)$ , we have  $X_i^t \notin \vec{X}$ ; the interval  $(a, b)$  is *intervention-free* if all dynamical variables are intervention-free on  $(a, b)$ . Given a context  $\mathbf{u}$  and an intervention  $I = \vec{X} \leftarrow \vec{x}$ , let  $\mathbf{F}(\mathbf{u}, I)$  consist of all outcomes  $\mathbf{v}$  that satisfy the following conditions. Define the functions  $\mathcal{X}_i$  for  $1 \leq i \leq n$  by taking  $\mathcal{X}_i(t) = \mathbf{v}[X_i^t]$ . We impose the following constraints:

**ODE1.** The outcome  $\mathbf{v}$  agrees with  $I$ , that is,  $\mathbf{v}[\vec{X}] = \vec{x}$ .

**ODE2.** For all  $i$ ,  $\mathcal{X}_i$  is left-continuous except when intervened on; that is,  $\mathcal{X}_i$  is left-continuous at all points  $t$  such that  $X_i^t \notin \vec{X}$ .

**ODE3.** For all intervention-free intervals  $(a, b) \subseteq [0, T]$ , the functions  $\mathcal{X}_1(t), \mathcal{X}_2(t), \dots, \mathcal{X}_n(t)$  solve the initial-value

problem on  $(a, b)$ . That is, for all  $1 \leq i \leq n$ ,  $\mathcal{X}_i$  is right-continuous at  $a$ , differentiable on  $(a, b)$ , and its derivative  $\dot{\mathcal{X}}_i$  satisfies  $\dot{\mathcal{X}}_i(t) = F_i(\mathcal{X}_1(t), \dots, \mathcal{X}_m(t))$  for all  $t \in (a, b)$ .

These conditions require, rather straightforwardly, that the outcomes agree with the differential equations except where the modeler has intervened (and are appropriately continuous). No fancy limits or partial functions are required.

This set of allowed interventions is rather rich. In practical scenarios, we are typically interested in interventions where

- a variable is set to a certain value at a certain instant in time, or
- a variable is set to a certain value and held at that value throughout an interval of time,

as well as finite compositions of these interventions.

Interventions on intervals interact particularly well with the initial-value problems that arise in solving differential equations. For these interventions, it makes sense to demand that outcomes of the GSEM satisfy a stronger version of condition ODE3 above. Suppose that there are two dynamical variables  $\mathcal{X}$  and  $\mathcal{Y}$ . We have  $\dot{\mathcal{X}} = \mathcal{Y}$  and  $\dot{\mathcal{Y}} = \mathcal{X}$ . We are interested in outcomes on the interval  $[0, 1]$ . If we intervene to set  $X_t$  to 0 on the whole interval  $[0, 1]$ , applying ODE3, we would find that any assignment whatsoever to the  $Y_t$  can be extended to a valid outcome (along with  $X_t = 0$ ). This is not very useful for modeling purposes. It is more useful to require that the differential equation for  $Y$  still hold, and remove the differential equation for  $X$ . This is the import of condition ODE3' below. A variable is *set to a constant*  $k$  during an open interval  $(a, b)$  if for all  $t \in (a, b)$ ,  $\vec{x}[X_i^t] = k$ . An open interval  $(a, b)$  is *intervention-constant* if for all  $1 \leq i \leq m$ ,  $\mathcal{X}_i$  is either intervention-free on  $(a, b)$  or set to a constant  $k$  (for some  $k \in \mathbb{R}$ ) during  $(a, b)$ .

**ODE3'.** For every intervention-constant open interval  $(a, b)$ , if  $\mathcal{X}_i$  is intervention-free on  $(a, b)$ , then  $\mathcal{X}_i$  is right-continuous at  $a$ , differentiable on  $(a, b)$ , and its derivative  $\dot{\mathcal{X}}_i$  satisfies  $\dot{\mathcal{X}}_i(t) = F_i(\mathcal{X}_1, \dots, \mathcal{X}_m)(t)$  for all  $t \in (a, b)$ .

Notice that ODE3' implies ODE3: if all variables are intervention-free on  $(a, b)$ , then the outcomes must satisfy all the differential equations on  $(a, b)$ .

We now show how to find the unique outcome in a GSEM  $M_{ODE}$  satisfying ODE1, ODE2, and ODE3' for a large class of interventions of practical interest that we denote  $\mathcal{I}_{intervals}$ .  $\mathcal{I}_{intervals}$  consists of all finite compositions  $I_1; I_2; \dots; I_m$ , where each  $I_j$  is either a point intervention  $X_i^t \leftarrow k$  or an interval intervention  $\{X_i^t \mid t \in (a, b)\} \leftarrow k$ , which we abbreviate as  $X_i(a, b) \leftarrow k$  for readability. (We similarly use the abbreviations  $X_i[a, b)$ ,  $X_i(a, b]$ , and  $X_i[a, b]$ .) Note that intervening on each of these sets can be achieved by composing two or three point or interval interventions.

Fix  $I = I_1; \dots; I_m \in \mathcal{I}_{intervals}$ . Let  $t_1 < t_2 < \dots < t_l$  be the endpoints of the intervals  $I_1, \dots, I_m$ . (The endpoint of  $X_i^t \leftarrow k$  is  $t$  and the endpoints of  $X_i(a, b) \leftarrow k$  are  $a$  and  $b$ .) It is easy to see that each interval  $(t_1, t_2)$  is intervention-constant. Thus, we can find outcomes of the model step by step. The following algorithm finds an outcome of  $M_{ODE}$

under intervention  $I = \vec{X} \leftarrow \vec{x} \in \mathcal{I}_{intervals}$  with initial conditions  $\mathbf{u} = (X_1^0, \dots, X_n^0)$ . Note that we take the ability to solve initial-value problems and store their solutions as primitive. For convenience, we define  $t_0 = 0$ .

#### Algorithm SOLVE-ODE-GSEM

1. For  $1 \leq i \leq n$ , define  $\mathcal{X}_i(0) = X_i^0$ .
2. For  $i = 1, \dots, l$ :
  - (a) For  $j = 1, \dots, n$ , if  $\mathcal{X}_j$  is set to a constant  $k$  on  $(t_{i-1}, t_i)$ , define  $\mathcal{X}_j(t) = k$  for all  $t \in (t_{i-1}, t_i)$ .
  - (b) Define the remaining (intervention-free) dynamical variables on  $(t_{i-1}, t_i)$  so that for  $j = 1, \dots, n$ , if  $\mathcal{X}_j$  is intervention-free, then  $\mathcal{X}_j$  is right-continuous at  $t_{i-1}$ , differentiable on  $(t_{i-1}, t_i)$ , and its derivative  $\dot{\mathcal{X}}_j$  satisfies  $\dot{\mathcal{X}}_j(t) = F_j(\mathcal{X}_1(t), \dots, \mathcal{X}_m(t))$  for all  $t \in (t_{i-1}, t_i)$ . If there is no way to do this, output “No solution”.
  - (c) For  $j = 1, \dots, n$ , define  $\mathcal{X}_j(t_i)$  as follows.
    - (i) If  $X_j^{t_i} \in \vec{X}$ , define  $\mathcal{X}_j(t_i) = \vec{x}[X_j^{t_i}]$ .
    - (ii) If  $X_j^{t_i} \notin \vec{X}$ , define  $\mathcal{X}_j(t_i) = \lim_{t \rightarrow t_i^-} \mathcal{X}_j(t)$ .
3. Define the functions  $\mathcal{X}_i$ ,  $1 \leq i \leq n$ , on  $(t_l, \infty)$  so that  $\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_n$  solve the initial-value problem on  $(t_l, \infty)$ , as defined in ODE3. Again, if there is no way to do this, output “No solution”.
4. Output the outcome  $\mathbf{v}$  defined by  $\mathbf{v}[X_i^t] = \mathcal{X}_i(t)$  for all  $1 \leq i \leq n$ ,  $t > 0$ .

If all initial-value problems arising in steps 2(b) and 3 are uniquely solvable, then SOLVE-ODE-GSEM outputs the unique outcome  $\mathbf{v} \in M_{ODE}(\mathbf{u}, I)$ . In the general case where some initial-value problems have multiple (or no) solutions, SOLVE-ODE-GSEM is under-specified, because it may pick any valid solution in steps 2(b) and 3. In this case, SOLVE-ODE-GSEM outputs all the outcomes (and only the outcomes)  $M_{ODE}(\mathbf{u}, I)$ . In particular, if the model has no outcome for intervention  $I$  given initial conditions  $\mathbf{u}$ , every execution of the algorithm outputs “No solution”.

**Theorem 4.1:** *The set of outcomes output by valid executions of SOLVE-ODE-GSEM are exactly the outcomes  $M_{ODE}(\mathbf{u}, I)$ .*

While GSEMs satisfying ODE3 but not ODE3' don't in general have meaningful outcomes under interval interventions, they do under point interventions; in fact, SOLVE-ODE-GSEM finds the outcomes of such GSEMs under finite compositions of point interventions.

We conclude this section by showing how a textbook dynamical system—an LC circuit—can be modeled as a GSEM. An LC circuit consists of a voltage source, a capacitor, and an inductor. The dynamical variable of interest is the charge  $Q(t)$  on the capacitor; the voltage  $V$ , capacitance  $C$ , and inductance  $L$  are fixed parameters (although we encode them as dynamical variables with zero derivatives). The differential equations governing this circuit's behavior are

$$\begin{aligned}\dot{Q}(t) &= K(t) \\ \dot{K}(t) &= \frac{V}{L} - \frac{1}{LC}Q(t) \\ \dot{V}(t) &= \dot{C}(t) = \dot{L}(t) = 0,\end{aligned}$$

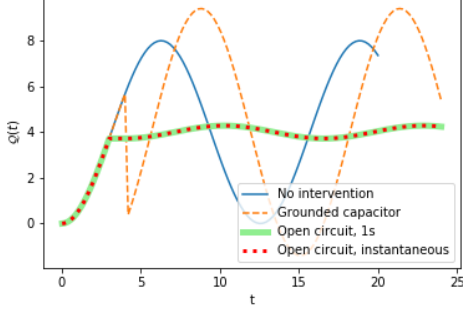


Figure 1: LC circuit outcomes under different interventions.

where  $K$  is the current. The solutions of these differential equations (for  $Q(t)$ ) take the form

$$Q(t) = \mathcal{VC} + A \cos(\omega t + B),$$

where  $\omega = 1/\sqrt{\mathcal{LC}}$ , and  $A$  and  $B$  are determined by the initial conditions on  $Q$  and  $K$ :  $B = \arctan(\frac{K(0)}{\omega(\mathcal{VC} - Q(0))})$  and  $A = \frac{Q(0) - \mathcal{VC}}{\cos(B)}$ . Note that these expressions make sense for all initial conditions except when  $\mathcal{VC} - Q(0) = 0$ ; in this case the solution is instead  $Q(t) = \mathcal{VC}$  for all  $t$ .

It follows from the explicit forms of the solutions given above that all initial-value problems involving these differential equations have unique solutions. It is similarly easy to see that if any of the differential equations (for variable  $\mathcal{X}$ ) is replaced with  $\dot{\mathcal{X}} = 0$ , all initial-value problems involving the resulting modified system of differential equations also have unique solutions. Thus, for all contexts  $\mathbf{u}$  and  $I \in \mathcal{I}_{\text{intervals}}$ , SOLVE-ODE-GSEM returns the unique outcome  $\mathbf{v} \in M(\mathbf{u}, I)$ . Suppose that the initial conditions of the circuit are given by the context  $\mathbf{u} = \{Q_0 = 0, K_0 = 0, V_0 = 2, C_0 = 2, L_0 = 2\}$ . The unique outcome  $\mathbf{v} \in M(\mathbf{u}, \emptyset)$  under the empty intervention is shown in Figure 4 (blue curve).

Now suppose that the capacitor breaks down if  $Q$  exceeds 6. Figure 1 shows that in the absence of intervention,  $Q$  exceeds 6 just after  $t = 4$ . To prevent this, the operators ground the capacitor at time 4 (i.e., make the intervention  $Q_4 \leftarrow 0$ ). However, this does not help. As shown in Figure 4 (orange curve), this initially reduces  $Q$ , but eventually it exceeds 6. Next, the operators try opening the circuit at time 3 and closing it again at time 4 (i.e., performing the intervention  $K(3, 4) \leftarrow 0$ ). This results in the circuit entering an operating regime where  $Q$  is nearly constant (the green curve in Figure 4), and never exceeds 6. We can show that, in agreement with intuition, the first intervention is not a cause of the capacitor not breaking down, and neither is the second (because it is not the minimal intervention required to bring about the outcome). But a sub-intervention  $K_4 \leftarrow 0$  of the second intervention is a cause (corresponding to the red curve in Figure 4). See Appendix C for details.

## 5 Related modeling techniques

We designed GSEMs as an extension of SEMs that can model systems with a continuous-time component. Many other

causal or mechanistic modeling schemes have been proposed for such systems in the literature. In this section, we review three of these schemes: causal constraints models [Blom *et al.*, 2019], hybrid automata [Alur *et al.*, 1992], and counterfactual traces in rule-based models [Laurent *et al.*, 2018], and discuss their relationship to GSEMs. (For space reasons, in this abstract, we defer the discussion of hybrid automata and rule-based models to the appendix.)

### 5.1 Causal constraints models

At a coarse-grained level of modeling, such as when describing equilibrium solutions of dynamical systems, constraints between variables are natural objects of study. In order to describe equilibrium solutions and functional laws (roughly, dependencies that cannot be violated via intervention), [Blom *et al.*, 2019] introduced the notion of a causal constraints model (CCM). These models are composed of a set of constraints, each of which are active only under selected interventions; their outcomes are the solutions of the active constraints.

A CCM  $M$  can be viewed as a pair  $(\mathcal{S}, \mathcal{C})$ , where  $\mathcal{S} = (\mathcal{U}, \mathcal{V}, \mathcal{R}, \mathcal{I})$  is a signature and  $\mathcal{C}$  is a set of causal constraints. Each constraint  $C \in \mathcal{C}$  is a pair  $(f_C, a_C)$ , where  $f_C : \mathcal{R}(\mathcal{U}) \times \mathcal{R}(\mathcal{V}) \rightarrow \{0, 1\}$  and  $a_C \subseteq \mathcal{I}$ . Given a context  $\mathbf{u}$  and an intervention  $\vec{X} \leftarrow \vec{x} \in \mathcal{I}$ , the outcomes  $M(\mathbf{u}, \vec{X} \leftarrow \vec{x})$  are all assignments  $\mathbf{v}$  to the variables of  $\mathcal{V}$  such that, for all  $C \in \mathcal{C}$  with  $I \in a_C$ , we have  $f_C(\mathbf{u}, \mathbf{v}) = 1$  and  $\mathbf{v}[\vec{X}] = \vec{x}$ . Causal constraints models are equivalent to GSEMs.

**Theorem 5.1:** *For all GSEMs, there is an equivalent CCM, and vice versa.*

Of course, if we restrict to CCMs that satisfy the axioms in  $AX^+$  (resp.  $AX_{rec}^+$ ), we can prove an analogue of Theorem 5.1 for GSEMs satisfying  $AX^+$  (resp.  $AX_{rec}^+$ ).

CCMs were designed for characterizing equilibrium solutions of dynamical systems. They seem less well suited for our intended applications, since they do not simplify the process of reasoning from the model specification (i.e., constraints) to solutions. Thus, even though they are equivalent in expressive power, we believe that CCMs and GSEMs will find complementary applications.

## 6 Conclusion

While SEMs are a popular modeling framework in many application areas, they have a restrictive form that makes working with infinitely many variables difficult. This has led to attempts to construct application-specific causal models in the study of ordinary differential equations [Blom *et al.*, 2019] and molecular biology [Laurent *et al.*, 2018]. GSEMs can capture all these applications, while retaining enough of the features of SEMs to allow the definitions of notions like actual cause to carry over without change, and being easy to use in practice. Converting a given dynamical model to a GSEM essentially reduces to setting up allowed interventions, as we demonstrate in examples of ordinary differential equation models, rule-based models, and hybrid systems. Determining actual causes is also typically straightforward. GSEMs are extremely flexible, and can be restricted to form classes paralleling SEMs and acyclic SEMs. Appropriately restricted

finite GSEMs are equivalent to finite SEMs when all interventions are allowed, but have greater expressive power when some interventions are not allowed. This permits previously inaccessible modeling tasks, for example, modeling quantum systems in terms of observables.

## References

- [Alur *et al.*, 1992] R. Alur, C. Courcoubetis, T. A. Henzinger, and P.-H. Ho. Hybrid automata: An algorithmic approach to the specification and verification of hybrid systems. In *Hybrid Systems*, pages 209–229. Springer, 1992.
- [Anonymous, 2021r] Anonymous. Reasoning about causal models with infinitely many variables. Unpublished manuscript, also submitted to IJCAI21, 2021r.
- [Beckers and Halpern, 2019] S. Beckers and J. Y. Halpern. Abstracting causal models. In *Proc. Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)*, 2019. The full version appears at [arxiv.org/abs/1812.03789](https://arxiv.org/abs/1812.03789).
- [Beckers and Vennekens, 2018] S. Beckers and J. Vennekens. A principled approach to defining actual causation. *Synthese*, 195(2):835–862, 2018.
- [Blom *et al.*, 2019] T. Blom, S. Bongers, and J. M. Mooij. Beyond structural causal models: causal constraints models. In *Proc. 35th Conference on Uncertainty in Artificial Intelligence (UAI 2019)*, 2019.
- [Galles and Pearl, 1998] D. Galles and J. Pearl. An axiomatic characterization of causal counterfactuals. *Foundation of Science*, 3(1):151–182, 1998.
- [Glymour and Wimberly, 2007] C. Glymour and F. Wimberly. Actual causes and thought experiments. In J. Campbell, M. O’Rourke, and H. Silverstein, editors, *Causation and Explanation*, pages 43–67. MIT Press, Cambridge, MA, 2007.
- [Halpern and Pearl, 2005] J. Y. Halpern and J. Pearl. Causes and explanations: a structural-model approach. Part I: Causes. *British Journal for Philosophy of Science*, 56(4):843–887, 2005.
- [Halpern, 2000] J. Y. Halpern. Axiomatizing causal reasoning. *Journal of A.I. Research*, 12:317–337, 2000.
- [Halpern, 2015] J. Y. Halpern. A modification of the Halpern-Pearl definition of causality. In *Proc. 24th International Joint Conference on Artificial Intelligence (IJCAI 2015)*, pages 3022–3033, 2015.
- [Halpern, 2016] J. Y. Halpern. *Actual Causality*. MIT Press, Cambridge, MA, 2016.
- [Henzinger, 2000] T. A. Henzinger. The theory of hybrid automata. In *Verification of Digital and Hybrid Systems*, pages 265–292. Springer, 2000.
- [Hitchcock, 2001] C. Hitchcock. The intransitivity of causation revealed in equations and graphs. *Journal of Philosophy*, XCVIII(6):273–299, 2001.
- [Hitchcock, 2007] C. Hitchcock. Prevention, preemption, and the principle of sufficient reason. *Philosophical Review*, 116:495–532, 2007.
- [Laurent *et al.*, 2018] J. Laurent, J. Yang, and W. Fontana. Counterfactual resimulation for causal analysis of rule-based models. In *Proc. Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI ’18)*, pages 1882–1890, 2018.
- [Rubenstein *et al.*, 2017] P. K. Rubenstein, S. Weichwald, S. Bongers, J. M. Mooij, D. Janzing, M. Grosse-Wentrup, and B. Schölkopf. Causal consistency of structural equation models. In *Proc. 33rd Conference on Uncertainty in Artificial Intelligence (UAI 2017)*, 2017.
- [Weslake, 2015] B. Weslake. A partial theory of actual causation. *British Journal for the Philosophy of Science*, 2015. To appear.
- [Woodward, 2003] J. Woodward. *Making Things Happen: A Theory of Causal Explanation*. Oxford University Press, Oxford, U.K., 2003.
- [Young and Mohlenkamp, 2017] T. Young and M. J. Mohlenkamp. *Introduction to Numerical Methods and MATLAB Programming for Engineers*. Ohio University, 2017. Available at <https://openlibra.com/en/book/introduction-to-numerical-methods-and-matlab-programming-for-engineers>.

## A An Axiom System for Causal Reasoning

We now review the axiom systems considered by Halpern [2000] for reasoning about causality. Note that there are two slight differences between our presentation and that of Halpern. First, as we mentioned earlier, we have weakened the language of causal formulas so that primitive causal formulas are no longer parameterized by contexts. Thus, our language has formulas such as  $[\vec{Y} \leftarrow \vec{y}](X = x)$  rather than  $[\vec{Y} \leftarrow \vec{y}](X(\mathbf{u}) = x)$ . Second, the list of axioms given below does not include two of Halpern’s axioms, which he called D10 and D11. D11 is a technical axiom that was needed only to reason about formulas with contexts (to reduce to formulas that mentioned only one context); D10 says that there are unique outcomes, and is redundant in acyclic systems. A minor modification of Halpern’s proof shows that the axiom systems  $AX^+$  and  $AX_{rec}^+$  defined below (which are identical to the system Halpern called  $AX^+$  and  $AX_{rec}^+$ , respectively, except that they omit the axioms D10 and D11) are sound and complete for SEMs and acyclic SEMs, respectively, with respect to the language that we are considering (just as Halpern’s versions of  $AX^+$  and  $AX_{rec}^+$  were sound and complete for his language); the proof is essentially identical to Halpern’s, so we omit it here. To axiomatize acyclic SEMs, following Halpern, we define  $Y \rightsquigarrow Z$ , read “ $Y$  affects  $Z$ ”, as an abbreviation for the formula

$$\forall \vec{X} \subseteq \mathcal{V}, \vec{x} \in \mathcal{R}(\vec{X}), y \in \mathcal{R}(y), z \neq z' \in \mathcal{R}(Z) \\ ([\vec{X} \leftarrow \vec{x}](Z = z) \wedge [\vec{X} \leftarrow \vec{x}, Y \leftarrow y](Z = z'));$$

that is,  $Y$  affects  $Z$  if there is some setting of some endogenous variables  $\vec{X}$  for which changing the value of  $Y$  changes the value of  $Z$ . This definition is used in axiom D6 below, which characterizes acyclicity.

**Definition A.1:**  $AX^+$  consists of axiom schema D0-D5 and D7-D9, and inference rule MP.  $AX_{rec}^+$  results from adding D6 to  $AX^+$ .

- D0. All instances of propositional tautologies.
- D1.  $[\vec{Y} \leftarrow \vec{y}](X = x \Rightarrow X \neq x')$  if  $x, x' \in \mathcal{R}(X)$ ,  $x \neq x'$  (functionality)
- D2.  $[\vec{Y} \leftarrow \vec{y}](x \in \mathcal{R}(X))$  (definiteness)
- D3.  $\langle \vec{X} \leftarrow \vec{x} \rangle (W = w \wedge \vec{Y} = \vec{y}) \Rightarrow \langle \vec{X} \leftarrow \vec{x}; W \leftarrow w \rangle (\vec{Y} = \vec{y})$  (composition)
- D4.  $[\vec{W} \leftarrow \mathbf{w}; X \leftarrow x](X = x)$  (effectiveness)
- D5.  $(\langle \vec{X} \leftarrow \vec{x}; Y \leftarrow y \rangle (W = w \wedge \vec{Z} = \vec{z}) \wedge \langle \vec{X} \leftarrow \vec{x}; W \leftarrow w \rangle (Y = y \wedge \vec{Z} = \vec{z})) \Rightarrow \langle \vec{X} \leftarrow \vec{x} \rangle (W = w \wedge Y = y \wedge \vec{Z} = \vec{z})$ , where  $\vec{Z} = \mathcal{V} - (\vec{X} \cup \{W, Y\})$  (reversibility)
- D6.  $(X_0 \rightsquigarrow X_1 \wedge \dots \wedge X_{k-1} \rightsquigarrow X_k) \Rightarrow \neg(X_k \rightsquigarrow X_0)$  (recursiveness)
- D7.  $([\vec{X} \leftarrow \vec{x}]\varphi \wedge [\vec{X} \leftarrow \vec{x}](\varphi \Rightarrow \psi)) \Rightarrow [\vec{X} \leftarrow \vec{x}]\psi$  (distribution)
- D8.  $[\vec{X} \leftarrow \vec{x}]\varphi$  if  $\varphi$  is a propositional tautology (generalization)

D9.  $\langle \vec{Y} \leftarrow \vec{y} \rangle true \wedge (\langle \vec{Y} \leftarrow \vec{y} \rangle (X = x) \Rightarrow \langle Y \leftarrow y \rangle (X \neq x'))$ , if  $x \neq x'$ . (unique outcomes for  $\mathcal{V} - \{X\}$ )

MP. From  $\varphi$  and  $\varphi \Rightarrow \psi$ , infer  $\psi$  (modus ponens)

## B Proofs

**Theorem B.1:**  $M_{shell}$  satisfies all the axioms in  $AX^+$ .

**Proof:** D0, D1, D2, D7 and D8 are trivial. No joint interventions are allowed, so the only way to instantiate D3 is to have  $\vec{X} = W = S_1$  (or symmetrically,  $\vec{X} = W = S_2$ ). But if  $\vec{X} = W$ , then  $\vec{X} \leftarrow \vec{x}; W \leftarrow w = W \leftarrow w$  and D3 follows trivially by eliminating the conjunction. D4 (effectiveness) holds by inspection. D5 holds for the same reason as D3. Finally, D9 cannot be instantiated because no complete interventions are allowed. ■

The following theorem is needed to prove Theorem 3.4.

**Theorem B.2:** If  $M$  and  $M'$  are causal models (either SEMs or GSEMs) with a common signature  $\mathcal{S} = (\mathcal{U}, \mathcal{V}, \mathcal{R}, \mathcal{I}_{univ})$ , where  $\mathcal{V}$  is finite and  $\mathcal{R}(X)$  is finite for all  $X \in \mathcal{V}$ , that both satisfy the axioms in  $AX^+$  and have the same outcomes under complete interventions—that is, for all  $\mathbf{u} \in \mathcal{R}(\mathcal{U})$  and  $X \in \mathcal{V}$ , if  $\vec{Y} = \mathcal{V} \setminus X$ , then for all  $\vec{y} \in \mathcal{R}(\vec{Y})$ ,  $M(\mathbf{u}, \vec{Y} \leftarrow \vec{y}) = M'(\mathbf{u}, \vec{Y} \leftarrow \vec{y})$ —then  $M$  and  $M'$  agree on all causal formulas.

**Proof:** Fix an arbitrary context  $\mathbf{u}$ .  $M$  satisfies axiom D9, so for every variable  $X \in \mathcal{V}$ , and for every assignment  $\vec{y} \in \mathcal{R}(\vec{Y})$  to the variables  $\vec{Y} = \mathcal{V} \setminus \{X\}$ , there is a unique  $x \in \mathcal{R}(X)$  such that  $M, \mathbf{u} \models [\vec{Y} \leftarrow \vec{y}](X = x)$ . Using this fact, we can define a SEM  $M''$  with signature  $\mathcal{S}$  as follows. Define  $\mathcal{F}_X''(\mathbf{u}, \vec{y})$  to be the unique  $x$  such that  $M, \mathbf{u} \models [\mathcal{V} \setminus \{X\} \leftarrow \vec{y}](X = x)$ . Let  $C$  be the set of all formulas  $\varphi = [\mathcal{V} \setminus \{X\} \leftarrow \vec{y}](X = x)$  such that  $M, \mathbf{u} \models \varphi$ . By assumption,  $C$  is also the set of all such formulas  $\varphi$  for which  $M', \mathbf{u} \models \varphi$ . Let  $\chi$  be the conjunction of all the formulas in  $C$ . Since there are finitely many variables, and all ranges are finite, this set of formulas is finite, and so taking the conjunction makes sense. We know that  $M$  and  $M'$  satisfy all axioms of  $AX^+$ , and both models satisfy  $\chi$ . This means that if  $\chi \Rightarrow \psi$  is provable in  $AX^+$ , then  $M$  and  $M'$  both satisfy  $\psi$ . We now show that, for all formulas  $\psi$ , either  $\chi \Rightarrow \psi$  or  $\chi \Rightarrow \neg\psi$  is provable in  $AX^+$ . This means that either both  $M, \mathbf{u} \models \psi$  and  $M', \mathbf{u} \models \psi$  (if  $\chi \Rightarrow \psi$  is provable in  $AX^+$ ), or both  $M, \mathbf{u} \not\models \psi$  and  $M', \mathbf{u} \not\models \psi$  (if  $\chi \Rightarrow \neg\psi$  is provable in  $AX^+$ ). That is,  $M$  and  $M'$  agree on all causal formulas. Note that  $\chi$  is false in all SEMs over  $\mathcal{S}$  other than models that agree with the  $M''$  that we defined using  $\chi$  in context  $\mathbf{u}$ . Thus, if  $(M'', \mathbf{u}) \models \psi$ , then  $\chi \Rightarrow \psi$  is valid; and if  $(M'', \mathbf{u}) \models \neg\psi$ , then  $\chi \Rightarrow \neg\psi$  is valid. Since  $AX^+$  is a sound and complete axiomatization, it follows that either  $\chi \Rightarrow \psi$  or  $\chi \Rightarrow \neg\psi$  is provable, as desired. ■

**Theorem 2.1:** If  $M$  and  $M'$  are causal models over the same signature  $\mathcal{S}$  that, given a context and intervention, return a finite set of outcomes, then  $M$  and  $M'$  have the same outcomes (that is, for all  $\mathbf{u} \in \mathcal{R}(\mathcal{U})$  and  $I \in \mathcal{I}$ ,  $M(\mathbf{u}, I) = M'(\mathbf{u}, I)$ ) if and only if they satisfy the same set of causal formulas (that is, for all  $\mathbf{u} \in \mathcal{R}(\mathcal{U})$  and  $\psi \in \mathcal{L}(\mathcal{S})$ ,  $M, \mathbf{u} \models \psi \Leftrightarrow M', \mathbf{u} \models \psi$ ).



**Proof:** Let  $M$  and  $M'$  be causal models with the same set of solutions. It suffices to consider the primitive causal formulas  $[\vec{Y} \leftarrow \vec{y}](X = x)$ , since the truth of other formulas in  $\mathcal{L}(\mathcal{S})$  are derived from these. Recall that  $M, \mathbf{u} \models [\vec{Y} \leftarrow \vec{y}](X = x)$  iff for all outcomes  $\mathbf{v} \in M(\mathbf{u}, \vec{Y} \leftarrow \vec{y})$ ,  $\mathbf{v}[X] = x$ . But  $M(\mathbf{u}, \vec{Y} \leftarrow \vec{y}) = M'(\mathbf{u}, \vec{Y} \leftarrow \vec{y})$ , so  $M, \mathbf{u} \models [\vec{Y} \leftarrow \vec{y}](X = x)$  if and only if  $M', \mathbf{u} \models [\vec{Y} \leftarrow \vec{y}](X = x)$ . Conversely, suppose that  $M$  and  $M'$  satisfy the same set of causal formulas. Suppose for contradiction that there exists some  $\mathbf{u}$  and  $I$  with  $M(\mathbf{u}, I) \neq M'(\mathbf{u}, I)$ . Then without loss of generality, there is an outcome  $\mathbf{v}$  in  $M(\mathbf{u}, I)$  that is not in  $M'(\mathbf{u}, I)$ . This outcome must differ from each of the finitely many outcomes  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} = M'(\mathbf{u}, I)$  in at least one variable; that is, there must be variables  $X_1, X_2, \dots, X_n \in \mathcal{V}$  with  $\mathbf{v}[X_i] \neq \mathbf{v}_i[X_i]$  for  $i \in 1, 2, \dots, n$ . Consider the causal formula  $\varphi = \langle I \rangle (\bigwedge_{1 \leq i \leq n} X_i = \mathbf{v}[X_i])$ . We have that  $M, \mathbf{u} \models \varphi$ , since  $\mathbf{v} \in M(\mathbf{u}, I)$ . However, it is not true that  $M', \mathbf{u} \models \varphi$ , because no outcome  $\mathbf{v}_i$  of  $M'$  satisfies  $\bigwedge_{1 \leq i \leq n} X_i = \mathbf{v}[X_i]$ . This contradicts the assumption that  $M$  and  $M'$  satisfy the same set of causal formulas; hence  $M(\mathbf{u}, I) = M'(\mathbf{u}, I)$  for all  $\mathbf{u} \in \mathcal{R}(U)$ ,  $I \in \mathcal{I}$ . ■

**Theorem 3.1:** *For all SEMs  $M$ , there is a GSEM  $M'$  such that  $M \equiv M'$ .*

**Proof:** Given a SEM  $M = ((\mathcal{U}, \mathcal{V}, \mathcal{R}, \mathcal{I}), \mathcal{F})$ , define  $M' = ((\mathcal{U}, \mathcal{V}, \mathcal{R}, \mathcal{I}), \mathcal{F}')$ , where for all  $\vec{X} \leftarrow \vec{x} \in \mathcal{I}$  and  $\mathbf{u} \in \mathcal{U}$ ,  $\mathcal{F}'(\mathbf{u}, \vec{X} \leftarrow \vec{x}) = M(\mathbf{u}, \vec{X} \leftarrow \vec{x})$ . Since  $M'(\mathbf{u}, \vec{X} \leftarrow \vec{x}) = \mathcal{F}'(\mathbf{u}, \vec{X} \leftarrow \vec{x})$ ,  $M'$  is equivalent to  $M$  by definition. ■

**Theorem 3.2:** *For all SEMs  $M$  and interventions  $I, J \in \mathcal{I}$  such that  $I; J \in \mathcal{I}$ , we have that  $M_I(\mathbf{u}, J) = M(\mathbf{u}, I; J)$ .*

**Proof:** We prove the equivalent statement  $(M_I)_J(\mathbf{u}) = M_{I;J}(\mathbf{u})$ . Since the outcomes of SEMs are determined by the structural equations, it suffices to show that the structural equations of  $(M_I)_J$  are the same as those of  $M_{I;J}$ . Let  $X \in \mathcal{V}$  be arbitrary and consider the structural equation  $\mathcal{F}_X$ . Without loss of generality, let  $I = \vec{Y} \leftarrow \vec{y}$  and  $J = \vec{Z} \leftarrow \vec{z}$ . There are three cases to consider:  $X \notin \vec{Y} \cup \vec{Z}$ ,  $X \in \vec{Z}$ , and  $X \in \vec{Y} \setminus \vec{Z}$ . The first case is trivial;  $\mathcal{F}_X$  is unmodified in both models. In the second case, letting  $\mathbf{s} \in \mathcal{R}(\mathcal{V} \setminus \{X\})$  denote an arbitrary input to  $\mathcal{F}_X$ , in  $(M_I)_J$ , we have that  $\mathcal{F}_X(\mathbf{s}) = \vec{Z}[X]$ . But by the definition of  $I; J$ , we also have  $\mathcal{F}_X(\mathbf{s}) = \vec{Z}[X]$  in  $M_{I;J}$ . In the third case, in  $(M_I)_J$ ,  $\mathcal{F}_X(\mathbf{s}) = \vec{Y}[X]$ , since in  $M_I$ ,  $\mathcal{F}_X(\mathbf{s}) = \vec{Y}[X]$ , and applying the intervention  $J$  does not affect  $\mathcal{F}_X$  since  $X \notin \vec{Z}$ . ■

**Theorem 3.3:** *Suppose that  $M$  and  $M'$  are causal models with  $M \equiv M'$ . Then for all  $I \in \mathcal{I}$ , we have that  $M_I \equiv M'_I$ .*

**Proof:** Clearly  $M_I$  and  $M'_I$  have the same signatures. It remains to show that for all contexts  $\mathbf{u}$  and all intervention  $J$  allowed in  $M_I$ , we have that  $M_I(\mathbf{u}, J) = M'_I(\mathbf{u}, J)$ . Applying the definition and the fact that  $M \equiv M'$ , we have that  $M'_I(\mathbf{u}, J) = M'(\mathbf{u}, I; J) = M(\mathbf{u}, I; J)$ . Therefore, it suffices to show  $M(\mathbf{u}, I; J) = M_I(\mathbf{u}, J)$ , which is exactly Theorem 3.2. ■

**Theorem 3.4:** *For all finite GSEMs over a signature  $\mathcal{S}$  such that  $\mathcal{I} = \mathcal{I}_{univ}$  and all the axioms of  $AX^+$  are valid, there is an equivalent SEM, and vice versa.*

**Proof:** Given a SEM  $M$ , define a GSEM  $M'$  with the same signature via  $\mathbf{F}'(\mathbf{u}, I) = M(\mathbf{u}, I)$ , as in Theorem 3.1. This GSEM is clearly equivalent to  $M$ . Furthermore, all the axioms in  $AX^+$  are valid in  $M$ . This follows from the facts that (1) equivalent causal models have the same outcomes (by definition), (2) finite causal models with the same outcomes satisfy the same causal formulas (Theorem 2.1), and (3)  $M$  is a SEM, so all the axioms in  $AX^+$  are valid in  $M$ . Conversely, given a finite GSEM  $M'$  in which all the axioms of  $AX^+$  are valid, the GSEM must have unique solutions for  $\mathcal{V} \setminus X$  (D9). That is, for each context  $\mathbf{u} \in \mathcal{R}(U)$  and each variable  $X \in \mathcal{V}$ , if we define  $\vec{Y} = \mathcal{V} \setminus X$ , for every  $\vec{y} \in \mathcal{R}(Y)$ , there is a unique  $x \in \mathcal{R}(X)$  such that  $M', \mathbf{u} \models [\vec{Y} \leftarrow \vec{y}](X = x)$ . Here we use the fact that  $\mathcal{I} = \mathcal{I}_{univ}$  to ensure that the relevant instances of D9 are in the language. We can use this property to define the structural equations of the SEM  $M$ . That is, define a SEM  $M$  with the same signature by defining  $\mathcal{F}_X(\mathbf{u}, \vec{y}) = x$ , where  $x$  is the value guaranteed above. We must show that  $M$  has the same outcomes as  $M'$ . But this is just Theorem B.2. ■

**Theorem 3.5:** *For all finite GSEMs over a signature  $\mathcal{S}$  such that  $\mathcal{I} = \mathcal{I}_{univ}$  and all the axioms  $AX_{rec}^+$  are valid, there is an equivalent acyclic SEM, and vice versa.*

**Proof:** Given a finite GSEM  $M'$  satisfying  $AX_{rec}^+$ , Theorem 3.4 guarantees the existence of an equivalent SEM  $M$ . Since  $M$  is equivalent to  $M'$ ,  $M$  satisfies  $AX_{rec}^+$ . This implies that  $M$  is acyclic. To prove this, suppose not. Then there is  $k > 1$  and endogenous variables  $V_1, \dots, V_k$  having cyclic dependencies; that is,  $V_{i+1}$  is not independent of  $V_i$  for  $i = 1, \dots, k-1$ , and  $V_1$  is not independent of  $V_k$ . But it is easy to see that if  $Y$  is not independent of  $X$ , then  $X$  affects  $Y$ , i.e.,  $X \rightsquigarrow Y$ . Thus,  $V_k \rightsquigarrow V_1 \wedge V_1 \rightsquigarrow V_2 \wedge \dots \wedge V_{k-1} \rightsquigarrow V_k$ . This is the negation of an instance of D6. Hence not all the axioms of  $AX_{rec}^+$  are valid in  $M$ , a contradiction. For the converse, given an acyclic SEM  $M$ , Theorem 3.1 guarantees the existence of an equivalent GSEM  $M'$ . This equivalent GSEM satisfies the same formulas as  $M$ , so it satisfies  $AX_{rec}^+$ . ■

**Theorem 4.1:** *The set of outcomes output by valid executions of SOLVE-ODE-GSEM are exactly the outcomes  $M_{ODE}(\mathbf{u}, I)$ .*

**Proof:** We walk through the algorithm's execution and show that whenever it defines a dynamical variable (and thus a model variable, via the translation in step 4), it can make all the choices compatible with ODE1, ODE2, and ODE3', and cannot make any other choices:

- In step 1,  $\mathcal{X}_i(0) = X_i^0$  is the only choice consistent with the right-continuity requirement of ODE3'.
- In step 2(a),  $\mathcal{X}_j(t) = k$  is the only choice consistent with ODE1. It is compatible with ODE2 and ODE3, since ODE2 and ODE3 require nothing of intervened points.
- In step 2(b) and step 3, the possible settings for the intervention-free variables are exactly the settings al-

lowed by ODE3'. They are compatible with ODE1, since the intervention-free variables are not intervened on in  $(a, b)$ , and compatible with ODE2, since solutions to initial value problems are always continuous.

- In step 2(c)(i),  $\mathcal{X}_j(t_i) = \vec{x}[X_j^{t_i}]$  is the only choice consistent with ODE1. It is compatible with ODE2 and ODE3 since, again, ODE2 and ODE3 require nothing of intervened points.
- In step 2(c)(ii),  $\mathcal{X}_j(t_i) = \lim_{t \rightarrow t_i^-} \mathcal{X}_j(t)$  is the only choice that maintains left-continuity (is consistent with ODE2). It is compatible with ODE1, since  $\mathcal{X}_j$  is not intervened on at time  $t_i$ , and compatible with ODE3, since by construction, there is no intervention-constant open interval containing  $t_i$ . Finally, the limit always exists, because the values of  $\mathcal{X}_j$  on  $(t_{i-1}, t_i)$  were set in step 2(b), so  $\mathcal{X}_j$  is continuous on the open interval  $(t_{i-1}, t_i)$ .

■

**Theorem 5.1:** *For all GSEMs, there is an equivalent CCM, and vice versa.*

**Proof:** Given a GSEM  $M = (\mathcal{S}, \mathbf{F}, \mathcal{I})$ , define a CCM  $M' = (\mathcal{S}, \mathcal{C}, \mathcal{I})$  as follows. For every intervention  $\vec{X} \leftarrow \vec{x} \in \mathcal{I}$ ,  $\mathcal{C}$  contains a constraint  $C$  such that  $A_C = \{\vec{X} \leftarrow \vec{x}\}$ , and for every context  $\mathbf{u}$ ,  $f_C(\mathbf{u}, \mathbf{v}) = 1$  if  $\mathbf{v} \in \mathbf{F}(\mathbf{u}, I)$ , and 0 otherwise. Since the GSEM satisfies effectiveness,  $\mathbf{v} \in \mathbf{F}(\mathbf{u}, I) \Rightarrow \mathbf{v}[\vec{X}] = \vec{x}$ , so the outcomes of the resulting CCM are exactly the outcomes of the GSEM. Conversely, given a CCM  $M' = (\mathcal{S}, \mathcal{C}, I)$ , define a GSEM  $M = (\mathcal{S}, \mathbf{F}, \mathcal{I})$  via  $\mathbf{F}(\mathbf{u}, I) = M'(\mathbf{u}, I)$ ; it is immediate that  $M$  and  $M'$  have the same outcomes. ■

## C Actual causes

One important application of causal modeling is to deducing the *actual cause(s)* of  $X = x$ , that is (roughly speaking) the specific reasons that  $X$  takes value  $x$  in a given context  $\mathbf{u}$  and outcome  $\mathbf{v}$ .

Many definitions of actual cause in SEMs have been proposed (e.g., [Beckers and Vennekens, 2018; Glymour and Wimberly, 2007; Hitchcock, 2001; 2007; Weslake, 2015; Woodward, 2003]). For definiteness, we use that of Halpern and Pearl [2005], as later modified by Halpern [2016], except that we further modify it to deal with allowed interventions (which were not considered in earlier definitions).

**Definition C.1:** [Actual cause] Given a causal model  $M$ ,  $\mathbf{u} \in \mathcal{R}(U)$ , and  $\mathbf{v} \in \mathcal{R}(V)$ ,  $\vec{X} = \vec{x}$  is an *actual cause* of the event  $\varphi$  in  $(M, \mathbf{u}, \mathbf{v})$  if the following three conditions hold:

**AC1.**  $\mathbf{v} \models \vec{X} = \vec{x}$  and  $\mathbf{v} \models \varphi$ .

**AC2.** There is some  $\vec{W} \subseteq V$  and a setting  $\vec{x}'$  of  $\vec{X}$  such that  $\vec{W} \leftarrow \mathbf{v}[\vec{W}]; \vec{X} \leftarrow \vec{x}' \in \mathcal{I}$  and  $(M, \mathbf{u}) \models \langle \vec{W} \leftarrow \mathbf{v}[\vec{W}]; \vec{X} \leftarrow \vec{x}' \rangle \neg \varphi$ .

**AC3.** No proper subset of  $\vec{X}$  satisfies conditions AC1 and AC2.

Intuitively, this definition captures the fact that when reasoning about counterfactuals in a concrete scenario, we often want to fix some details  $\vec{W}$  to the values  $\mathbf{v}[\vec{W}]$  that they actually had in that scenario; see [Halpern, 2015] for more examples and motivation. That paper did not consider allowed interventions. To deal with allowed interventions, we insist that the intervention  $\vec{W} \leftarrow \mathbf{v}[\vec{W}]; \vec{X} \leftarrow \vec{x}'$  appearing in AC2 above is allowed.

Using this formal definition of causality, we can verify the claims made in Section 4 that (1) the intervention  $Q_4 \leftarrow 0$  is not an (actual) cause of the capacitor not breaking down, but (2) the intervention  $K_4 \leftarrow 0$  is. Recall that the capacitor breaks down if  $Q$  exceeds 6, so the capacitor not breaking down corresponds to the statement  $\varphi = \forall t > 4, Q(t) \leq 6$ .<sup>2</sup> Claim 1 is obvious, because under the first intervention,  $Q_t$  does exceed 6 (at, say,  $t = 10$ ). Thus  $\mathbf{v} \models \neg \varphi$ , where  $\mathbf{v}$  is the outcome under  $Q_4 \leftarrow 0$ , which violates AC1. For Claim 2, note that on intervention-free intervals, the solutions are periodic in time, with period  $2\pi/\omega = 2\pi\sqrt{L_o C_o} = 4\pi \approx 12.6$ . Since  $Q$  does not exceed 6 on the interval  $(4, 4 + 12.6)$ , it will never exceed 6. Thus  $\mathbf{v} \models \varphi$ , where  $\mathbf{v}$  is the outcome  $K_4 \leftarrow 0$ , satisfying AC1. AC2 is satisfied, because if we instead set  $K_t \leftarrow K(t)$  for  $t \in (3, 4)$  where  $K(t)$  is the value  $K_t$  had under the empty intervention,  $Q_t$  does exceed 6 (at, say,  $t = 6$ ). (Here we are taking  $\vec{W} = \emptyset$ .) Finally, AC3 is satisfied, because the intervention  $K_4 \leftarrow 0$  is on a single variable, and therefore minimal.

## D More on Related Modeling Techniques

### D.1 Hybrid automata

Hybrid automata are a well-developed class of models for systems that have both continuous and discrete components [Alur *et al.*, 1992]; for example, a thermostat controlling a heater to keep the temperature within a certain tolerance of a set point. The state variables for this system are both continuously varying in time (the temperature) and discretely changing in time (whether the heater is on). In this section, we show how to construct a GSEM model corresponding to an arbitrary hybrid automaton. We demonstrate our construction on a simple example from [Henzinger, 2000] and show how the resulting GSEM can be used to answer causal questions.

Mathematically, a hybrid automaton is a finite directed multigraph  $G = (V, E)$ , a set  $\mathcal{X} = \{\mathcal{X}_1, \dots, \mathcal{X}_n\}$  of real-valued *dynamical variables*,<sup>3</sup> and some predicates (discussed below).<sup>4</sup> The states  $v \in V$  are called *control modes*, and

<sup>2</sup> This statement  $\varphi$ , since it is universally quantified over all  $t > 4$ , cannot be expressed in the language of causal formulas  $\mathcal{L}(\mathcal{S})$  we defined in Section 2. However, we do not view this as a serious problem. The definition of actual causality (Definition C.1) allows arbitrary  $\varphi$ , and it still makes sense if we take  $\varphi$  to be a formula in a richer language containing the first-order quantification we need here.

<sup>3</sup> In the literature, these are usually just called variables; we call them dynamical variables to avoid confusing them with GSEM variables.

<sup>4</sup> There is also a set of events  $\Sigma$  used to disambiguate discrete transitions, but this is not relevant for us.

they represent the state of the discrete component of the system. The edges  $e = (u, v, n) \in E$ , where  $u$  and  $v$  are control modes and  $n$  indexes the edge among all edges from  $u$  to  $v$ , are called *control switches*, and they describe transitions between control modes. Recall that a multigraph is a graph which may have multiple edges between any given pair of nodes; different control switches between the same pair of control modes represent different modes of transition between them. For example, the heater may have multiple triggers that change its state from OFF to ON.

The semantics of a hybrid automaton is defined by predicates *init*, *flow*, *jump*, and *inv*. Possible starting configurations are given by  $\text{init}(v, \mathcal{X})$ . Possible continuous dynamics within a control mode are given by  $\text{flow}(v, \mathcal{X}, \dot{\mathcal{X}})$ , where  $\dot{\mathcal{X}}$  represents the vector of first derivatives. Possibly discontinuous changes are given by  $\text{jump}(e, \mathcal{X}, \mathcal{X}')$ , where  $\mathcal{X}'$  represents values at the conclusion of the change. Hard constraints are represented by  $\text{inv}(v, \mathcal{X})$ , which may be thought of as describing invariants of the different control modes. Hybrid automata are nondeterministic in general; all dynamics compatible with the automaton's predicates are possible.

Let  $A = (V, E, \mathcal{X}, \text{init}, \text{flow}, \text{jump}, \text{inv})$  be a hybrid automaton. We now construct a GSEM  $M$  corresponding to  $A$ . The endogenous variables are as follows. For each  $\mathcal{X}_i \in \mathcal{X}$ ,  $M$  has real-valued variables  $\{X_i^t \mid t \geq 0\}$  corresponding to the value of  $\mathcal{X}_i$  at time  $t$ .  $M$  also has  $V$ -valued variables  $\{S_t \mid t \geq 0\}$  corresponding to the control mode of the system at time  $t$ .  $M$  has a single exogenous variable with a single value, so there is only one (trivial) context.

The outcomes for intervention  $\vec{Y} \leftarrow \vec{y}$  are, similar to ODEs, all assignments to the variables that (1), agree with  $\vec{Y} \leftarrow \vec{y}$ , and (2), are otherwise compatible with the predicates of the automaton. More precisely,  $\mathbf{v} \in M(\mathbf{u}, \vec{Y} \leftarrow \vec{y})$  iff all the following conditions hold. For convenience, we define functions  $\mathcal{X}_i(t) = \mathbf{v}[X_i^t]$  for  $i = 1, \dots, n$ ,  $\mathcal{X}(t) = (\mathcal{X}_1(t), \dots, \mathcal{X}_n(t))$ , and  $\mathcal{D}(t) = \mathbf{v}[S_t]$ .

**HA1.**  $\mathbf{v}[\vec{Y}] = \vec{y}$ .

**HA2.**  $\text{init}(\mathbf{v}[S_t], \mathcal{X}(0))$  holds.

**HA3.** For all  $t$ ,  $\text{inv}(\mathbf{v}[S_0], \mathcal{X}(t))$  holds.

**HA4.** For all  $t \geq 0$ , if  $\mathcal{X}(t)$  is not continuous, at least one of (a) or (c) below holds; and if  $\mathcal{D}(t)$  is not continuous, at least one of (b) or (c) below holds.

(a)  $X_i^t \in \vec{Y}$  for some  $i = 1, \dots, n$ .

(b)  $S_t \in \vec{Y}$ .

(c) There is an edge  $e = (u, v, n) \in E$  such that  $\text{jump}(e, \lim_{s \rightarrow t^-} \mathcal{X}(s), \mathcal{X}(t))$  holds.

**HA5.** Defining *intervention-free* intervals in the same way as in Section 4, the following holds. For all intervention-free intervals  $(a, b)$  such that  $\mathcal{X}$  is continuous on  $(a, b)$ , we have that  $\mathcal{X}$  is right-continuous at  $a$ , differentiable on  $(a, b)$ , and its derivative  $\dot{\mathcal{X}}$  satisfies the flow condition  $\text{flow}(\mathbf{v}[S_t], \mathcal{X}(t), \dot{\mathcal{X}}(t))$  for all  $t \in (a, b)$ .

Note that HA5 is analogous to the condition ODE3 in Section

4.<sup>5</sup> The modeler is free to select a set of allowed interventions that fits the task at hand. In the example below, we choose  $\mathcal{I}$  to be the set of finite compositions of point and interval interventions on the dynamical variables and control mode variables, but the outcomes of  $M$  are well-defined for arbitrary interventions. This concludes the construction of  $M$ .

A simple thermostat and heater automaton given by [Henzinger, 2000] (with the flow conditions slightly simplified) is as follows. There is a single continuous variable  $\mathcal{T} \in \mathcal{X}$  and two control modes  $V = \{\text{OFF}, \text{ON}\}$ . The system starts in OFF with  $\mathcal{T} = 20$ , which is the desired set point. This defines  $\text{init}(v, \mathcal{X})$ . In OFF, the temperature drifts slowly downward; we have  $\dot{\mathcal{T}} = -0.1$ . An invariant of OFF is that  $\mathcal{T} \geq 18$ ; if at any point  $\mathcal{T} < 18$ , the system transitions to ON. Likewise, in ON, the temperature increases rapidly as  $\dot{\mathcal{T}} = 0.5$ , and an invariant of ON is  $\mathcal{T} \leq 22$ . These conditions define  $\text{flow}(v, \mathcal{X}, \dot{\mathcal{X}})$  and  $\text{inv}(v, \mathcal{X})$ . There are two control switches,  $e_{\text{on}}$  from OFF to ON, and  $e_{\text{off}}$  from ON to OFF. Finally, the heater can transition from OFF to ON when  $\mathcal{T} < 19$ , and from ON to OFF when  $\mathcal{T} > 21$ . Neither of these transitions affects the instantaneous value of  $\mathcal{T}$ ; that is,  $\mathcal{T}' = \mathcal{T}$ . These conditions define  $\text{jump}(e, \mathcal{X}, \mathcal{X}')$ .

Let us analyze the GSEM  $M$  corresponding to this automaton. In our case, since neither of the two possible discrete jumps change the value of  $\mathcal{T}$ , HA4 implies that  $\mathcal{T}$  must be continuous everywhere it is not intervened on. Interventions on  $\mathcal{T}$  are finite compositions of point and interval interventions. Furthermore, when  $\mathcal{T}$  is not intervened on, it cannot change very quickly; by HA5, it must obey the flow condition (either  $\dot{\mathcal{T}} = -0.1$  or  $\dot{\mathcal{T}} = 0.5$ ). Hence, given a time horizon  $\tau$ ,  $\mathcal{T}$  can cross the vertical lines  $\mathcal{T} = 19$  and  $\mathcal{T} = 21$  only finitely many times prior to  $\tau$ . The state of the heater is specified by the control mode  $S_t$ . By HA4, the heater switches ON only at times  $t$  when the state of the heater (i.e.,  $S_t$ ) is intervened on, or when  $\mathcal{T} < 19$ ; likewise, the heater switches OFF only when  $S_t$  is intervened on, or when  $\mathcal{T} > 21$ . Again, interventions on the state of the heater are finite compositions of point and interval interventions. It follows that the state of the heater changes only a finite number of times before any given time horizon  $\tau$ . Hence, it is meaningful to talk about the heater discretely changing state—before any given time  $\tau$ , the heater turns on at  $t_1, t_3$  and so on, and turns off at  $t_2, t_4$  and so on.

Now that we have some intuition for the behavior of  $M$ , we examine how  $M$  can be used to answer questions of actual cause. By HA2, the heater is initially OFF, and the temperature is initially  $\mathcal{T} = 20$ . In the absence of intervention or discrete jumps, the heater will stay OFF and the temperature will drop at the rate of 0.1 per unit time.

Consider an outcome  $\mathbf{v}$  where the heater does not turn on until, at  $t = \frac{18-20}{-0.1} = 20$ , it is required to do so by HA3;

<sup>5</sup>It is possible (and probably desirable) to strengthen HA5 analogously to the way we strengthened ODE3 to ODE3', so that interval interventions on one dynamical variable do not interfere with the flow conditions on another dynamical variable. We do not do this here, because it is not necessary for our simple example (which has only a single variable), and because doing this requires some knowledge of the structure of the predicate  $\text{flow}(v, \mathcal{X}, \dot{\mathcal{X}})$ .

specifically, by the invariant of OFF that  $\mathcal{T} \geq 18$ . If the heater had been on over any open subinterval  $(a, b)$  of  $[0, 20]$ , the temperature would have been higher than 18 by  $t = 20$  by at least  $0.5(b - a)$ . Hence, intuitively, the heater being off over any such subinterval should be considered a cause of  $T_{20} = 18$ . However, if we fix any subinterval  $(a, b) \subset [0, 20]$  and ask the formal question of whether  $S(a, b) = OFF$  is an actual cause of  $T_{20} = 18$  in  $(M, \mathbf{u}, \mathbf{v})$  (where  $S(a, b) = \{S_t \mid t \in (a, b)\}$ ), we run into problems.<sup>6</sup> AC1 and AC2 both hold, but AC3 does not. AC1 holds, because  $\mathbf{v}[S_t] = OFF$  for all  $t \in [0, 20]$ , and  $\mathbf{v}[T_{20}] = 18$ . AC2 holds, since if we choose  $\vec{W} = \emptyset$  and  $\vec{x}' = ON$  (recall that  $\vec{X} = S(a, b)$ ) we find that the outcome  $\mathbf{v}'$  of  $M$  under intervention  $\vec{X} \leftarrow \vec{x}'$  where the heater is on only during  $(a, b)$  has  $\mathbf{v}'[T_{20}] = 18 + 0.5(b - a) \neq 18$ . However, AC3 does not hold, because the open subinterval  $(a, b)$  contains other open subintervals for which AC1 and AC2 also hold, by the same arguments. This implies that there is no open interval on which the heater being off is an actual cause of  $T_{20} = 18$ .

This creates a dilemma. Since turning the heater on results in  $T_{20} > 18$ , a good definition of actual cause should provide for some cause. In this case, the resolution is that the equality  $H_t = OFF$  for any point  $t \in (0, 20)$  is in fact an actual cause of  $T_{20} = 18$ . This is because one of the solutions to  $H_t \leftarrow 1$  has  $H_s = ON$  for all  $s$  in some nonempty interval starting at  $t$ , which as before implies  $T_{20} > 18$ . So AC2 holds. AC1 clearly holds, and AC3 holds because  $H_t \leftarrow 1$  is a point intervention, therefore minimal. We believe this resolution to the dilemma is always possible in hybrid automata (if point interventions are allowed), since we see no way of defining a hybrid automaton such that when the control mode is intervened on at a point in time, the control mode does not remain at the intervened value for some nonzero amount of time in some solution (although we have not attempted to prove this).

However, we see no reason for this resolution to work in general. Other models of dynamical systems may not respond in the same way to intervention. This resolution even fails for our example hybrid automaton, if it is modified so that point interventions are not allowed. In these cases, the definition of actual causality presented in Section C fails to provide for any cause of  $T_{20} > 18$  involving only the heater state. The issue is that AC3 requires a minimal cause; but minimal causes do not exist in general when causes can involve infinitely many variables. It is an open problem to find a new definition of actual cause that handles infinitely many variables well in general. One potential solution is to broaden the set of things that count as causes in infinitary settings. In the definition presented in Section C, only conjunctions of equalities  $X = x$  can be actual causes. One could consider expanding possible causes to include infinite disjunctions over these equalities, for example, the statement that there is *some* nonempty inter-

val on which the heater is off:

$$\exists(a, b) \forall t \in (a, b) S_t = OFF.$$

However, we do not pursue this approach further in this paper.

Note that if we ask instead whether  $S_t(a, b) = OFF$  is an actual cause of  $T_{20} \leq 19$  instead, there are no problems. The answer is yes, iff  $b - a = 2$ . This agrees with the natural intuition that the heater being off for a sufficiently short time is not enough to cause the temperature to be low. (There is nothing special about 19; the solution for any value  $x > 18$  is similar.)

## D.2 Rule-based models

A rule-based model is a dynamical system that transitions probabilistically between states, with the transition defined by rewrite rules. In this section, we construct a GSEM corresponding to the generic rule-based model given by Laurent et al. [2018].

Laurent et al. [2018] show how a rule-based model can describe a reaction between a set  $S$  of *substrates* and a set  $K$  of *kinases*. The state of the mixture at time  $t$  is a binary relation  $\text{Bound} \subseteq S \times K$  that specifies which substrates are bound to which kinases, and a unary relation  $\text{Phos} \subseteq S \cup K$  that defines which substrates and kinases have a phosphate group attached. Chemical interactions between groups of molecules are intended to take place spontaneously, in an analogous fashion to radioactive decay. For example, if at time  $t$  there is a substrate  $s \in S$  and a kinase  $k \in K$  such that  $(s, k) \in \text{Bound}$  and  $s \notin \text{Phos}$ , then at time  $t + \Delta t$ , where  $\Delta t$  is drawn from an exponential distribution with a time constant that depends on the rule being applied,  $s$  will gain a phosphate group (unless in the meantime some other rule has changed the state of the mixture so that the precondition for  $s$  and  $k$  no longer holds). These updates are called *events*. Interventions correspond to blocking some interactions from taking place at specific times, for example, “between  $t = 1$  and  $t = 2$ , even if  $s$  and  $k$  satisfy the above conditions,  $s$  cannot gain a phosphate group.”

Laurent et al. explain how to simulate these dynamics using the following algorithm. For every possible target  $\text{Targ}$  of a rule  $r$ —in the example above, this would be every substrate-kinase pair  $(s, k)$ —sample from a Poisson process with parameter  $\tau / \ln(2)$  to obtain a schedule of times when the rule applies to this target. Then, starting with the initial mixture and moving through time, whenever any rule applies according to the schedule, check if the rule’s condition—in our example  $(s, k) \in \text{Bound} \wedge s \notin \text{Phos}$ —is satisfied for the target, and that the rule is not currently blocked by an intervention. If these conditions hold, update the mixture using the rule’s mapping (e.g.  $\text{Phos} \mapsto \text{Phos} \cup \{s\}$ ); otherwise, do nothing.

This algorithm can immediately be described in a GSEM model. In the example above, for each time  $t \in [0, \infty)$  we would have binary endogenous variables  $\text{Bound}_t^{(s, k)}$  for each  $s \in S, k \in K$ , along with variables  $\text{Phos}_t^x$  for each  $x \in S \cup K$ . The exogenous variables correspond to the firing schedule; we have timestamped variables  $T_t^{r, \text{Targ}}$ , one for each rule  $r$ , each target  $\text{Targ}$  compatible with that rule. (There are also exogenous variables describing the initial state of the

<sup>6</sup>Similar to Footnote 2, there is a technical issue here, because the event  $S(a, b) = OFF$  (which is an infinite conjunction of equalities) is not in the language  $\mathcal{L}(S)$ , even though the intervention  $S(a, b) \leftarrow OFF$  is. However, again, we do not view this as a problem, since the definition of actual causality makes perfect sense for this formula.

mixture.) In order to match the intervention model of [Laurent *et al.*, 2018], we add additional binary variables  $B_t^{r, \text{Targ}}$ . Intuitively,  $B_t^{r, \text{Targ}} = 1$  means that the firing of rule  $r$  applied to target Targ at time  $t$  is blocked. Finally, for bookkeeping, we have binary endogenous variables of the form  $X_t^{r, \text{Targ}}$  that model whether rule  $r$  actually fired on target Targ at time  $t$ . The unique outcome is specified in the obvious way:  $X_t^{r, \text{Targ}}$  is true exactly if, at time  $t$ , Targ satisfies the condition of  $r$ ,  $T_t^{r, \text{Targ}}$  is true, and  $B_t^{r, \text{Targ}}$  is false. If  $X_t^{r, \text{Targ}}$  is true, then at time  $t$  the state (i.e., the relations Bound and Phos) gets updated using the rule’s mapping. Interventions such as the one above can simply be described by setting some of the  $B_t^{r, \text{Targ}}$  to false; we take the set of allowed interventions  $\mathcal{I}$  to be all interventions of this form. The *trace*  $T(\mathbf{u}, I)$  is simply the (countable) sequence of variables  $X_t^{r, \text{Targ}}$  (in ascending order of  $t$ ) for which  $X_t^{r, \text{Targ}} = 1$ .

Laurent *et al.* [2018] defined notions of *enablement* and *prevention*. Enablement and prevention happen at the level of *events*, or updates to the mixture. Every event  $e$  corresponds to a variable  $X_t^{r, \text{Targ}}$ ; it *occurs* if  $X_t^{r, \text{Targ}} = 1$ . We can think of the relations Bound and Phos as binary vectors; each entry in these vectors is called a *site*. For any given event to occur, certain sites must have certain values. Hence, intuitively, given two events  $e, e'$ ,  $e$  enables  $e'$  if  $e$  is the last event before  $e'$  that modifies some site to the value that is needed for  $e$  to occur. Likewise,  $e$  prevents  $e'$  (roughly) if  $e$  is the last event before  $e'$  to set a site  $s$ , and it sets  $s$  to a value such that  $e$  cannot occur. Given a context  $\mathbf{u}$  and an intervention  $I \in \mathcal{I}$ , they considered the difference between the trace  $T(\mathbf{u})$  and the trace  $T(\mathbf{u}, I)$ . They showed that for every element of the first sequence absent from the intervened sequence, a chain of enablements and preventions could be traced back from that element to an element that was directly blocked by  $I$ . That is, enablements and preventions were sufficient to explain why each element of  $T(\mathbf{u})$  no longer in  $T(\mathbf{u}, I)$  was missing.

The notion of actual cause complements this analysis. For example, it follows from the sufficiency of enablements and preventions just discussed that if one rule firing is the actual cause of another rule firing, then a chain of enablements and preventions can be traced back from the trace entry for the second rule to the trace entry for the first. More precisely, if  $B_t^{r, \text{Targ}} = 0$  is an actual cause of  $X_{t'}^{r', \text{Targ}'} = 1$  in context  $\mathbf{u}$ , then a chain of enablements and preventions can be traced back from  $X_{t'}^{r', \text{Targ}'}$  to  $X_t^{r, \text{Targ}}$  in the pair of traces  $T(\mathbf{u}), T(\mathbf{u}, B_t^{r, \text{Targ}} \leftarrow 1)$ . Without going into the formalism of [Laurent *et al.*, 2018], a sketch of the proof of this claim is as follows. The actual cause statement implies that  $X_{t'}^{r', \text{Targ}'}$  is in  $T(\mathbf{u})$  but not in  $T(\mathbf{u}, B_t^{r, \text{Targ}} \leftarrow 1)$ , because the intervention  $B_t^{r, \text{Targ}} \leftarrow 1$  is the only one that can satisfy AC2 and AC3. (The other blocking variables take value zero in both outcomes  $M(\mathbf{u}, \emptyset)$  and  $M(\mathbf{u}, B_t^{r, \text{Targ}} \leftarrow 1)$ , so setting them to 0 is redundant and violates AC3.) The only element blocked by  $B_t^{r, \text{Targ}} \leftarrow 1$  is  $X_t^{r, \text{Targ}}$ . Hence, a chain of enablements and preventions can be traced back from  $X_{t'}^{r', \text{Targ}'}$  to  $X_t^{r, \text{Targ}}$ .

Actual cause and the GSEM machinery can also be used to answer questions not addressed by the analysis of Laurent *et al.* [2018]. We can ask counterfactual questions like “What would the state at  $t = 7$  be if every kinase gained a phosphate group at time  $t = 5$ ?” (potentially corresponding to the addition of a test tube’s worth of phosphate solution) or “Is the fact that substrate  $s$  was bound to kinase  $k$  at  $t = 1$  the actual cause of kinase  $k$  gaining a phosphate group at  $t = 2$ ?” For this reason, we believe that GSEMs are a useful addition to the rule-based causal modeling toolkit developed by Laurent *et al.*