

---

# THE MINIMAL SEARCH SPACE FOR CONDITIONAL CAUSAL BANDITS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Causal knowledge can be used to support decision-making problems. This has been recognized in the causal bandits literature, where a causal (multi-armed) bandit is characterized by a causal graphical model and a target variable. The arms are then interventions on the causal model, and rewards are samples of the target variable. Causal bandits were originally studied with a focus on hard interventions. We focus instead on cases where the arms are *conditional interventions*, which more accurately model many real-world decision-making problems by allowing the value of the intervened variable to be chosen based on the observed values of other variables. This paper presents a graphical characterization of the minimal set of nodes guaranteed to contain the optimal conditional intervention, which maximizes the expected reward. We then propose an efficient algorithm with a time complexity of  $O(|V| + |E|)$  to identify this minimal set of nodes. We prove that the graphical characterization and the proposed algorithm are correct. Finally, we empirically demonstrate that our algorithm significantly prunes the search space and substantially accelerates convergence rates when integrated into standard multi-armed bandit algorithms.

## 1 INTRODUCTION

Lattimore et al. (2016) introduced a class of problems termed *causal bandit* problems, where actions are interventions on a causal model, and rewards are samples of a chosen reward variable  $Y$  belonging to the causal model. They focus on hard interventions, where the intervened variables are set to specific values, without considering the values of any other variables. We will refer to this as a hard-intervention causal bandit problem. They propose a best-arm identification algorithm that utilizes observations of the non-intervened variables in the causal model to accelerate learning of the best arm as compared to standard multi-armed bandit (MAB) algorithms. Causal bandits have applications across a broad range of domains, particularly in scenarios requiring the selection of an intervention on a causal system. These include computational advertising and context recommendation (Bottou et al., 2013; Zhao et al., 2022), biochemical and gene interaction networks (Meinshausen et al., 2016; Basharin, 1959), epidemiology (Joffe et al., 2012), and drug discovery (Michoel & Zhang, 2023). Most of the work in causal bandits (see Section 7) focuses on developing MAB algorithms which incorporate knowledge about the causal graph. Lee & Bareinboim (2018), in contrast, use the fact that the causal graph is known not to develop yet another MAB algorithm, but to reduce the set of nodes (*i.e.* variables) of the causal graph on which hard interventions should be examined, thereby reducing the search space for hard-intervention causal bandit problems. This reduction of the search space significantly improves and scales the applicability of existing causal MAB algorithms.

It is recognized in the MAB literature that, for many if not most applications, actions are taken in a context, that is, with available information (Lattimore & Szepesvári, 2020; Agarwal et al., 2014; Dudik et al., 2011; Jagerman et al., 2020; Langford & Zhang, 2007). *E.g.*, content recommendation based on the user’s demographic characteristics, such as age, gender, nationality and occupation. Similarly, in causality, conditional interventions — where a variable  $X$  is set to a value  $g(\mathbf{Z}_X)$  through some function  $g$  after observing a set of variables (a context)  $\mathbf{Z}_X$  — are more realistic than hard or soft<sup>1</sup> interventions in many real-world scenarios. Conditional interventions were first introduced in Pearl (1994) based on the argument that “In general, interventions may involve complex policies in

---

<sup>1</sup>In a soft intervention, the intervened variable keeps its direct causes (Peters et al., 2017).

054 which a variable  $X$  is made to respond in a specified way to some set  $\mathbf{Z}_X$  of other variables.” Shpitser  
055 & Pearl (2012) motivate their interest in conditional interventions by providing the concrete example  
056 of a doctor selecting treatments based on observed symptoms and medical test results  $\mathbf{Z}_X$  to improve  
057 the patient’s health condition. The doctor performs interventions of the form  $do(X_i = x_i)$ , but “the  
058 specific values of the treatment variables are not known in advance, but instead depend on symptoms  
059 and test results performed ‘on the fly’ via policy functions  $g_i$ ” (Shpitser & Pearl, 2012). Formally,  
060 this is denoted  $do(X_i = g(Z_{X_i}))$ . See the paragraph on conditional interventions in Section 2 for  
061 further motivation and details about  $\mathbf{Z}_X$

062 **Novelty and contributions:** This work, like that of Lee & Bareinboim (2018), leverages the causal  
063 graph to *reduce the search space of the MAB problem, thereby accelerating MAB algorithms applied*  
064 *to it and effectively serving as a pre-processing step for (causal) MAB problems.* While Lee &  
065 Bareinboim (2018) study causal bandits with multi-node hard interventions in the presence of latent  
066 confounders, we focus on single-node conditional interventions under the assumption of no latent  
067 confounders. As discussed in Section 2, *restricting to single-node interventions in fact makes the*  
068 *problem more challenging*, as does considering conditional rather than hard interventions. Therefore,  
069 our work addresses a fundamentally different and non-comparable problem from that of Lee &  
070 Bareinboim (2018). Because the single-node intervention problem without latent confounders is  
071 already highly non-trivial, we leave latent confounders to future work, making our study a necessary  
072 step toward the general case. The setting we study remains widely applicable — for instance, to the  
073 examples discussed in Section 2. Explicitly, our work is novel because we consider the case where (i)  
074 the arms are *conditional interventions* (which generalize both hard and soft interventions); and (ii)  
075 the interventions are *single-node interventions*. This is the first time the minimal search space for a  
076 causal bandits problem with non-hard interventions is fully characterized. Such a characterization  
077 has also not been done for single-node interventions (of any kind). Our contributions are as follows:  
078 (a) we establish a graphical characterization of the minimal set of nodes guaranteed to contain the  
079 optimal node on which to perform a conditional intervention; and (b) we propose an algorithm which  
080 finds this set, given only the causal graph, with a time complexity of  $O(|\mathbf{V}| + |E|)$ , that is, linear in  
081 the number of nodes and edges of the causal graph. As a supplementary result, we also show that,  
082 perhaps surprisingly, the exact same minimal set would hold for the optimization problem of selecting  
083 an atomic (*i.e.* single-node and hard) intervention in a deterministic causal model. We provide proofs  
084 for the graphical characterization and correctness of the algorithm, as well as experiments that assess  
085 the fraction of the search space that can be expected to be pruned using our method, in both randomly  
086 generated and real-world graphs, and demonstrate, using well-known real-world models, that our  
087 intervention selection can significantly improve a classical MAB algorithm. Note that, if the true  
088 causal graph is unknown and instead a family of candidate graphs is available, the C4 algorithm  
089 can simply be applied to each candidate graph, and the results combined by taking the union of the  
090 resulting minimal search spaces. All proofs of the results presented in the paper can be found in the  
091 appendix. The code repository with the experiments can be found in the supplementary material.

## 092 2 PRELIMINARIES

093 **Graphs and causal models** We will make use of Directed Acyclic Graphs (DAGs). The main  
094 concepts of DAGs and notation used in this paper are reviewed in Appendix A. Furthermore, we  
095 operate within the Pearlian graphical framework of causality, where causal systems are modeled  
096 using Structural Causal Models (SCMs) (Peters et al., 2017; Pearl, 2009). An *SCM*  $\mathfrak{C}$  is a tuple  
097  $(\mathbf{V}, \mathbf{N}, \mathcal{F}, p_{\mathbf{N}})$ , where  $\mathbf{V} = (V_1, \dots, V_n)$  and  $\mathbf{N} = (N_{V_1}, \dots, N_{V_n})$  are vectors of random variables.  
098 The exogenous variables are pairwise independent, and are distributed according to the *noise distri-*  
099 *bution*  $p_{\mathbf{N}}$ , while each endogenous variable  $V_i$  is a deterministic function  $f_{V_i}$  of its noise variable  
100  $N_{V_i}$  and a (possibly empty) set of other endogenous variables  $\text{Pa}(V_i)$ , called the parents of  $V_i$ . The  
101  $V_i$  and  $N_{V_i}$  are called *endogenous* and *exogenous* (or *noise*) variables, respectively.  $R_V$  denotes  
102 the range of the random variable  $V$ .  $\mathcal{F}$  is a set of functions  $f_{V_i} : R_{\text{Pa}(V_i)} \times R_{N_{V_i}} \rightarrow R_{V_i}$ , termed  
103 *structural assignments*. The endogenous variables together with  $\mathcal{F}$  characterize a DAG called the  
104 *causal graph*  $G^{\mathfrak{C}} := (\mathbf{V}, E)$  of  $\mathfrak{C}$ , whose edge set is  $E = \{(P, X) : X \in \mathbf{V}, P \in \text{Pa}(X) \setminus \{X\}\}$ .  
105 We denote by  $\mathfrak{C}(G)$  the set of SCMs whose causal graph is  $G$ . Having an SCM allows us to model  
106 interventions: intervening on a variable changes its structural assignment  $f_X$  to a new one, say  $\tilde{f}_X$ .  
107 This intervention is then denoted  $do(f_X = \tilde{f}_X)$ . In the simplest type of interventions, called *atomic*  
*interventions*, a variable  $X$  is set to a chosen value  $x$ , thus replacing the structural assignment  $f_X$

of  $X$  with a constant function setting it to  $x$ . Such an intervention is denoted  $do(X = x)$ , and the SCM resulting from performing this intervention is denoted  $\mathcal{C}^{do(X=x)}$ . The joint distribution over the endogenous variables resulting from the atomic intervention  $do(X = x)$  is denoted  $p^{do(X=x)}$  and called the *post-intervention distribution* for this intervention. Each realization  $\mathbf{n} \in R_{\mathbf{N}}$  of the noise variables will be called a *unit*. A *deterministic SCM* is an SCM for which the noise distribution is a point mass distribution with all its mass on some (known) unit  $\mathbf{n} \in R_{\mathbf{N}}$ . Finally, nodes are denoted by upper case letters, sets of nodes by boldface letters, and variable values by lower case letters. We will make use of the fact that the structural assignments of the ancestors of an endogenous variable  $X$  (including its own structural assignment) can be composed to express  $X$  as a function  $f_X(\mathbf{n})$  of the vector  $\mathbf{n}$  of exogenous variables values. We call this<sup>2</sup>the *unrolled assignment* of  $X$ .

**Conditional interventions** Given an SCM  $\mathcal{C} = (\mathbf{V}, \mathbf{N}, \mathcal{F}, p_{\mathbf{N}})$  with causal graph  $G$ ,  $X \in \mathbf{V}$ ,  $\mathbf{Z}_X \subseteq \mathbf{V} \setminus \{X\}$ , and a (any) function  $g: R_{\mathbf{Z}_X} \rightarrow R_X$  (which we call a *policy for  $X$* ), the *conditional intervention on  $X$  given  $\mathbf{Z}_X$  for the policy  $g$* , denoted  $do(X = g(\mathbf{Z}_X))$ , is the intervention where the value of  $X$  is determined by that of  $\mathbf{Z}_X$  through  $g$  (Pearl, 2009). The precise conditioning set  $\mathbf{Z}_X$  for each  $X$  is pre-determined by the specific problem or application, or by the practitioner. In order to systematically study conditional interventions, we will need to make some assumptions of what nodes can reasonably be in  $\mathbf{Z}_X$ , i.e. what variables can we expect to have knowledge of at the time of applying the policy  $g$  to intervene on  $X$ . As noted in Pearl (1994; 2009), the nodes in  $\mathbf{Z}_X$  cannot be descendants of  $X$  in  $G$ . Hence,  $\mathbf{Z}_X \subseteq \mathbf{V} \setminus \text{De}(X)$ . On the other hand, all (proper) ancestors of  $X$  are realized before  $X$ . Since we will be dealing with the case with no latent variables, we can assume that all ancestors of  $X$  are observed, and can be used by a policy  $g$  to set  $X$  to a value  $g(\mathbf{Z}_X)$ . Thus, we assume<sup>3</sup>that  $\text{An}(X) \setminus \{X\} \subseteq \mathbf{Z}_X$ . We will then focus on the case where for each  $X$ , the conditioning set  $\mathbf{Z}_X$ , chosen by the practitioner, obeys the inclusion relations  $\text{An}(X) \setminus \{X\} \subseteq \mathbf{Z}_X \subseteq \mathbf{V} \setminus \text{De}(X)$ . Furthermore, we focus on cases where the context that is available for an intervention is also available for later interventions. As an example, consider the case where a traffic controller needs to intervene on the delay  $D_{i,s}$  of a train  $i$  at a train station  $s$  (for example by forcing it to wait for 5 extra minutes before departing). Clearly, all delays  $D_{i',s'}$  of all train/station pairs affecting  $D_{i,s}$  have already been observed, and can therefore be used when selecting  $D_{i,s}$ . As another example, similar to the one used in Pearl (1994) when first introducing conditional interventions, consider the situation where a doctor must decide, over a period of three weeks, whether and when to intervene on the weight, blood pressure or renal blood flow of a patient, in order to improve the patient’s kidney function. The goal is to maximize kidney function (variable Kidney3) at the end of the third week. Due to side-effects, the patient can only be prescribed medication for one week. The causal graph for this situation can be found in Figure 4, Appendix B. Notice that at the time of intervening on a node  $X_i$ , the doctor can use information about all measurements made until then. For instance, all the data available when performing an intervention on the renal flow on week 1 (node RenalFlow1) will also be available when intervening on the renal flow on week 2 (node RenalFlow2). Mathematically, this last assumption can be written  $W \in \text{An}(X) \Rightarrow \mathbf{Z}_W \subseteq \mathbf{Z}_X$ . We then say that  $\mathbf{Z}_X$  is an *observable conditioning set for  $X$* .

**Conditional causal bandits** Recall that a MAB problem consists of an agent pulling an arm  $a \in \mathcal{A}$  at each round  $t$ , resulting in a reward sample  $Y_t$  from an unknown distribution associated to the pulled arm (Lattimore & Szepesvári, 2020). We denote the mean reward for arm  $a$  by  $\mu_a$  and the mean reward for the best arm by  $\mu^* = \max_{a \in \mathcal{A}} \mu_a$ . The objective is to maximize the total reward obtained over all  $T$  rounds. Equivalently, this can be framed as minimizing the cumulative regret  $\text{Reg}_T = T\mu^* - \sum_{t=1}^T \mathbb{E}[Y_t]$ . We now introduce a novel type of (causal) MAB problem. Consider the setting where the bandit’s reward is a (endogenous) variable  $Y$  in an SCM  $\mathcal{C} = (\mathbf{V}, \mathbf{N}, \mathcal{F}, p_{\mathbf{N}})$ , and the arms are the conditional interventions  $do(X = g(\mathbf{Z}_X))$ , where  $X \in \mathbf{V} \setminus \{Y\}$ . Furthermore, the agent has knowledge of the causal graph  $G$  of  $\mathcal{C}$ , but not of the structural assignments  $\mathcal{F}$  or the noise distribution  $p_{\mathbf{N}}$ . We call this a *single-node conditional-intervention causal bandit*, or simply *conditional causal bandit*. The reward distribution for arm  $do(X = g(\mathbf{Z}_X))$  is the post-intervention distribution  $p_Y^{do(X=g(\mathbf{Z}_X))}$ , and is unknown to the agent, since it has no knowledge of  $\mathcal{F}$ . Notice that selecting an arm can be subdivided in (i) choosing a node  $X$  to intervene on; and (ii) choosing a policy  $g$ , i.e. choosing a value to set  $X$  to given the observed variables  $\mathbf{Z}_X$ . *We do not impose any*

<sup>2</sup>The formal definition can be found in Appendix C.

<sup>3</sup>We are not claiming that all variables in  $\text{An}(X) \setminus \{X\}$  need to be in  $\mathbf{Z}_X$  for the best decision to be made, or for our results to hold, but that we *can* always include them in  $\mathbf{Z}_X$  under the assumptions of our problem.

restrictions on the function  $g$ . The conditioning sets  $\mathbf{Z}_X$  are specified in advance, as described in the paragraph on conditional interventions above. In this paper, we find the minimal set of nodes that need to be considered by the agent in step (i). The value of  $X$  chosen in step (ii) can be selected by an MAB algorithm.

As stressed in Section 1, the novelty of our problem lies in the fact that we deal with *conditional interventions* that are *single-node*. Both of these characteristics of our problem complicate the analysis. Unsurprisingly, searching over conditional interventions is more complicated than over hard or soft interventions. Perhaps more unexpectedly, single-node interventions also make a search for a minimal search space more involved. Indeed, if one allows for interventions on arbitrary sets, one simply needs to intervene on all the parents  $\text{Pa}(Y)$  of  $Y$  (Lee & Bareinboim, 2018). Since in our case the agent cannot do this whenever  $|\text{Pa}(Y)| > 1$ , the minimal search space will, as we will see, be complex even without unobserved confounding. That said, the assumption that there is no unobserved confounding is a limitation of this paper, and a natural next step for future work (see Section 7).

### 3 CONDITIONAL-INTERVENTION SUPERIORITY

In this section, we will define a preorder  $\succeq_Y^c$  of “conditional-intervention superiority” on nodes of an SCM. If  $X \succeq_Y^c W$ , then  $W$  can never be a better node than  $X$  to intervene on with a conditional intervention<sup>4</sup>. We will then show that, perhaps surprisingly, this relation is equivalent to another superiority relation, defined in terms of atomic interventions in a deterministic SCM.

**Definition 1** (Conditional-Intervention Superiority). *Let  $X, W, Y$  be nodes of a DAG  $G$ .  $X$  is conditional-intervention superior to  $W$  relative to  $Y$  in  $G$ , denoted  $X \succeq_Y^c W$ , if for all SCM with causal graph  $G$  there is a policy  $g$  for  $X$  such that for every observable conditioning sets  $\mathbf{Z}_X$  and  $\mathbf{Z}_W$  for  $X$  and  $W$  and all policies  $h$  for  $W$ ,*

$$\mathbb{E}_{\mathbf{n}} \bar{f}_Y^{\text{do}(X=g(\mathbf{Z}_X))}(\mathbf{n}) \geq \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{\text{do}(W=h(\mathbf{Z}_W))}(\mathbf{n}). \quad (1)$$

A similar relation can be defined for atomic interventions in deterministic SCMs, where the vector  $\mathbf{N}$  of exogenous variables is fixed to a *known* value  $\mathbf{n}$  (see Section 2).

**Definition 2** (Deterministic Atomic-Intervention Superiority). *Let  $X, W, Y$  be nodes of a DAG  $G$ .  $X$  is deterministically atomic-intervention superior to  $W$  relative to  $Y$ , denoted  $X \succeq_Y^{\text{det},a} W$ , if for every SCM  $\mathfrak{C}$  with causal graph  $G$  and every unit  $\mathbf{n}$  there is  $x \in R_X$  such that no atomic intervention on  $W$  results in a larger  $Y$  than the value of  $Y$  resulting from setting  $X = x$ . That is, for all  $(\mathfrak{C}, \mathbf{n}) \in \mathfrak{C}(G) \times R_{\mathbf{N}}$ :*

$$\exists x \in R_X: \forall w \in R_w, \bar{f}_Y^{\text{do}(X=x)}(\mathbf{n}) \geq \bar{f}_Y^{\text{do}(W=w)}(\mathbf{n}). \quad (2)$$

We extend Definitions 1 and 2 for sets of nodes in the obvious way:  $\mathbf{X}$  is superior to  $\mathbf{W}$  if every node in  $\mathbf{W}$  is inferior to some node in  $\mathbf{X}$ .

**Definition 3.** *Let now  $\mathbf{X}, \mathbf{W}$  be sets of nodes of  $G$ .  $\mathbf{X}$  is conditional-intervention superior (respectively deterministic atomic intervention superior) to  $\mathbf{W}$ , also denoted  $\mathbf{X} \succeq_Y^c \mathbf{W}$  (respectively  $\mathbf{X} \succeq_Y^{\text{det},a} \mathbf{W}$ ), if  $\forall W \in \mathbf{W}, \exists X \in \mathbf{X}$  such that  $X \succeq_Y^c W$  (respectively  $X \succeq_Y^{\text{det},a} W$ ).*

The two relations  $\succeq_Y^c, \succeq_Y^{\text{det},a}$  actually coincide (both for nodes and sets of nodes).

**Proposition 4** (Conditional vs Atomic superiority). *Let  $X, W, Y$  be nodes in a DAG  $G$ . Then  $X$  is conditional-intervention superior to  $W$  relative to  $Y$  in  $G$  if and only if  $X$  is deterministic atomic-intervention superior to  $W$  relative to  $Y$  in  $G$ . That is,  $X \succeq_Y^c W \Leftrightarrow X \succeq_Y^{\text{det},a} W$ .*

Since these two relations are equivalent, we henceforth refer simply to interventional superiority without further specification, and use the symbol  $\succeq_Y$  when distinguishing them is not necessary. We will use Proposition 4 to simplify our problem. Since deterministic atomic interventions are easier to reason about, we use them in formulating proposals for the minimal search space and in our proofs.

<sup>4</sup>The relation between nodes introduced by Lee & Bareinboim (2018) is similar, but pertains to multi-node hard interventions.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

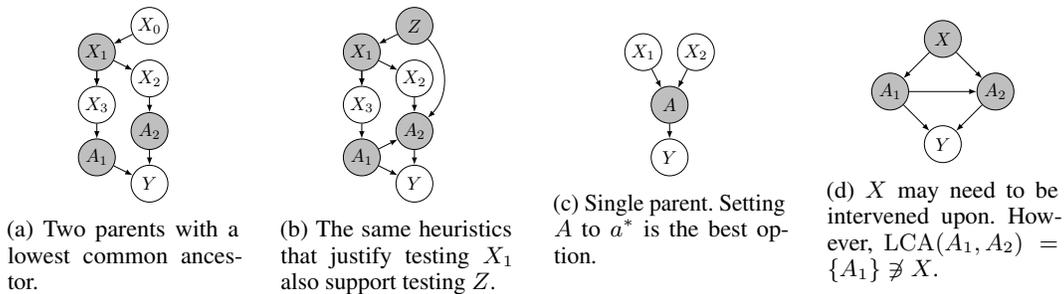


Figure 1: Examples illustrating heuristics behind the graphical characterization of the mGISS. The gray nodes are the elements of the mGISS relative to  $Y$ .

#### 4 GRAPHICAL CHARACTERIZATION OF THE MGISS

**Goal** Our aim is to develop a method to identify, based on a causal graph  $G$ , the smallest set of nodes that are “worth testing” when attempting to maximize  $Y$  by performing one single-node intervention. Specifically, regardless of the structural causal model  $\mathcal{C}$  associated with  $G$ , we want to ensure that the optimal intervention can be discovered within this selected set of nodes. We define this set as follows:

**Definition 5** (GISS and mGISS). *Let  $G$  be a DAG with set of nodes  $\mathbf{V}$ . A globally interventionally superior set (GISS) of  $G$  relative to  $Y$ , is a subset  $\mathbf{U}$  of  $\mathbf{V} \setminus \{Y\}$  satisfying  $\mathbf{U} \succeq_Y (\mathbf{V} \setminus \{Y\}) \setminus \mathbf{U}$ . A minimal globally interventionally superior set (mGISS) is a GISS which is minimal with respect to set inclusion.*

This set is unique, so that we can talk of *the* minimal globally interventionally superior set.

**Proposition 6** (Uniqueness of the mGISS). *Let  $G$  be a DAG and  $Y$  a node of  $G$  with at least one parent. The minimal globally interventionally superior set of  $G$  relative to  $Y$  is unique. We denote it by  $\text{mGISS}_Y(G)$ .*

**Intuition** Since the value of  $Y$  is completely determined by the values of its parents  $A_1, \dots, A_m$ , along with the fixed value  $n_Y$  of a noise variable that cannot be intervened upon (see Definition 2), we aim to induce the parents to acquire the combination of values  $(a_1^*, \dots, a_m^*)$  that maximizes  $Y$  when  $N_Y = n_Y$ . If this is not possible to achieve using a single intervention, we aim to obtain the best combination possible. Clearly, the parents of  $Y$  themselves need to be tested by bandit algorithms: there may be one parent on which  $Y$  is highly dependent, in such a way that there is a value of that parent which will maximize  $Y$ . In the particular case where  $Y$  has a single parent  $A$ , that node is the only node worth intervening on, since all other nodes can only influence  $Y$  through  $A$ . Indeed, if  $a^* \in R_A$  is the value of  $A$  which maximizes  $Y$ , it is not necessary to try to find an intervention on ancestors of  $A$  which results in  $A = a^*$ : just set  $A = a^*$  directly (Figure 1c). If  $Y$  has two or more parents, it is possible that a single intervention on one of the  $A_i$  does not yield the best possible outcome. Instead, a better configuration (potentially even the ideal case  $(a_1^*, \dots, a_m^*)$ ) may be achieved by intervening on a common ancestor of some or all of the  $A_i$  (Figure 1a). Notice that  $X_0$  is also a common ancestor of  $A_1, A_2$ , but one is never better off intervening on  $X_0$  than on  $X_1$ . This seems to indicate that testing interventions on, for instance, all lowest common ancestors (LCAs, see Appendix A) of the parents of  $Y$ , and only them, is necessary. While this works in Figure 1a, it fails for a graph such as Figure 1d, where  $X$  needs to be tested and yet it is not in  $\text{LCA}(A_1, A_2) = \{A_1\}$ . This suggests that we need to define a stricter notion of common ancestor to make progress in characterizing  $\text{mGISS}_Y(G)$ .

**Definition 7** (Lowest Strict Common Ancestors of a Pair of Nodes). *The node  $V \in \mathbf{V}$  is a strict common ancestor of  $X, Y \in \mathbf{V}$  if  $V$  is a common ancestor of  $X, Y$  from which both  $X$  and  $Y$  can be reached from  $V$  with paths  $V \dashrightarrow X$  and  $V \dashrightarrow Y$  not containing  $Y$  and  $X$ , respectively. The set of strict common ancestors of  $X, Y$  is denoted  $\text{SCA}(X, Y)$ . Furthermore,  $V$  is a lowest strict common ancestor of  $X, Y \in \mathbf{V}$  if  $V$  is a minimal element of  $\text{SCA}(X, Y)$  with respect to the ancestor partial order  $\preceq$ . The set of lowest strict common ancestors of  $X, Y$  is denoted  $\text{LSCA}(X, Y)$ .*

**Definition 8** (Lowest Strict Common Ancestors of a Set). *Let  $\mathbf{U} \subseteq \mathbf{V}$  and  $V \in \mathbf{V} \setminus \mathbf{U}$ . The node  $V$  is a lowest strict common ancestor of  $\mathbf{U}$  if it is a lowest strict common ancestor of some pair of nodes*

270  $U, U'$  in  $\mathbf{U}$ . The set of lowest strict common ancestors is denoted  $\text{LSCA}(\mathbf{U})$ . That is,  
 271 
$$\text{LSCA}(\mathbf{U}) := \{V \in \mathbf{V} \setminus \mathbf{U} : \exists U, U' \in \mathbf{U} \text{ s.t. } V \in \text{LSCA}(U, U')\}.$$
 (3)  
 272

273 Our heuristic argument so far suggests that we need to test the parents of  $Y$  and their LSCAs. However,  
 274 there are additional nodes that must be considered: the reasoning for testing the lowest strict common  
 275 ancestors of the parents can be repeated. For instance, in Figure 1b, the best possible configuration of  
 276 the  $A_i$  may be achieved by intervening on  $Z$ . Such an intervention could result in a combination of  
 277 values of  $X_1$  and  $A_2$  that leads to the best possible combinations of  $A_1$  and  $A_2$ . This suggests that the  
 278  $\text{mGISS}_Y(G)$  should be determined by recursively finding all the LSCAs of the parents of  $Y$ , then the  
 279 LSCAs of that set, and so on, ultimately resulting in what we call the ‘‘LSCA closure of the parents  
 280 of  $Y$ ’’, denoted  $\mathcal{L}^\infty(\text{Pa}(Y))$ . In the remainder of this section, we formally define  $\mathcal{L}^\infty(\text{Pa}(Y))$ , find a  
 281 simple graphical characterization for it, and prove that it indeed equals  $\text{mGISS}_Y(G)$ .

282 **Definition 9** (LSCA closure). For every  $i \in \mathbb{N}$  we define the  $i^{\text{th}}$ -order LSCA set  $\mathcal{L}^i(\mathbf{U})$  of  $\mathbf{U} \subseteq \mathbf{V}$  as  
 283 follows:

$$\mathcal{L}^0(\mathbf{U}) := \mathbf{U}, \text{ and } \mathcal{L}^i(\mathbf{U}) := \text{LSCA}(\mathcal{L}^{i-1}(\mathbf{U})) \cup \mathcal{L}^{i-1}(\mathbf{U}). \quad (4)$$

284 The LSCA closure  $\mathcal{L}^\infty(\mathbf{U})$  of  $\mathbf{U}$  is given by<sup>5</sup>

$$\mathcal{L}^\infty(\mathbf{U}) := \mathcal{L}^{k^*}(\mathbf{U}), \text{ where } k^* = \min\{i \in \mathbb{N} : \mathcal{L}^i(\mathbf{U}) = \mathcal{L}^{i+1}(\mathbf{U})\}. \quad (5)$$

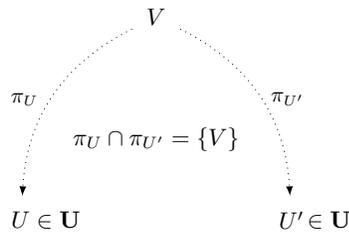
287 *Example 10.* Consider the graph in Figure 1b and set  $\mathbf{U} = \{A_1, A_2\}$ . Then,  $\mathcal{L}^0(\mathbf{U}) =$   
 288  $\{A_1, A_2\}$ ,  $\mathcal{L}^1(\mathbf{U}) = \{X_1, A_1, A_2\}$ ,  $\mathcal{L}^2(\mathbf{U}) = \mathcal{L}^3(\mathbf{U}) = \{Z, X_1, A_1, A_2\} = \mathcal{L}^\infty(\mathbf{U})$ .

289 We will introduce the notion of ‘‘ $\Lambda$ -structures’’ (Figure 2a), which provides an alternative, elegant,  
 290 simple graphical characterization of  $\mathcal{L}^\infty(\text{Pa}(Y))$ . It will also be instrumental in the proofs of the  
 291 main results of this paper.

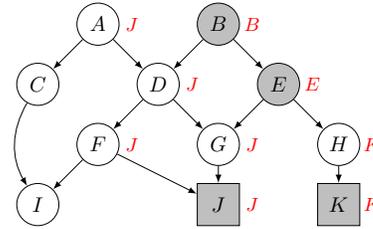
292 **Definition 11** ( $\Lambda$ -structure). Let  $V, A, B \in \mathbf{V}$ . Furthermore, let  $\pi_A : V \dashrightarrow A$ ,  $\pi_B : V \dashrightarrow B$  be  
 293 paths. The tuple  $(V, \pi_A, \pi_B)$  is a  $\Lambda$ -structure over  $(A, B)$  if  $\pi_A$  and  $\pi_B$  only intersect at  $V$ . Now,  
 294 let  $\mathbf{U}, \mathbf{W} \subseteq \mathbf{V}$ . The node  $V$  is said to form a  $\Lambda$ -structure over  $(\mathbf{U}, \mathbf{W})$  if there are nodes  $U \in \mathbf{U}$   
 295 and  $W \in \mathbf{W}$ , and paths  $\pi_U : V \dashrightarrow U$ ,  $\pi_W : V \dashrightarrow W$  such that  $(V, \pi_U, \pi_W)$  is a  $\Lambda$ -structure over  
 296  $(U, W)$ . Denote by  $\Lambda(\mathbf{U}, \mathbf{W})$  the set of all nodes forming a  $\Lambda$ -structure over  $(\mathbf{U}, \mathbf{W})$ .

297 Notice that, if  $V \in \mathbf{U} \cap \mathbf{W}$ , then trivially  $V \in \Lambda(\mathbf{U}, \mathbf{W})$ : just take the trivial paths  $\pi = \pi' = (V)$ .

298 **Theorem 12** (Simple Graphical Characterization of LSCA Closure). A node  $V \in \mathbf{V}$  is in the LSCA  
 299 closure  $\mathcal{L}^\infty(\mathbf{U})$  of  $\mathbf{U} \subseteq \mathbf{V}$  if and only if  $V$  forms a  $\Lambda$ -structure over  $(\mathbf{U}, \mathbf{U})$ . I.e.  $\mathcal{L}^\infty(\mathbf{U}) = \Lambda(\mathbf{U}, \mathbf{U})$ .



302  
303  
304  
305  
306  
307  
308  
309 (a) A  $\Lambda$ -structure over  $(\mathbf{U}, \mathbf{U})$ . Theorem 12 states that the LSCA closure  $\mathcal{L}^\infty(\mathbf{U})$  of a set  $\mathbf{U}$  is the set of all such structures.



310  
311  
312 (b) Illustration of the connectors in a graph. The square nodes belong to  $\mathbf{U}$ , the connector of each node is written in red next to its node, and the LSCA closure  $\mathcal{L}^\infty(\mathbf{U})$  consists of the gray nodes.

313  
314 Figure 2

315 We are now ready for the main result of this paper.

316 **Theorem 13** (Superiority of the LSCA Closure). Let  $G$  be a causal graph and  $Y$  a node of  $G$  with at  
 317 least one parent. Then, the LSCA closure  $\mathcal{L}^\infty(\text{Pa}(Y))$  of the parents of  $Y$  is the minimal globally  
 318 interventionally superior set  $\text{mGISS}_Y(G)$  of  $G$  relative to  $Y$ .

319 We emphasize that, due to Proposition 4, this graphical characterization of the  $\text{mGISS}_Y(G)$  is valid  
 320 both for conditional interventions in a probabilistic causal model as for atomic interventions in a  
 321 deterministic causal model (i.e. a causal model with known  $\mathbf{n}$ ).  
 322

323 <sup>5</sup>Notice that the existence of the  $k^*$  is guaranteed, since by construction  $\mathcal{L}^i(\mathbf{U}) \subseteq \mathcal{L}^{i+1}(\mathbf{U}) \subseteq \mathbf{V}$  for all  $i \in \mathbb{N}$  and  $\mathbf{V}$  is finite.

---

324 5 ALGORITHM TO FIND THE MINIMAL GLOBALLY INTERVENTIONALLY  
325 SUPERIOR SET  
326

---

327 **Algorithm 1 C4**

---

329 1: **input:** DAG  $G = (\mathbf{V}, E)$ , set of nodes  $\mathbf{U} \subseteq \mathbf{V}$   
330 2: **output:** The closure  $\mathcal{L}^\infty(\mathbf{U})$   
331 3:  $S \leftarrow \mathbf{U}$  ▷ initialize closure  
332 4:  $c[V] \leftarrow V$  for  $V \in \mathbf{U}$  ▷ initialize connectors  
333 5: **for**  $V \in \mathbf{V} \setminus \mathbf{U}$  in reverse topological order **do**  
334 6:    $\mathbf{C} \leftarrow \{c[V'] : V' \in \text{Ch}(V) \cap \text{An}(\mathbf{U})\}$   
335 7:   **if**  $|\mathbf{C}| = 1$  **then**  
336 8:      $c[V] \leftarrow X$  where  $\mathbf{C} = \{X\}$   
337 9:   **else if**  $|\mathbf{C}| > 1$  **then**  
338 10:      $c[V] \leftarrow V, S \leftarrow S \cup \{V\}$  ▷  $V$  is added to closure  
339 11: **return**  $S$

---

340 The Closure Computation via Children with Multiple Connectors (C4) Algorithm (Algorithm 1)  
341 computes the closure  $\mathcal{L}^\infty(\mathbf{U})$  in  $O(|\mathbf{V}| + |E|)$  time, using *connectors* (illustrated in Figure 2b):

342 **Definition 14** (Connector). *Let  $G = (\mathbf{V}, E)$  be a DAG,  $\mathbf{U} \subseteq \mathbf{V}$ ,  $V \in \text{An}(\mathbf{U})$ . The  $\mathbf{U}$ -connector  
343  $c[V]$  of  $V$  in  $G$  is defined recursively. Let  $\mathbf{C} = \{c[V'] : V' \in \text{Ch}(V) \cap \text{An}(\mathbf{U})\}$  be the set of  
344  $V$ 's children's connectors. If  $V \in \mathbf{U}$ , then  $c[V] := V$ . If  $V \notin \mathbf{U}$ : if  $|\mathbf{C}| = 1$  and  $\mathbf{C} = \{X\}$  then  
345  $c[V] := X$ , otherwise  $c[V] := V$ .*

346 Lemma 15 illuminates the connector's relation to  $\mathcal{L}^\infty(\mathbf{U})$ :  $c[V]$  "connects"  $V$  to  $\mathcal{L}^\infty(\mathbf{U})$  in that it  
347 is the first node in  $\mathcal{L}^\infty(\mathbf{U})$  in any path from  $V$  to  $\mathcal{L}^\infty(\mathbf{U})$ . Thus,  $c[V]$  mediates all influence that  $V$   
348 exerts over  $\mathcal{L}^\infty(\mathbf{U})$ .

349 **Lemma 15.** *Let  $G = (\mathbf{V}, E)$  be a DAG,  $\mathbf{U} \subseteq \mathbf{V}$ ,  $V \in \text{An}(\mathbf{U})$ .  $c[V]$  is the unique node s.t. a path  
350  $\pi_{c[V]} : V \dashrightarrow c[V]$  exists where  $\pi_{c[V]} \cap \mathcal{L}^\infty(\mathbf{U}) = \{c[V]\}$  (if  $V$  is its own connector, the path is  
351 trivial). This is equivalent to: for every node  $X \in \mathcal{L}^\infty(\mathbf{U})$  and path  $\pi_X : V \dashrightarrow X$ ,  $c[V]$  is the  
352 maximal element of  $\pi_X \cap \mathcal{L}^\infty(\mathbf{U})$  w.r.t. the ancestor partial order  $\preceq$ .*

353 Crucially, Lemma 15 implies that  $V \in \mathcal{L}^\infty(\mathbf{U}) \Leftrightarrow c[V] = V$ . Intuitively, if all children of  $V$  have the  
354 same connector  $X$  (i.e.  $\mathbf{C} = \{X\}$ ), then  $V$  can only influence  $\mathbf{U}$  via  $X$ , making  $X$  interventionally  
355 superior to  $V$ , and thus  $V \notin \mathcal{L}^\infty(\mathbf{U})$ . On the other hand, if  $V$ 's children have multiple connectors (i.e.  
356  $|\mathbf{C}| > 1$ ), then interventions on  $V$  can influence all those connectors, so  $V$  is a potentially worthwhile  
357 candidate for intervention, and thus  $V \in \mathcal{L}^\infty(\mathbf{U})$ . This establishes correctness of C4, which finds all  
358 nodes satisfying  $c[V] = V$  in linear time.

359 **Theorem 16.** *C4 correctly computes  $\mathcal{L}^\infty(\mathbf{U})$ , and runs in  $O(|\mathbf{V}| + |E|)$  time.*

360  
361  
362 **6 EXPERIMENTAL RESULTS**

363 We evaluate C4 on both random and real graphs. Additionally, we examine the impact of our method  
364 on the cumulative regret of a bandit algorithm.

365 **Search space reduction in random graphs** We applied the C4 algorithm to randomly generated  
366 DAGs using the Erdős-Rényi model for  $N$  graphs and probability  $p$  (Erdős & Rényi, 1959) adapted  
367 to DAG-generation<sup>6</sup>. We generated 1000 graphs using 20, 100, 300, and 500 nodes, and varying the  
368 expected (total) degree of nodes from 2 to 11 in steps of 3. For each graph  $G$ , we set the target  $Y$  to be  
369 the node with the most ancestors, used C4 to compute  $\mathcal{L}^\infty(\text{Pa}(Y)) = \text{mGISS}_Y(G)$ , and calculated  
370 the fraction of nodes in  $\text{An}(Y) \setminus \{Y\}$  that remain in  $\text{mGISS}_Y(G)$ . The results revealed that, for a  
371 given number of nodes, graphs with lower expected degrees benefit more from our method (i.e. their  
372  $\text{mGISS}_Y(G)$  correspond to smaller fractions of  $\text{An}(Y) \setminus \{Y\}$ ). Furthermore, for a fixed expected  
373 degree, our method is more effective for higher numbers of nodes. For example, for graphs with 500  
374 nodes, our method is more effective for higher numbers of nodes. For example, for graphs with 500  
375 nodes, our method is more effective for higher numbers of nodes. For example, for graphs with 500  
376 nodes, our method is more effective for higher numbers of nodes.

---

377 <sup>6</sup>After fixing a total order  $\preceq$  on the nodes, each pair of nodes  $V, u$  with  $V \preceq u$  is assigned an edge  $(V, u)$   
with probability  $p$ . The value  $p$  can be used to control the expected degree.

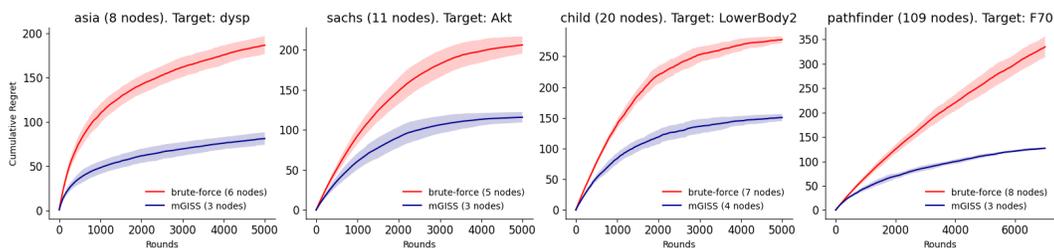


Figure 3: Comparison of cumulative regret curves for node selection using a UCB-based bandit algorithm for conditional interventions, with (mGISS) and without (brute-force) pruning the search space. These curves were obtained by averaging over 500 runs for the `bnlearn` datasets `asia`, `sachs` and `child`, and over 300 runs for `pathfinder`. For every dataset, pruning the search space with the C4 algorithm results in faster convergence and smaller values of regret.

nodes, the mGISS retained, on average, 17%, 29%, 62% and 77% of the nodes, for expected degrees of 2, 5, 8 and 11, respectively. Moreover, graphs with an expected degree of 5 saw these numbers decrease from 70% at 20 nodes to 47%, 35% and 29% for 100, 300 and 500 nodes, respectively. The complete results are presented in Figure 5 (Appendix H). These results are not surprising: if the average degree is small compared to the number of nodes, the edge density is small, in which case we expect fewer  $\Lambda$ -structures to form over  $\text{Pa}(Y)$ . Graphs modeling real-world systems tend to have low average degrees, as can be seen in the graphs from the popular Bayesian network repository `bnlearn`. Therefore, we expect our method to be especially effective in those graphs. We test this below.

**Search space reduction in real-world graphs** We tested our method in most graphs from the `bnlearn` repository<sup>7</sup>, as well as on a graph representing the causal relationships between train delays in a segment of the railway system of the Netherlands (see Appendix H). For each graph, we set  $Y$  to be the node with most ancestors<sup>8</sup>. The results are presented in the bar plot of Figure 6 (Appendix H). This confirmed that realistic models with larger graphs tended to benefit more from our method, with a reduction of over 90% of the search space for some of the largest models. Notice also that these models indeed have relatively small average degrees, all below 4.0. From this, we conclude that we can expect our method to be useful when reducing the search space of conditional causal bandit tasks in real-world causal models, especially when they are large.

**Impact on conditional intervention bandits** We present empirical evidence that restricting the node search space to the mGISS allows a straightforward UCB-based<sup>9</sup> algorithm (which we call `CondIntUCB`) for conditional causal bandits to converge more rapidly to better nodes. As explained in Section 2, on each round the algorithm must (i) choose which node  $X$  to intervene on; and (ii) choose the value for  $X$ , given its conditioning set  $Z_X$ <sup>10</sup>. Choice (i) employs UCB over nodes, while choice (ii) utilizes a UCB instance specific to the conditioning set value. In other words, for each realization of  $Z_X$  (each context) there is a UCB. This is identical to what is described in Lattimore & Szepesvári (2020, §18.1) for contextual bandits with one bandit per context. The cumulative regret<sup>11</sup> is computed with respect to node choice, since we want to see how our node selection method affects the quality of node choice by `CondIntUCB`. We use 4 real-world datasets from the `bnlearn` repository, and again choose the node of each dataset with the most ancestors as the target<sup>8</sup>. These datasets were selected because their graphical structures are non-trivial<sup>12</sup> and both  $\text{An}(Y)$  and  $\text{mGISS}_Y(G)$  are sufficiently small to allow experimentation with our setup. For each dataset, we run `CondIntUCB` up to 500 times and plot the two average cumulative regret curves along with their standard deviations, corresponding to using all nodes (brute-force) and the mGISS nodes (Figure 3). The total number of

<sup>7</sup>All which can be imported in Python using the library `pgmpy`.

<sup>8</sup>We also require  $Y$  to have more than one parent, to avoid the trivial case with  $|\text{mGISS}_Y(G)| = 1$ .

<sup>9</sup>The Upper Confidence Bound (UCB) algorithm is widely used. See *e.g.* Lattimore & Szepesvári (2020).

<sup>10</sup>For simplicity, we use the smallest observable conditioning set  $Z_X = \text{An}(X) \setminus \{X\}$  (see Section 2).

<sup>11</sup>For the computation of regret, we use the estimated best arm, defined as the arm that most runs concluded to be the best at the end of training.

<sup>12</sup>In contrast, the `cancer` dataset, for example, only has nodes whose mGISS is either all of the node’s ancestors or a single node.

---

432 rounds is chosen as to observe (near) convergence. These results show that cumulative regret curves  
433 can be significantly improved—meaning that better nodes are selected earlier for applying conditional  
434 interventions—if the search space over nodes is pruned using our C4 algorithm.

## 436 7 RELATED WORK AND CONCLUSION

437  
438  
439 Recent works in “contextual causal bandits” address interventions that account for context, bearing  
440 resemblance to our problem. However, our problem remains distinct. In a  $K$ -arm contextual bandit  
441 problem, each round is associated with a context that determines the reward distributions of the  $K$   
442 arms. The agent uses the context to select one of the  $K$  arms. A general approach to solving such  
443 problems is to maintain a separate standard bandit algorithm for each context. More efficient solutions  
444 typically rely on assumptions about relationships between contexts (Lattimore & Szepesvári, 2020).  
445 In contrast, a conditional causal bandit problem involves, in each round, an intervened node  $X$  and  
446 an observed context that is a sample of  $\mathbf{Z}_X$ . This context determines the reward distributions of the  
447  $K = |R_X|$  possible atomic interventions on  $X$ , and the agent chooses among these according to a  
448 policy. Thus, a conditional causal bandit problem can be interpreted as a collection of contextual  
449 bandits, one for each node  $X$  in a causal graph. In particular, conditional causal bandits are not simply  
450 particular cases of contextual bandits. In this paper, we leverage the structure of the causal graph  
451 to eliminate certain nodes, *i.e.*, to exclude some of these contextual bandits from consideration. In  
452 Madhavan et al. (2024), the term “contexts” is used in a very different way to the one used in our  
453 paper, actually referring to different graphs as opposed to different variable values. Subramanian &  
454 Ravindran (2022; 2024) tackle the scenario in which an intervention is performed, with knowledge  
455 of a given set of context variables, on a *pre-chosen* variable  $X$  that has an edge into  $Y$  (and no  
456 other outgoing edges). This approach can be understood as selecting a conditional intervention for a  
457 predefined node from a very simple graph. In contrast, in our setting we need to choose what variable  
458 to intervene on to begin with, and there are no restrictions on the causal graph.

458 As mentioned in Section 1, Lattimore et al. (2016) introduced the original causal bandit problems,  
459 which involve hard interventions in causal models. Subsequent works (Sen et al., 2017; Yabe et al.,  
460 2018; Lu et al., 2020; Nair et al., 2021; Sawarni et al., 2023; Maiti et al., 2022; Feng & Chen, 2023)  
461 proposed algorithms for variants of causal bandits with both hard and soft interventions, budget  
462 constraints, and unobserved confounders.

463 All of the works described above proposed algorithms which aim at accelerating learning by utilizing  
464 knowledge of the causal model. As explained in Section 1, this contrasts with our work, which, like  
465 the work of Lee & Bareinboim (2018; 2019), uses knowledge of the causal graph to find a minimal  
466 search space (over the nodes) for causal bandits. And, while the latter focus on multi-node, hard  
467 interventions, we focus on single-node, conditional interventions.

468 The work of Lee & Bareinboim (2020) presents an interesting connection to our work. Given a  
469 causal graph, they study the sets of pairs (node, context(node)) (referred to as “scopes”) that may  
470 correspond to an optimal (multi-node) intervention policy where each node  $X$  in a scope is intervened  
471 on according to a policy  $\pi_X(X \mid \text{context}(X))$ . This is a challenging problem, and they do not provide  
472 a full characterization of these optimal scopes, instead deriving a set of rules that can be used to  
473 compare certain pairs of scopes. In this paper, we instead address the single-node intervention case,  
474 and assume that the problem sets the conditioning set  $\mathbf{Z}_X$  (context) to use and impose only minimal  
475 restrictions on what  $\mathbf{Z}_X$  can be, focusing instead on choosing the nodes that can yield the best results.

476 To conclude, in this paper we introduced the conditional causal bandit problem, where the agent only  
477 has knowledge of the causal graph  $G$ , the arms are conditional interventions, and the reward variable  
478 belongs to  $G$ . The theoretical contributions include a rigorous, simple graphical characterization  
479 of the minimal set of nodes which is guaranteed to contain the node with the optimal conditional  
480 intervention, and the C4 algorithm, which computes this set in linear time. Empirical results validate  
481 that our approach substantially prunes the search space in both real-world and sparse randomly-  
482 generated graphs. Furthermore, integrating mGISS with a UCB-based conditional bandits algorithm  
483 showcased improved cumulative regret curves. While Lee & Bareinboim (2020) consider multi-node  
484 interventions, it would be interesting in future work to adapt their ideas to the single-node case  
485 to identify the smallest  $\mathbf{Z}_X$  sets for which the best policy can still be found. Addressing latent  
confounding would also require substantially more research and is thus left as future work. On the  
practical side, instead of combining C4 with the simple CondIntUCB, one could replace CondIntUCB

---

486 with any other conditional bandit algorithm that leverages the model’s causal structure. As discussed  
487 earlier in this section, no such algorithm currently exists. Nevertheless, we expect that combining C4  
488 with any future algorithm for causal bandits with conditional interventions will be advantageous, as it  
489 reduces the number of arms that need to be considered.

490  
491

## REPRODUCIBILITY STATEMENT

492  
493  
494  
495  
496  
497

All experiments and results described in Section 6 can be reproduced using the code in the repository submitted alongside the paper. The experiments are simple to run, and instructions are included in the repository itself. All theoretical results are proved in the appendix.

## REFERENCES

498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In *International Conference on Machine Learning*, pp. 1638–1646. PMLR, 2014.

Georgij P Basharin. On a statistical estimate for the entropy of a sequence of independent random variables. *Theory of Probability & Its Applications*, 4(3):333–336, 1959.

Michael A Bender, Martin Farach-Colton, Giridhar Pemmasani, Steven Skiena, and Pavel Sumazin. Lowest common ancestors in trees and directed acyclic graphs. *Journal of Algorithms*, 57(2): 75–94, 2005.

Léon Bottou, Jonas Peters, Joaquin Quiñero-Candela, Denis X Charles, D Max Chickering, Elon Portugaly, Dipankar Ray, Patrice Simard, and Ed Snelson. Counterfactual reasoning and learning systems: The example of computational advertising. *Journal of Machine Learning Research*, 14 (11), 2013.

Miroslav Dudik, Daniel Hsu, Satyen Kale, Nikos Karampatziakis, John Langford, Lev Reyzin, and Tong Zhang. Efficient optimal learning for contextual bandits. *arXiv preprint arXiv:1106.2369*, 2011.

P. Erdős and A. Rényi. On random graphs. i. *Publicationes Mathematicae*, 6(3–4):290–297, 1959.

Shi Feng and Wei Chen. Combinatorial causal bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 7550–7558, 2023.

Rolf Jagerman, Ilya Markov, and Maarten De Rijke. Safe exploration for optimizing contextual bandits. *ACM Transactions on Information Systems (TOIS)*, 38(3):1–23, 2020.

Michael Joffe, Manoj Gambhir, Marc Chadeau-Hyam, and Paolo Vineis. Causal diagrams in systems epidemiology. *Emerging themes in epidemiology*, 9:1–18, 2012.

John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. *Advances in Neural Information Processing Systems*, 20, 2007.

Finnian Lattimore, Tor Lattimore, and Mark D Reid. Causal bandits: Learning good interventions via causal inference. *Advances in Neural Information Processing Systems*, 29, 2016.

Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.

Sanghack Lee and Elias Bareinboim. Structural causal bandits: Where to intervene? *Advances in Neural Information Processing Systems*, 31, 2018.

Sanghack Lee and Elias Bareinboim. Structural causal bandits with non-manipulable variables. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 4164–4172, 2019.

Sanghack Lee and Elias Bareinboim. Characterizing optimal mixed policies: Where to intervene and what to observe. *Advances in Neural Information Processing Systems*, 33:8565–8576, 2020.

---

540 Yangyi Lu, Amirhossein Meisami, Ambuj Tewari, and William Yan. Regret analysis of bandit prob-  
541 lems with causal background knowledge. In *Conference on Uncertainty in Artificial Intelligence*,  
542 pp. 141–150. PMLR, 2020.

543

544 Rahul Madhavan, Aurghya Maiti, Gaurav Sinha, and Siddharth Barman. Causal contextual bandits  
545 with adaptive context. *arXiv preprint arXiv:2405.18626*, 2024.

546

547 Aurghya Maiti, Vineet Nair, and Gaurav Sinha. A causal bandit approach to learning good atomic  
548 interventions in presence of unobserved confounders. In *Uncertainty in Artificial Intelligence*, pp.  
549 1328–1338. PMLR, 2022.

550

551 Nicolai Meinshausen, Alain Hauser, Joris M Mooij, Jonas Peters, Philip Versteeg, and Peter Bühlmann.  
552 Methods for causal inference from gene perturbation experiments and validation. *Proceedings of  
the National Academy of Sciences*, 113(27):7361–7368, 2016.

553

554 Tom Michoel and Jitao David Zhang. Causal inference in drug discovery and development. *Drug  
discovery today*, 28(10):103737, 2023.

555

556 Vineet Nair, Vishakha Patil, and Gaurav Sinha. Budgeted and non-budgeted causal bandits. In  
557 *International Conference on Artificial Intelligence and Statistics*, pp. 2017–2025. PMLR, 2021.

558

559 Judea Pearl. A probabilistic calculus of actions. In *Uncertainty in Artificial Intelligence*, pp. 454–462.  
Elsevier, 1994.

560

561 Judea Pearl. *Causality*. Cambridge university press, 2009.

562

563 Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: Foundations  
and learning algorithms*. The MIT Press, 2017.

564

565 Ayush Sawarni, Rahul Madhavan, Gaurav Sinha, and Siddharth Barman. Learning good interventions  
566 in causal graphs via covering. In *Uncertainty in Artificial Intelligence*, pp. 1827–1836. PMLR,  
567 2023.

568

569 Rajat Sen, Karthikeyan Shanmugam, Alexandros G Dimakis, and Sanjay Shakkottai. Identifying  
570 best interventions through online importance sampling. In *International Conference on Machine  
Learning*, pp. 3057–3066. PMLR, 2017.

571

572 Ilya Shpitser and Judea Pearl. Identification of conditional interventional distributions. *arXiv preprint  
arXiv:1206.6876*, 2012.

573

574 Chandrasekar Subramanian and Balaraman Ravindran. Causal contextual bandits with targeted  
575 interventions. In *International Conference on Learning Representations*, 2022.

576

577 Chandrasekar Subramanian and Balaraman Ravindran. Causal contextual bandits with one-shot data  
578 integration. *Frontiers in Artificial Intelligence*, 7:1346700, 2024.

579

580 Akihiro Yabe, Daisuke Hatano, Hanna Sumita, Shinji Ito, Naonori Kakimura, Takuro Fukunaga, and  
581 Ken-ichi Kawarabayashi. Causal bandits with propagating inference. In *International Conference  
on Machine Learning*, pp. 5512–5520. PMLR, 2018.

582

583 Yan Zhao, Mitchell Goodman, Sameer Kanase, Shenghe Xu, Yannick Kimmel, Brent Payne, Saad  
584 Khan, and Patricia Grao. Mitigating targeting bias in content recommendation with causal bandits’.  
585 In *Proc. ACM Conference on Recommender Systems Workshop on Multi-Objective Recommender  
Systems*, Seattle, WA., 2022.

586

587

588

589

590

591

592

593

## A DIRECTED ACYCLIC GRAPHS

All graphs in this paper are directed acyclic graphs (DAGs). Every path is assumed to be directed. A path  $\pi$  in a graph  $G = (\mathbf{V}, E)$  is a tuple of nodes such that each node  $X$  in the path has an outgoing arrow from  $X$  to the next node in the tuple<sup>13</sup>. For  $X \in \mathbf{V}$ , we denote by  $\text{Pa}(X)$ ,  $\text{Ch}(X)$ ,  $\text{De}(X)$  and  $\text{An}(X)$  the sets of parents, children, descendants and ancestors of  $X$ , respectively. We denote by  $\pi: X \dashrightarrow Y$  a path starting at node  $X$  and ending at node  $Y$ , and  $\overset{\circ}{\pi}$  denotes the path formed by the inner nodes of  $\pi$ . By abuse of notation, we often perform set operations such as  $\pi_1 \cap \pi_2$  between paths, which implicitly means that these operations are performed on the sets of nodes belonging to the paths. Tuples with a single node are also considered to be paths, and are said to be *trivial*. Also, if  $B \in \pi: X \dashrightarrow Y$ , then the paths  $\pi|_Z: Z \dashrightarrow Y$  and  $\pi|_Z: X \dashrightarrow Z$  are the paths resulting from removing from  $\pi$  all nodes before and after  $Z$ , respectively. Every node is an ancestor of itself, so that the relation  $\preceq$  defined by  $X \preceq Y \iff Y \in \text{An}(X)$  is a partial order. Given a set  $\mathbf{U}$  of nodes, we denote by  $\max_{\preceq}[\mathbf{U}]$  the set of maximal elements of  $\mathbf{U}$  with respect to  $\preceq$ . We call this the *ancestor partial order*. If there is a non-trivial path from  $X$  to  $Y$ , then  $Y$  is said to be *reachable* from  $X$ . The set of common ancestors of nodes  $X$  and  $Y$  is denoted  $\text{CA}(X, Y) = \text{An}(X) \cap \text{An}(Y) = \{Z \in \mathbf{V} : Z \preceq X \wedge Z \preceq Y\}$ . Finally, the *degree* of a node in a DAG is the sum of the incoming and outgoing arrows of that node.

We also make use of a lesser-known graph theory concept, relevant for this paper: the “lowest common ancestors” of nodes  $(X, Y)$ . These are common ancestors that don’t reach any other common ancestors, intuitively making them the “closest” to  $(X, Y)$ .

**Definition 17** (Lowest Common Ancestors in a DAG (Bender et al., 2005)). *Let  $X, Y$  be nodes of a DAG  $G = (\mathbf{V}, E)$ . A lowest common ancestor (LCA) of  $X$  and  $Y$  is a minimal element of  $\text{CA}(X, Y)$  with respect to the ancestor partial order  $\preceq$ . The set of all lowest common ancestors of  $X$  and  $Y$  is denoted  $\text{LCA}(X, Y)$ .*

For example, in Figure 1a,  $\text{LCA}(A_1, A_2) = \{X_1\}$ , whereas in Figure 1b,  $\text{LCA}(A_1, A_2) = \{A_1\}$ .

## B THE KIDNEY FUNCTION EXAMPLE

Recall the kidney function example discussed in Section 2. The variables  $\text{Weight}_N$ ,  $\text{BPN}$  and  $\text{RenalFlow}_N$  are the weight, blood pressure, and renal blood flow of the patient at the end of week  $N$  (equivalently, at the start of week  $N + 1$ ). All are measured at the end of each week. The doctor can intervene on one of these variables *using the measured values as context for the intervention*, in order to optimize the kidney function of the patient at the end of week 3 (Kidney3). We model this situation with the causal graph depicted in Figure 4. Making use of Theorem 13, we see that the minimal set of nodes which needs to be tested in this case is  $\{\text{RenalFlow}_2, \text{Weight}_2, \text{Weight}_1, \text{Weight}_0\}$ .

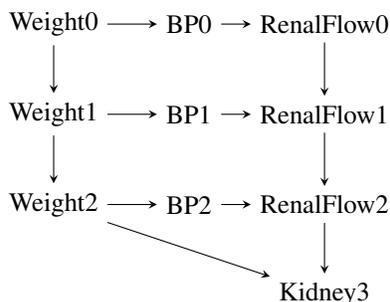


Figure 4: Causal graph for the kidney function example from Section 2. The doctor can intervene on any node  $X$  except  $\text{Kidney}_3$ , making use of the measured variables  $\mathbf{Z}_X$  until (and including) that week, thus including in particular  $\text{An}(X) \setminus \{X\}$ .

<sup>13</sup>Since all DAGs we are considering in this paper come from SCMs, there is at most one arrow between any two nodes, so that a tuple of nodes is enough to define a path. For a general graph one would have to specify a list of edges.

---

## C UNROLLED ASSIGNMENTS

The structural assignments of an SCM can be utilized to express any endogenous variable as a function of the exogenous variables only. This is achieved by composing the assignments until reaching the exogenous variables. Our proofs will rely on these functions, which we will refer to as “unrolled assignments”, since we “unroll” the expressions for the endogenous variables until only exogenous variables are left. We define them formally by induction as follows:

**Definition 18** (Unrolled Assignment). *We define the unrolled assignment  $\bar{f}_X: R_{\mathbf{N}} \rightarrow R_X$  of any (exogenous or endogenous) variable  $X$  from an SCM  $\mathfrak{C} = (\mathbf{V}, \mathbf{N}, \mathcal{F}, p_{\mathbf{N}})$  by induction. For  $X = N_i \in \mathbf{N}$ , define  $\bar{f}_X(\mathbf{n}) := n_i$ . Now, let  $\trianglelefteq$  be a topological order on  $G$  where the first elements are the endogenous variables with no endogenous parents. Let  $S$  be the poset  $(\mathbf{V}, \trianglelefteq)$ . In ascending order, take  $X \in S$ , and define:*

$$\bar{f}_X(\mathbf{n}) := \begin{cases} f_X(n_X), & \text{if } \text{Pa}(X) = \emptyset \\ f_X(\bar{f}_{\text{Pa}(X)}(\mathbf{n}), n_X), & \text{otherwise} \end{cases}, \quad (6)$$

where  $\bar{f}_{\text{Pa}(X)}(\mathbf{n}) = (\bar{f}_{\text{Pa}(X)_1}(\mathbf{n}), \dots, \bar{f}_{\text{Pa}(X)_{m_X}}(\mathbf{n}))$  and  $m_X = |\text{Pa}(X)|$ .

Additionally, we can consider  $X$  as a function of both exogenous variables and a chosen endogenous variable  $B$ . To achieve this, we substitute the assignments until we reach either  $B$  or the exogenous variables, thereby “unrolling” the dependencies until we reach the exogenous variables or we are blocked by  $B$ .

**Definition 19** (Blocked Unrolled Assignment). *Let  $X, B$  endogenous variables from an SCM  $\mathfrak{C} = (\mathbf{V}, \mathbf{N}, \mathcal{F}, p_{\mathbf{N}})$ . We define the unrolled assignment  $\bar{f}_X[B]: R_B \times R_{\mathbf{N}} \rightarrow R_X$  of  $X$  blocked by  $B$  by induction. Let  $S$  be the poset from Definition 18. In ascending order, take  $X \in S$ , and define:*

$$\bar{f}_X[B](B, \mathbf{n}) := \begin{cases} \bar{f}_X(\mathbf{n}), & \text{if } X \notin \text{De}(B) \\ B, & \text{if } X = B \\ f_X(\bar{f}_{\text{Pa}(X)}[B](B, \mathbf{n}), n_X) & \text{otherwise} \end{cases}, \quad (7)$$

where  $\bar{f}_{\text{Pa}(X)}[B](\mathbf{n}) = (\bar{f}_{\text{Pa}(X)_1}[B](B, \mathbf{n}), \dots, \bar{f}_{\text{Pa}(X)_{m_X}}[B](B, \mathbf{n}))$  and  $m_X = |\text{Pa}(X)|$ .

*Remark 20.* Strictly speaking,  $\bar{f}_X$  is not a function of all the values of all the noise variables, but only of the exogenous variables  $N_W$  associated with endogenous variables  $W$  that  $Y$  depends on. Similarly,  $\bar{f}_X[B]$  is also not a function of all the values of all the noise variables. Namely, if  $X$  only depends on an endogenous variable  $W$  through  $B$ , then  $n_W$  will never appear in the expression for  $\bar{f}_X[B]$ , and the same holds in case  $B = W$ . A more accurate notation would reflect these facts, writing the unrolled assignments as functions of the specific noise variables that can affect them, rather than as functions of all noise variables. We opted not to adopt this notation to avoid complicating the notation and conceptual simplicity of these quantities.

The following lemma relates blocked unrolled assignments with atomic interventions, and will be used to prove Theorem 13.

**Lemma 21.** *Let  $X \in \mathbf{V}$  and  $Y \in \mathbf{V} \cup \mathbf{N}$ . Then  $\bar{f}_Y[X](x, \mathbf{n}) = \bar{f}_Y^{do(X=x)}(\mathbf{n})$ .*

*Proof.* Let  $X$  be an endogenous variable. We want to prove that the expression holds for any variable  $Y$ . We will prove this by induction. Let  $\trianglelefteq$  be a topological order on the nodes of  $G^*$ . Note that the first elements with respect to this order are the exogenous variables, i.e.  $N \trianglelefteq Z$  whenever  $N \in \mathbf{N}$  and  $Z \in \mathbf{V}$ . The result is true for the exogenous variables. Indeed, for  $Y \in \mathbf{N}$  we have that  $\bar{f}_Y[X](x, \mathbf{n}) = \bar{f}_Y(\mathbf{n}) = Y = \bar{f}_Y^{do(X=x)}(\mathbf{n})$ , since  $Y \notin \text{De}(X) \cup \{X\}$  and  $Y$  is exogenous (both in the pre- and post-intervention structural causal models). This establishes the base case of the induction. Now let  $Y$  be endogenous. For the inductive step, we will prove that, if the result is true for the parents  $\text{Pa}_{G^*}(Y)$  of  $Y$  in  $G^*$  (induction hypothesis), then it is also true for  $Y$ . Assume the antecedent (induction hypothesis). There are three possibilities:  $Y \in \text{De}(X) \setminus \{X\}$ ,  $Y = X$  or

702  $Y \notin \text{De}(X)$ . In case  $Y \in \text{De}(X) \setminus \{X\}$ :

$$\begin{aligned}
703 \quad \bar{f}_Y[X](x, \mathbf{n}) &= f_Y(\bar{f}_{\text{Pa}(Y)}[X](x, \mathbf{n}), n_Y). \\
704 & \\
705 \quad &\stackrel{\text{I.H.}}{=} f_Y(\bar{f}_{\text{Pa}(Y)}^{\text{do}(X=x)}(\mathbf{n}), n_Y) \\
706 & \\
707 \quad &= f_Y^{\text{do}(X=x)}(\bar{f}_{\text{Pa}(Y)}^{\text{do}(X=x)}(\mathbf{n}), n_Y) \\
708 & \\
709 \quad &\stackrel{\text{def}}{=} \bar{f}_Y^{\text{do}(X=x)}(\mathbf{n}),
\end{aligned} \tag{8}$$

710 where in the third equality we used that  $f_Y^{\text{do}(X=x)} = f_Y$ . If instead  $Y = X$ , then one simply has  
711  $\bar{f}_Y[X](x, \mathbf{n}) = \bar{f}_X[X](x, \mathbf{n}) \stackrel{\text{def}}{=} x$ . Furthermore,  $\bar{f}_Y^{\text{do}(X=x)}(\mathbf{n}) = \bar{f}_X^{\text{do}(X=x)}(\mathbf{n}) = f_X^{\text{do}(X=x)}(\mathbf{n}) =$   
712  $x$ , where the second equality holds simply because  $X$  has no non-exogenous parents in the  
713 post-intervention graph. Finally, if  $Y \notin \text{De}(X)$ , then  $\bar{f}_Y[X](x, \mathbf{n}) = \bar{f}_Y(\mathbf{n})$  by definition. And  
714  $\bar{f}_Y^{\text{do}(X=x)}(\mathbf{n}) = f_Y^{\text{do}(X=x)}(\bar{f}_{\text{Pa}(Y)}^{\text{do}(X=x)}(\mathbf{n}), n_Y) = f_Y(\bar{f}_{\text{Pa}(Y)}(\mathbf{n}), n_Y)$ , where in the last equality we  
715 used that  $X \notin \text{An}(Y) \Rightarrow \bar{f}_{\text{Pa}(Y)}^{\text{do}(X=x)}(\mathbf{n}) = \bar{f}_{\text{Pa}(Y)}(\mathbf{n})$ . This establishes the inductive step: if the  
716 results holds for the first  $j \geq |\mathbf{N}|$  variables with respect to  $\preceq$ , then it also holds for the variable  $j + 1$ ,  
717 since its parents are among the first  $j$  variables.  $\square$

719 The following lemma shows how one can chain (blocked) unrolled assignments when there is a node  
720  $Z$  present in all paths from the blocking node  $B$  to  $Y$ . This result is consistent with the intuition that,  
721 if all paths from  $B$  to  $Y$  must go through  $Z$ , then knowing the value of  $Z$  is enough to compute  $Y$ .

722 **Lemma 22.** *If all paths from  $B$  to  $Y$  must include  $Z$ , then  $\bar{f}_Y[B](b, \mathbf{n}) = \bar{f}_Y[Z](\bar{f}_Z[B](b, \mathbf{n}), \mathbf{n})$ .*

724 *Proof.* Let  $S$  be the poset whose elements are all the descendants  $A$  of  $B$  for which all paths from  $B$   
725 to  $A$  must go through  $Z$ , and the partial order is a topological order  $\preceq$ . Denote the elements of  $S$  by  
726  $W_i$ , where  $i \in \{0, \dots, m-1\}$  corresponds to the position of  $W_i$  in the order  $\preceq$ . We will prove the  
727 result by induction on a topological order. Notice that  $Y \in S$ . Thus, we can just show the result for  
728 all  $W_i$ . We start with the base case  $W_0$ . By definition:

$$729 \quad \bar{f}_{W_0}[Z](\bar{f}_Z[B](b, \mathbf{n}), \mathbf{n}) = f_{W_0}(\bar{f}_{\text{Pa}(W_0)}[Z](\bar{f}_Z[B](b, \mathbf{n}), \mathbf{n}), n_{W_0}). \tag{9}$$

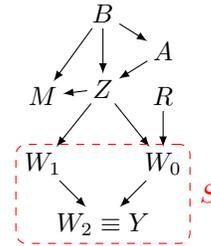
730 Recall that  $\text{Pa}(W_0) = (\text{Pa}(W_0)_1, \dots, \text{Pa}(W_0)_{m_0})$ . Hence, we want to check that  
731  $\bar{f}_{\text{Pa}(W_0)_i}[Z](\bar{f}_Z[B](b, \mathbf{n}), \mathbf{n}) = \bar{f}_{\text{Pa}(W_0)_i}[B](b, \mathbf{n})$  for all  $i$ , since in that case the right hand side  
732 of Equation (9) becomes  $f_{W_0}(\bar{f}_{\text{Pa}(W_0)_i}[B](b, \mathbf{n}), n_{W_0}) \stackrel{\text{def}}{=} \bar{f}_{W_0}[B](b, \mathbf{n})$ .

733 If  $\text{Pa}(W_0)_i = Z$ , then by definition of blocked unrolled assignment  $\bar{f}_{\text{Pa}(W_0)_i}[Z](\bar{f}_Z[B](b, \mathbf{n}), \mathbf{n}) =$   
734  $\bar{f}_Z[Z](\bar{f}_Z[B](b, \mathbf{n}), \mathbf{n}) = \bar{f}_Z[B](b, \mathbf{n})$ .

736 If  $\text{Pa}(W_0)_i \neq Z$ , then  $\text{Pa}(W_0)_i$  cannot be a descendant of  $B$ . Indeed,  
737  $W_0$  must have no parent that is a descendant of  $B$ , except maybe for  
738  $Z$ . That is:  $\text{Pa}(W_0) \cap \text{De}(B) \subseteq \{Z\}$ . Otherwise, either that parent  
739 would be in  $S$  and thus equal to  $W_k$  for some  $k > 0$ , or it would be  
740 in  $\text{De}(B) \setminus (S \cup \{Z\})$ , so that there would be a path  $B \dashrightarrow W_0$  not  
741 crossing  $Z$  — both cases contradict the definition of  $W_0$ . Hence, we  
742 only need to consider the case where  $\text{Pa}(W_0)_i \notin \text{De}(B)$ . In particular,  
743  $\text{Pa}(W_0)_i \notin \text{De}(Z)$ . Then:

$$744 \quad \bar{f}_{\text{Pa}(W_0)_i}[Z](\bar{f}_Z[B](b, \mathbf{n}), \mathbf{n}) = \bar{f}_{\text{Pa}(W_0)_i}(\mathbf{n}) = \bar{f}_{\text{Pa}(W_0)_i}[B](b, \mathbf{n}). \tag{10}$$

745 This shows the result for the base case  $W_0$ . Now, assume it to be true for all  $W_j$  with  $j \leq k$  (induction  
746 hypothesis). Equation (9) still holds for  $W_{k+1}$ . Now, each parent  $\text{Pa}(W_{k+1})_i$  must either be equal  
747 to  $W_j$  for some  $j < k + 1$ , or not a descendant of  $B$  (for the same reason as for the parents of  $W_0$ ).  
748 In the latter case, Equation (10) still holds for  $\text{Pa}(W_{k+1})_i$ . Hence, we only need to check that, for  
749  $\text{Pa}(W_{k+1})_i = W_j$  (with  $j < k + 1$ ), we have that  $\bar{f}_{W_j}[Z](\bar{f}_Z[B](b, \mathbf{n}), \mathbf{n}) = \bar{f}_{W_j}[B](b, \mathbf{n})$ . But this  
750 is just the induction hypothesis.  $\square$



## 752 D CONDITIONAL SUPERIORITY VS DETERMINISTIC ATOMIC SUPERIORITY

754 We will show that conditional intervention superiority is equivalent to deterministic atomic interven-  
755 tion superiority. This result will help prove results about the former by making use of the former,  
which is mathematically simpler and easier to reason about.

756 *Notation.* We denote by  $G^*$  the graph resulting from adding to a causal graph  $G$  the exogenous  
757 variables as nodes, and an edge  $N_{X_i} \rightarrow X_i$  for each exogenous variable  $N_{X_i}$ .  
758

759 **Lemma 23** (Conditional Intervention vs Atomic Intervention). *Let  $A$  be a set of endogenous variables  
760 of an SCM  $\mathfrak{C}$  and let  $X, Y$  be endogenous variables of  $\mathfrak{C}$  not in  $A$ . When evaluated at a setting  $\mathbf{n}$ , the  
761 unrolled assignment of  $Y$  after a conditional intervention  $do(X = g(A))$  coincides with the unrolled  
762 assignment of  $Y$  after the atomic intervention  $do(X = g(\bar{f}_A(\mathbf{n})))$ . That is:*

$$763 \bar{f}_Y^{do(X=g(A))}(\mathbf{n}) = \bar{f}_Y^{do(X=g(\bar{f}_A(\mathbf{n})))}(\mathbf{n}).$$

766 *Proof.* This result can be proved by induction in a similar way to Lemma 21.

767 Let  $X$  be an endogenous variable. We want to prove that the expression holds for any variable  $Y$ . We  
768 will prove this by induction on a topological order  $\triangleleft$  on the nodes of  $G^*$  such that the first elements  
769 are precisely the exogenous variables, *i.e.*  $N \triangleleft Z$  whenever  $N \in \mathbf{N}$  and  $Z \in \mathbf{V}$ .

770 The result is true for the exogenous variables. Indeed, for  $Y \in \mathbf{N}$ , and making use of Lemma 21,  
771 we have that  $\bar{f}_Y^{do(X=g(\bar{f}_A(\mathbf{n})))}(\mathbf{n}) = \bar{f}_Y[X](g(\bar{f}_A(\mathbf{n})), \mathbf{n}) = \bar{f}_Y(\mathbf{n}) = Y = \bar{f}_Y^{do(X=g(A))}(\mathbf{n})$ , since  
772  $Y \notin \text{De}(X) \cup \{X\}$  and  $Y$  is exogenous (both in the pre- and post-intervention (both conditional and  
773 atomic) structural causal models). This establishes the base case of the induction.

774 Now let  $Y$  be endogenous. For the inductive step, we will prove that, if the result is true for the  
775 parents  $\text{Pa}_{G^*}(Y)$  of  $Y$  in  $G^*$  (induction hypothesis), then it is also true for  $Y$ . Assume the antecedent  
776 (induction hypothesis). There are three possibilities:  $Y \in \text{De}(X) \setminus \{X\}$ ,  $Y = X$  or  $Y \notin \text{De}(X)$ . In  
777 case  $Y \in \text{De}(X) \setminus \{X\}$ :

$$778 \begin{aligned} \bar{f}_Y^{do(X=g(\bar{f}_A(\mathbf{n})))}(\mathbf{n}) &\stackrel{\text{def}}{=} f_Y^{do(X=g(\bar{f}_A(\mathbf{n})))}(\bar{f}_{\text{Pa}(Y)}^{do(X=g(\bar{f}_A(\mathbf{n})))}(\mathbf{n}), n_Y) \\ &= f_Y(\bar{f}_{\text{Pa}(Y)}^{do(X=g(\bar{f}_A(\mathbf{n})))}(\mathbf{n}), n_Y) \\ &\stackrel{\text{I.H.}}{=} f_Y(\bar{f}_{\text{Pa}(Y)}^{do(X=g(A))}(\mathbf{n}), n_Y) \\ &= f_Y^{do(X=g(A))}(\bar{f}_{\text{Pa}(Y)}^{do(X=g(A))}(\mathbf{n}), n_Y) \\ &\stackrel{\text{def}}{=} \bar{f}_Y^{do(X=g(A))}(\mathbf{n}), \end{aligned} \quad (11)$$

787 where in the second and fourth equalities we used that  $f_Y^{do(X=g(\bar{f}_A(\mathbf{n})))} = f_Y = f_Y^{do(X=g(A))}$ . We  
788 also used that  $\text{Pa}(Y)$  is unchanged by these interventions. If instead  $Y = X$ , then one has:

$$789 \begin{aligned} \bar{f}_X^{do(X=g(A))}(\mathbf{n}) &\stackrel{\text{def}}{=} f_X^{do(X=g(A))}(\bar{f}_{\text{Pa}_{G^*}^{do(X=g(A))}(X)}^{do(X=g(A))}(\mathbf{n}), n_X) \\ &= f_X^{do(X=g(A))}(\bar{f}_A^{do(X=g(A))}(\mathbf{n}), n_X) \\ &= g(\bar{f}_A(\mathbf{n}), n_X), \end{aligned} \quad (12)$$

794 and also:

$$795 \begin{aligned} \bar{f}_X^{do(X=g(\bar{f}_A(\mathbf{n})))}(\mathbf{n}) &\stackrel{\text{def}}{=} f_X^{do(X=g(\bar{f}_A(\mathbf{n})))}(\bar{f}_{\text{Pa}_{G^*}^{do(X=g(\bar{f}_A(\mathbf{n})))}(X)}^{do(X=g(\bar{f}_A(\mathbf{n})))}(\mathbf{n}), n_X) \\ &= g(\bar{f}_A(\mathbf{n}), n_X). \end{aligned} \quad (13)$$

800 Finally, if  $Y \notin \text{De}(X)$ , then trivially  $\bar{f}_Y^{do(X=g(A))}(\mathbf{n}) = \bar{f}_Y(\mathbf{n})$  and  $\bar{f}_Y^{do(X=g(\bar{f}_A(\mathbf{n})))}(\mathbf{n}) = \bar{f}_Y(\mathbf{n})$ .

801 This establishes the inductive step: if the results holds for the first  $j \geq |\mathbf{N}|$  variables with respect to  
802  $\triangleleft$ , then it also holds for the variable  $j + 1$ , since its parents are among the first  $j$  variables.  $\square$   
803

804 **Lemma 24** (Superiority and Paths). *If  $X \succeq_Y^{\text{det}, a} W$ , then all paths  $W \dashrightarrow Y$  must include  $X$ .*  
805

806 *Proof.* If  $W \notin \text{An}(Y)$ , there are no paths from  $W$  to  $Y$  and the conclusion is vacuously true. We  
807 assume from now on that  $W \in \text{An}(Y)$ . Assume, for the sake of contradiction, that there is a path  
808  $\pi: W \dashrightarrow A \rightarrow Y$  in  $G$  without  $X$ , where  $A$  is a parent of  $Y$ . Let  $\text{pr}_\pi$  denote the operator that, given  
809 a node  $X$  in the path  $\pi$  (where  $X$  is different from  $W$ ), outputs the node that precedes  $X$  in that path.

810 Consider the SCM with graph  $G$  and structural assignments and noise distributions given by:

$$\begin{cases}
f_Y(A, \text{Pa}(Y) \setminus A, N_Y) = 2A + N_Y \cdot \mathbf{1}_{>0}(\sum_{Z \in \text{Pa}(Y) \setminus A} Z) \\
f_{C \in \pi \setminus W}(\text{Pa}(C), N_C) = \text{pr}_\pi(C) + N_C \cdot \mathbf{1}_{>0}(\sum_{Z \in \text{Pa}(C) \setminus \text{pr}_\pi(C)} Z) \\
f_W(\text{Pa}(W), N_W) = N_W \cdot \mathbf{1}_{>0}(\sum_{Z \in \text{Pa}(W)} Z) \\
f_{V \notin \pi}(\text{Pa}(V), N_V) = N_V \cdot \mathbf{1}_{>0}(\sum_{Z \in \text{Pa}(V)} Z) \\
N_V \sim \text{Ber}(\frac{1}{2})
\end{cases},$$

818 where  $\mathbf{1}_{>0}: \mathbb{R} \rightarrow \{0, 1\}$  is the unit step function, which maps values larger than 0 to 1, and all  
819 non-positive values to 0. Then,  $\bar{f}_Y^{do(W=1)}(\mathbf{0}) = 2\bar{f}_A^{do(W=1)}(\mathbf{0}) = 2$ , while, for every  $Z \in \mathbf{V} \setminus \pi$   
820 and  $z \in R_Z$ , we have  $\bar{f}_Y^{do(Z=z)}(\mathbf{0}) = 0$ . Since  $X$  is not in  $\pi$ , then in particular  $\bar{f}_Y^{do(Z=z)}(\mathbf{0}) = 0$   
821 for every  $x \in R_X$ . That is, for the setting  $\mathbf{n} = \mathbf{0}$ , there is no intervention on  $X$  that is better than  
822  $do(W = 1)$ , which contradicts the antecedent.  $\square$

825 We will also need a lemma similar to Lemma 30 but for conditional-intervention superiority. Its proof  
826 is similar to that of Lemma 30.

827 **Lemma 25.** *If  $X_1 \in \mathbf{V} \setminus \{Y\}$  and  $X_2 \notin \text{An}(Y)$ , then  $X_1 \succeq_Y^c X_2$ .*

828 *Proof.* For any SCM  $\mathfrak{C}$ , intervening on  $X_2 \notin \text{An}(Y)$  will give  $Y = \bar{f}_Y(\mathbf{n})$ . When intervening on  
829  $X_1$ , we can simply set it to  $\bar{f}_{X_1}(\mathbf{n})$  (the observational value for  $X_1$ ) to obtain the same value of  $Y$   
830 (and possibly we can do better). Notice that this is a (trivial) instance of conditional intervention.  
831 Thus  $X_1 \succeq_Y^c X_2$ .  $\square$

832 **Proposition 4** (Conditional vs Atomic superiority). *Let  $X, W, Y$  be nodes in a DAG  $G$ . Then*  
833  *$X$  is conditional-intervention superior to  $W$  relative to  $Y$  in  $G$  if and only if  $X$  is deterministic*  
834 *atomic-intervention superior to  $W$  relative to  $Y$  in  $G$ . That is,  $X \succeq_Y^c W \Leftrightarrow X \succeq_Y^{\text{det}, a} W$ .*

835 *Proof.* ( $\Rightarrow$ ): Assume  $X \succeq_Y^c W$ . Let  $\mathfrak{C} = (\mathbf{V}, \mathbf{N}, \mathcal{F}, p_{\mathbf{N}})$  be an SCM with causal graph  $G$   
836 and  $\mathbf{m} \in R_{\mathbf{N}}$ . Let  $g^* = \arg \max_g \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{do(X=g(\mathbf{Z}_X))}(\mathbf{n})$ . Then,  $\forall h, \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{do(X=g^*(\mathbf{Z}_X))}(\mathbf{n}) \geq$   
837  $\mathbb{E}_{\mathbf{n}} \bar{f}_Y^{do(W=h(\mathbf{Z}_W))}(\mathbf{n})$ . This holds in particular for  $p_{\mathbf{N}} = \delta(\mathbf{m})$ . Denoting by  $\mathcal{F}(\mathbf{A}, \mathbf{B})$  the set of  
838 functions with domain  $\mathbf{A}$  and codomain  $\mathbf{B}$ , we can then write:

$$\forall h \in \mathcal{F}(R_{\mathbf{Z}_W}, R_W), \bar{f}_Y^{do(X=g^*(\bar{f}_{\mathbf{Z}_X}(\mathbf{m})))}(\mathbf{m}) \geq \bar{f}_Y^{do(W=h(\bar{f}_{\mathbf{Z}_W}(\mathbf{m})))}(\mathbf{m}),$$

839 where we also used Lemma 23. Now, since every  $w \in R_W$  can be attained from  $\bar{f}_{\mathbf{Z}_W}(\mathbf{m})$  by simply  
840 choosing  $h$  to be the constant function which is always equal to  $w$ , then choosing  $x^* = g^*(\bar{f}_{\mathbf{Z}_X}(\mathbf{m}))$   
841 allows us to write:

$$\forall w \in R_W, \bar{f}_Y^{do(X=x^*)}(\mathbf{m}) \geq \bar{f}_Y^{do(W=w)}(\mathbf{m}).$$

842 This proves that  $X \succeq_Y^{\text{det}, a} W$ .

843  $\square \Rightarrow$

844 ( $\Leftarrow$ ): Assume now that  $X \succeq_Y^{\text{det}, a} W$ . If  $W \notin \text{An}(Y)$ , the result follows immediately from Lemma 25.  
845 Assume henceforth that  $W \in \text{An}(Y)$ . Let  $p_{\mathbf{N}} \in \mathcal{P}(\mathbf{N})$  and  $\mathcal{F}(G) = \{f_V: V \in G\}$ . We want to  
846 show that  $\max_g \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{do(X=g(\mathbf{Z}_X))}(\mathbf{n}) \geq \max_h \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{do(W=h(\mathbf{Z}_W))}(\mathbf{n})$ . From Lemma 23, we can write  
847 this as  $\max_g \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{do(X=g(\bar{f}_{\mathbf{Z}_X}(\mathbf{n})))}(\mathbf{n}) \geq \max_h \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{do(W=h(\bar{f}_{\mathbf{Z}_W}(\mathbf{n})))}(\mathbf{n})$ . Denote the expected value  
848 on the left-hand-side by  $\alpha(g)$ , and the one on the right-hand-side by  $\beta(h)$ . Assume, for the sake of  
849 contradiction, that there is  $h^*$  such that  $\beta(h^*) > \alpha(g)$  for all  $g$ . Define  $H(\mathbf{n}) = h^*(\bar{f}_{\mathbf{Z}_W}(\mathbf{n}))$ . Now,  
850 if  $W \notin \text{An}(Y)$ , we simply define  $g^*$  to output the observational value of  $X$ . If instead  $W \in \text{An}(Y)$ ,  
851 from Lemma 24, we know that  $X \in \text{De}(W)$  and all paths from  $W$  to  $Y$  go through  $X$ . We then

864 define<sup>14</sup>  $g^*(\bar{f}_{Z_X}(\mathbf{n})) = \bar{f}_X[W](h^*(\bar{f}_{Z_W}(\mathbf{n})), \mathbf{n})$ . Let  $G(\mathbf{n}) = g^*(\bar{f}_{Z_X}(\mathbf{n}))$ . Then:

$$\begin{aligned}
866 \quad \alpha(g^*) &= \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{\text{do}(X=G(\mathbf{n}))}(\mathbf{n}) \\
867 \quad &= \mathbb{E}_{\mathbf{n}} \bar{f}_Y[X](G(\mathbf{n}), \mathbf{n}) \\
868 \quad &= \mathbb{E}_{\mathbf{n}} \bar{f}_Y[X](\bar{f}_X[W](H(\mathbf{n}), \mathbf{n}), \mathbf{n}) \\
869 \quad &= \mathbb{E}_{\mathbf{n}} \bar{f}_Y[W](H(\mathbf{n}), \mathbf{n}) \\
870 \quad &= \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{\text{do}(W=H(\mathbf{n}))}(\mathbf{n}) \\
871 \quad &= \mathbb{E}_{\mathbf{n}} \bar{f}_Y^{\text{do}(W=H(\mathbf{n}))}(\mathbf{n}) \\
872 \quad &= \beta(h^*). \\
873 \quad &
\end{aligned}$$

874 where in the fourth equality we used Lemma 22, and in the fifth we used Lemma 21. This contradicts  
875 our assumption.

876  $\square \Leftarrow$

877  $\square$

878 As mentioned in the main text (Section 3), the superiority relation for atomic interventions in non-  
879 deterministic (general) SCMs defined in the natural way is *not* equivalent to  $\succeq_Y^c$ . Indeed, consider  
880 the following example:

881 *Example 26.* Consider the SCM given by  $Y = A \oplus W$ ,  $A = Z \oplus W$  and  $N_Z, N_W \sim \text{Bern}(1/2)$ ,  
882 where  $\oplus$  is the XOR operator and all variables are binary. Setting  $Z$  to 1 ensures that  $Y = 1$ ,  
883 so that  $\mathbb{E}_{\mathbf{n}} \bar{f}_Y^{\text{do}(Z=1)} = 1$ . No atomic intervention on  $A$  would accomplish this:  $\mathbb{E}_{\mathbf{n}} \bar{f}_Y^{\text{do}(A=0)} =$   
884  $\mathbb{E}_{\mathbf{n}} \bar{f}_Y^{\text{do}(A=1)} = \frac{1}{2}$ . Hence  $A \not\succeq_Y^a Z$ . However,  $\mathbb{E}_{\mathbf{n}} \bar{f}_Y^{\text{do}(A=g(W))} = 1 = \max R_Y$  if one uses the policy  
885  $g(0) = 1, g(1) = 0$ . Thus  $A \succeq_Y^c Z$ .  
886  
887  
888

## 889 E INTERVENTION SUPERIORITY RELATIONS ARE PREORDERS

890 It is straightforward to show that both interventional superiority relations are in fact preorders, and are  
891 partial orders if restricted to  $\text{An}(Y)$ .

892 **Proposition 27.** *The interventional superiority relation between nodes is a preorder in  $G$ . The*  
893 *interventional superiority relation between node sets is also a preorder.*

894 *Proof.* Let  $G$  be a DAG and let  $Y \in G$ . We will first prove the result for the interventional superiority  
895 relation on nodes.

896 **Reflexivity:** Let  $X$  be a node in  $G$  and  $\mathfrak{C} \in \mathfrak{C}(G)$ . For each setting  $n$ , the largest value of  $Y$  that  
897 can be achieved by intervening on  $X$  is attained when setting  $X$  to  $x^*(n) = \arg \max_x \bar{f}_Y^{\text{do}(X=x)}(n)$ .

898 Hence,  $\bar{f}_Y^{\text{do}(X=x^*(n))}(n) \geq \bar{f}_Y^{\text{do}(X=x)}(n)$  for all  $x \in R_X$ , so that  $X \succeq_Y^{\text{det},a} X$ .

899 **Transitivity:** assume that  $Z \succeq_Y^{\text{det},a} W$  and  $W \succeq_Y^{\text{det},a} X$ . Let  $\mathfrak{C} \in \mathfrak{C}(G)$  and  $n \in R_N$ . Then  
900  $\max_x \bar{f}_Y^{\text{do}(X=x)}(n) \leq \max_w \bar{f}_Y^{\text{do}(W=w)}(n) \leq \max_z \bar{f}_Y^{\text{do}(Z=z)}(n)$ . Hence  $Z \succeq_Y^{\text{det},a} X$ .

901 This establishes that  $\succeq_Y^{\text{det},a}$  is a preorder in  $G$ . We now show the result for node sets. Let  $\mathbf{X}, \mathbf{W}$  and  
902  $\mathbf{Z}$  be sets of nodes in  $G$ .

903 **Reflexivity:** let  $X \in \mathbf{X}$ . Since, by reflexivity of  $\succeq_Y^{\text{det},a}$  on nodes, we have that  $X \succeq_Y^{\text{det},a} X$ , it trivially  
904 follows that  $\mathbf{X} \succeq_Y^{\text{det},a} \mathbf{X}$ .

905 **Transitivity:** assume that  $\mathbf{Z} \succeq_Y^{\text{det},a} \mathbf{W}$  and  $\mathbf{W} \succeq_Y^{\text{det},a} \mathbf{X}$ . Let  $X \in \mathbf{X}$ . Then there is  $W \in \mathbf{W}$  such  
906 that  $W \succeq_Y^{\text{det},a} X$ . There is also  $Z \in \mathbf{Z}$  such that  $Z \succeq_Y^{\text{det},a} W$ . By transitivity of  $\succeq_Y^{\text{det},a}$  on nodes, it  
907 follows that  $Z \succeq_Y^{\text{det},a} X$ . Hence  $\mathbf{Z} \succeq_Y^{\text{det},a} \mathbf{X}$ .  $\square$

908 *Remark 28* (Interventional Superiority is not a partial order, and it is not total). One may have  
909 expected interventional superiority (both on nodes and on node sets) to be a partial order in  $G$ .  
910 However, they are merely preorders. That is, the antisymmetry property does not hold (unless, as  
911 shown in Appendix E.1, we restrict ourselves to the ancestors of  $Y$ ). To see this for  $\succeq_Y^{\text{det},a}$  on nodes,  
912 just notice that, if  $X, W \notin \text{An}(Y)$ , then trivially  $X \succeq_Y^{\text{det},a} W$  and  $W \succeq_Y^{\text{det},a} X$ , no matter what  $X$   
913  
914  
915  
916  
917

<sup>14</sup>Notice that, since  $Z_W \subseteq Z_X$ , this is well defined.

and  $W$  are. For node sets, consider the case where  $\mathbf{X} \subsetneq \mathbf{W}$ , but the best intervention lies in  $\mathbf{X}$ . Then  $\mathbf{X} \succeq_Y^{\text{det},a} \mathbf{W}$  and  $\mathbf{W} \succeq_Y^{\text{det},a} \mathbf{X}$ , even though  $\mathbf{X} \neq \mathbf{W}$ .

Notice also that  $\succeq_Y^{\text{det},a}$  on nodes cannot be a total preorder: just consider the graph  $A_1 \rightarrow Y \leftarrow A_2$ . One can have an SCM  $\mathcal{C}$  in which intervening on  $A_1$  can lead to larger values of  $Y$  than interventions on  $A_2$ . But one can also switch the structural assignments assignments of  $\mathcal{C}$ , which would lead to the opposite conclusion. This example also shows that  $\succeq_Y^{\text{det},a}$  on node sets also cannot be a total preorder.

## E.1 ... AND ARE PARTIAL ORDERS IN $\text{An}(Y)$

**Lemma 29** (Antisymmetry holds in  $\text{An}(Y)$ ). *For  $X_1 \in \text{An}(Y) \setminus \{Y\}$  with  $\pi_1$  a directed path from  $X_1$  to  $Y$  and  $X_2 \in \mathbf{V} \setminus \{Y\}$  with  $X_2 \succeq_Y^{\text{det},a} X_1$ , we have  $X_2 \in \pi_1$ .*

*Proof.* Take  $\mathcal{C}_1$  to be the SCM in which the nodes  $\mathbf{V}$  are binary variables, with structural assignments

$$V_i = \begin{cases} P_i & \text{if } P_i \rightarrow V_i \text{ is on the path } \pi_1, \\ 0 & \text{otherwise.} \end{cases}$$

Then intervening on  $X_1$  (by setting it to 1) or another node in  $\pi_1$  will give  $Y = 1$ , but intervening on a node not in  $\pi_1$  gives  $Y = 0$ . So  $X_2 \succeq_Y^{\text{det},a} X_1$  implies  $X_2 \in \pi_1$ .  $\square$

**Lemma 30** (Everything beats Non-Ancestors). *If  $X_1 \in \mathbf{V} \setminus \{Y\}$  and  $X_2 \notin \text{An}(Y)$ , then  $X_1 \succeq_Y^{\text{det},a} X_2$ .*

*Proof.* For any SCM  $\mathcal{C}$ , intervening on  $X_2 \notin \text{An}(Y)$  will give  $Y = \bar{f}_Y(\mathbf{n})$ . When intervening on  $X_1$ , we can set it to  $\bar{f}_{X_1}(\mathbf{n})$  (the observational value for  $X_1$ ) to obtain the same value of  $Y$  (and possibly we can do better). Thus  $X_1 \succeq_Y^{\text{det},a} X_2$ .  $\square$

**Lemma 31** (Non-Ancestors do not beat Ancestors). *If  $X_1 \in \text{An}(Y) \setminus \{Y\}$  and  $X_2 \notin \text{An}(Y)$ , then  $X_2 \not\succeq_Y^{\text{det},a} X_1$ .*

*Proof.*  $X_2$  is not on any directed path from  $X_1$  to  $Y$ , so it follows from Lemma 29 that  $X_2 \not\succeq_Y^{\text{det},a} X_1$ .  $\square$

**Proposition 32** (Antisymmetry). *For  $X_1 \in \text{An}(Y) \setminus \{Y\}$  and  $X_2 \in \mathbf{V} \setminus \{Y\}$ , if  $X_1 \succeq_Y^{\text{det},a} X_2$  and  $X_2 \succeq_Y^{\text{det},a} X_1$  then  $X_1 = X_2$ .*

*Proof.* Pick a directed path  $\pi_1$  from  $X_1$  to  $Y$ . By Lemma 29,  $X_2 \in \pi_1$ . Thus  $X_2 \in \text{An}(Y)$ . Pick  $\pi_2 = \pi_1|^{X_2}$ . Now by an analogous argument, we find that  $X_1 \in \pi_2$ , which implies  $X_1 = X_2$ .  $\square$

Notice that Proposition 32 is actually slightly stronger than antisymmetry within  $\text{An}(Y)$ , since  $X_1 \succeq_Y^{\text{det},a} X_2$  and  $X_2 \succeq_Y^{\text{det},a} X_1$  can only occur for distinct  $X_1, X_2$  if *both* are outside  $\text{An}(Y)$ .

## F PROOFS FOR THE MINIMAL GLOBALLY INTERVENTIONALLY SUPERIOR SET

### F.1 UNIQUENESS OF THE MGISS

**Lemma 33** (Elements of an mGISS are not Comparable). *Let  $\mathbf{A} \subseteq \mathbf{V}$  be an mGISS relative to  $Y$ . Let  $X, X' \in \mathbf{A}$  and  $X \neq X'$ . Then  $X' \not\succeq_Y^{\text{det},a} X$ .*

*Proof.* Assume  $X' \succeq_Y^{\text{det},a} X$  for the sake of contradiction. We will show that this implies that  $\mathbf{A} \setminus X$  is also a GISS. That is, that for every element of  $(\mathbf{V} \setminus Y) \setminus (\mathbf{A} \setminus X)$  there is an element of  $\mathbf{A} \setminus X$  which is superior to it. Let  $W \in (\mathbf{V} \setminus Y) \setminus (\mathbf{A} \setminus X)$ . If  $W = X$ , then  $X' \in \mathbf{A} \setminus X$  and  $X' \succeq_Y^{\text{det},a} X$ . If  $W \neq X$ , then  $W \in (\mathbf{V} \setminus Y) \setminus \mathbf{A}$ . Since  $\mathbf{A}$  is a GISS, we can pick  $\tilde{X} \in \mathbf{A}$  such that  $\tilde{X} \succeq_Y^{\text{det},a} W$ . In case  $\tilde{X} = X$ , we can choose instead  $X'$ . Indeed, since  $X' \succeq_Y^{\text{det},a} X$  and  $X \succeq_Y^{\text{det},a} W$ , we have by transitivity of  $\succeq_Y^{\text{det},a}$  (Proposition 27) that  $X' \succeq_Y^{\text{det},a} W$ . This shows that  $\mathbf{A} \setminus X \subseteq \mathbf{A}$  is a GISS, contradicting the minimality of  $\mathbf{A}$ .  $\square$

**Lemma 34.** *Let  $G$  be a DAG and  $Y$  a node of  $G$  with at least one parent. Then any minimal globally interventionally superior set of  $G$  relative to  $Y$  is a nonempty subset of  $\text{An}(Y)$ .*

*Proof.* Let  $\mathbf{U}$  be a minimal globally interventionally superior set of  $G$  with respect to  $Y$ . Suppose  $\mathbf{U} \cap \text{An}(Y) = \emptyset$ . Then  $((\mathbf{V} \setminus \{Y\}) \setminus \mathbf{U}) \cap \text{An}(Y) \neq \emptyset$ , and by Lemma 31, no element of  $\mathbf{U}$  is interventionally superior to those: contradiction. So  $\mathbf{U} \cap \text{An}(Y) \neq \emptyset$ .

Next suppose  $\mathbf{U} \not\subseteq \text{An}(Y)$ . Then  $\mathbf{U} \cap \text{An}(Y)$  is also a GISS: any of its elements is interventionally superior to  $\mathbf{U} \setminus \text{An}(Y)$  by Lemma 30. So again we have a contradiction, and we conclude that if  $\mathbf{U}$  is minimal,  $\mathbf{U} \subseteq \text{An}(Y)$ .  $\square$

**Proposition 6** (Uniqueness of the mGISS). *Let  $G$  be a DAG and  $Y$  a node of  $G$  with at least one parent. The minimal globally interventionally superior set of  $G$  relative to  $Y$  is unique. We denote it by  $\text{mGISS}_Y(G)$ .*

*Proof.* Let  $\mathbf{A}$  and  $\mathbf{B}$  be minimal globally interventionally superior sets of  $G$  with respect to  $Y$ . Assume, for the sake of contradiction, that  $\mathbf{B} \neq \mathbf{A}$ . By minimality of  $\mathbf{A}$ , we have  $\mathbf{B} \not\subseteq \mathbf{A}$ , so that  $\mathbf{B} \setminus \mathbf{A} \neq \emptyset$ . Let  $X \in \mathbf{B} \setminus \mathbf{A}$ . In particular,  $X \in (\mathbf{V} \setminus Y) \setminus \mathbf{A}$ . Hence,  $\exists Z \in \mathbf{A}$  s.t.  $Z \succeq_Y^{\text{det},a} X$ . Either  $Z \in \mathbf{A} \cap \mathbf{B}$  or  $Z \in \mathbf{A} \setminus \mathbf{B}$ . If  $Z \in \mathbf{A} \setminus \mathbf{B}$ , then in particular  $Z \in (\mathbf{V} \setminus Y) \setminus \mathbf{B}$ . Since  $\mathbf{B}$  is a GISS, there is  $X' \in \mathbf{B}$  such that  $X' \succeq_Y^{\text{det},a} Z$ . We claim that  $X' \neq X$ : Suppose for a contradiction that  $X' = X$ . Then by Proposition 32,  $X, Z \notin \text{An}(Y)$ . But this contradicts Lemma 34, so we conclude  $X' \neq X$ . By transitivity of  $\succeq_Y^{\text{det},a}$  (Proposition 27), it follows that  $X' \succeq_Y^{\text{det},a} X$ . Similarly, if  $Z \in \mathbf{A} \cap \mathbf{B}$ , one again has two elements  $Z$  and  $X$  of  $\mathbf{B}$  such that  $Z \succeq_Y^{\text{det},a} X$ . In both cases, this contradicts the assumption that  $\mathbf{B}$  is a GISS, as per Lemma 33.  $\square$

## F.2 THE LSCA CLOSURE AND $\Lambda$ -STRUCTURES

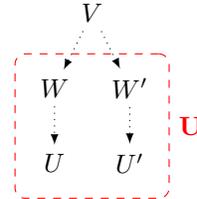
It will be useful to know that, in order to show that a node belongs to  $\mathcal{L}^\infty(\mathbf{U})$ , it suffices to prove that it belongs to the LSCA closure of a subset of  $\mathbf{U}$ . We show by induction that this is indeed the case.

**Lemma 35.** *If  $\mathbf{U}' \subseteq \mathbf{U}$ , then  $\mathcal{L}^\infty(\mathbf{U}') \subseteq \mathcal{L}^\infty(\mathbf{U})$ .  $\mathbf{U}' \subseteq \mathbf{U}$*

*Proof.* Recall that  $\mathcal{L}^\infty(\mathbf{U}) = \mathcal{L}^i(\mathbf{U})$  for some  $i \in \mathbb{N}$ . We will show the result by induction on  $i \in \mathbb{N}$ . The base case holds trivially:  $\mathcal{L}^0(\mathbf{U}') = \mathbf{U}' \subseteq \mathbf{U} = \mathcal{L}^0(\mathbf{U})$ . Now assume that  $\mathcal{L}^i(\mathbf{U}') \subseteq \mathcal{L}^i(\mathbf{U})$  for a given  $i \in \mathbb{N}$  (induction hypothesis). Let  $V \in \text{LSCA}(\mathcal{L}^i(\mathbf{U}'))$ . Then there are paths  $V \dashrightarrow X$ ,  $V \dashrightarrow Y$  with  $X, Y \in \mathcal{L}^i(\mathbf{U}')$  not containing  $Y$  and  $X$ , respectively. But  $X, Y$  are also in  $\mathcal{L}^i(\mathbf{U})$ , so that  $V \in \text{LSCA}(\mathcal{L}^i(\mathbf{U}))$ . Then  $\text{LSCA}(\mathcal{L}^i(\mathbf{U}')) \subseteq \text{LSCA}(\mathcal{L}^i(\mathbf{U}))$ . Using once more the induction hypothesis, it follows that  $\mathcal{L}^{i+1}(\mathbf{U}') = \text{LSCA}(\mathcal{L}^i(\mathbf{U}')) \cup \mathcal{L}^i(\mathbf{U}') \subseteq \text{LSCA}(\mathcal{L}^i(\mathbf{U})) \cup \mathcal{L}^i(\mathbf{U}) = \mathcal{L}^{i+1}(\mathbf{U})$ .  $\square$

**Lemma 36.** *Let  $\mathbf{U} \subseteq \mathbf{V}$ . If  $V \in \text{LSCA}(\mathbf{U}) \setminus \mathbf{U}$ , then  $V$  forms a  $\Lambda$ -structure over  $(\mathbf{U}, \mathbf{U})$ .*

*Proof.* Let  $V \in \text{LSCA}(\mathbf{U}) \setminus \mathbf{U}$ . By Definition 8, there are distinct  $U, U' \in \mathbf{U}$  for which there are paths  $\pi: V \dashrightarrow U$  and  $\pi': V \dashrightarrow U'$  whose interiors do not intersect  $\{U, U'\}$ . Now, let  $W$  (respectively  $W'$ ) be the first element in  $\pi$  (respectively  $\pi'$ ) in  $\mathbf{U}$ . Notice that  $W \neq W'$ , otherwise  $W = W'$  would be in  $\text{SCA}(U, U')$  and be reachable from  $V$ , so that  $V$  would not be a minimal element of  $\text{SCA}(U, U')$ . This would contradict  $V \in \text{LSCA}(U, U')$ . Similarly, the paths  $\pi|_W: V \dashrightarrow W$ ,  $\pi'|_{W'}: V \dashrightarrow W'$  resulting from restricting  $\pi$  cannot have interior intersections: such an intersection node  $\tilde{V}$  would be an SCA of  $U, U'$  reachable from  $V$ , so that  $V \notin \text{LSCA}(U, U')$  — again a contradiction. Therefore,  $V$  forms a  $\Lambda$ -structure over  $W, W'$ .  $\square$

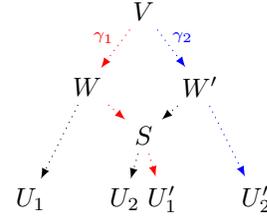


**Theorem 12** (Simple Graphical Characterization of LSCA Closure). *A node  $V \in \mathbf{V}$  is in the LSCA closure  $\mathcal{L}^\infty(\mathbf{U})$  of  $\mathbf{U} \subseteq \mathbf{V}$  if and only if  $V$  forms a  $\Lambda$ -structure over  $(\mathbf{U}, \mathbf{U})$ . I.e.  $\mathcal{L}^\infty(\mathbf{U}) = \Lambda(\mathbf{U}, \mathbf{U})$ .*

1026 *Proof.* Proof of  $\subseteq$ : If  $\mathcal{L}^\infty(\mathbf{U}) = \mathbf{U}$ , then the result is trivially true. We assume from now on that  
 1027  $\mathcal{L}^\infty(\mathbf{U}) \supseteq \mathbf{U}$ . We will prove that  $V \in \mathcal{L}^\infty(\mathbf{U}) \Rightarrow V \in \Lambda(\mathbf{U}, \mathbf{U})$  by induction with respect to  
 1028 a chosen strict reverse topological order  $<$  (*i.e.*  $V' \in \text{An}(V) \setminus \{V\} \Rightarrow V < V'$ ). The base case  
 1029 is  $V_0 \in \mathbf{U}$ , since an element of  $\mathbf{U}$  will be the first element of  $\mathcal{L}^\infty(\mathbf{U})$  for any chosen  $<$ . In this  
 1030 case, we can simply take the trivial paths  $\pi = \pi' = (V_0)$ . Then  $V_0 \in \Lambda(\mathbf{U}, \mathbf{U})$ . Now, assume  
 1031 that  $V \in \mathcal{L}^\infty(\mathbf{U}) \setminus \mathbf{U}$  and that the implication holds for all  $W \in \mathcal{L}^\infty(\mathbf{U})$  such that  $W < V$   
 1032 (induction hypothesis). Let  $W, W'$  be<sup>15</sup>distinct elements of  $\mathcal{L}^\infty(\mathbf{U})$  such that  $V \in \text{LSCA}(W, W')$ .

1033 In particular,  $W, W' < V$ . By Lemma 36 applied to  $\{W, W'\}$ , there are paths  $V \xrightarrow{\alpha} W$ ,  $V \xrightarrow{\alpha'} W'$   
 1034 intersecting only at  $V$ . Furthermore, by the induction hypothesis we have that  $W, W' \in \Lambda(\mathbf{U}, \mathbf{U})$ , so  
 1035 that there are paths  $W \xrightarrow{\pi_1} U_1$ ,  $W \xrightarrow{\pi_2} U_2$ ,  $W' \xrightarrow{\pi'_1} U'_1$ ,  $W' \xrightarrow{\pi'_2} U'_2$  such that  $U_1, U_2, U'_1, U'_2 \in \mathbf{U}$ ,  
 1036  $\pi_1 \cap \pi_2 = \{W\}$  and  $\pi'_1 \cap \pi'_2 = \{W'\}$ .  
 1037

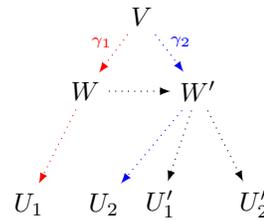
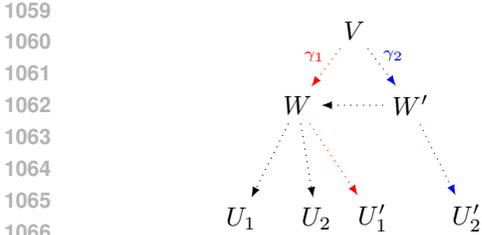
1038 Let  $\mathbf{S} = (\alpha \cup \pi_1 \cup \pi_2) \cap (\alpha' \cup \pi'_1 \cup \pi'_2)$  and  $\preceq$  be a chosen topological  
 1039 order. If  $\mathbf{S} = \emptyset$ , we can just take  $\gamma = \pi_1 \circ \alpha : V \dashrightarrow U_1$  and  
 1040  $\gamma' = \pi'_1 \circ \alpha' : V \dashrightarrow U'_1$  to form a  $\Lambda$ -structure for  $V$  over  $(\mathbf{U}, \mathbf{U})$ .  
 1041 Assume from now on that  $\mathbf{S} \neq \emptyset$ . Let  $S$  be the first element of  
 1042  $\mathbf{S}$  with respect to  $\preceq$ . Since  $\alpha \cap \alpha' = \emptyset$ , there are three options:  
 1043 either (i)  $S \in \pi_i \cap \alpha' \setminus \{W'\}$  for some  $i$ ; (ii)  $S \in \pi'_i \cap \alpha \setminus \{W\}$   
 1044 for some  $i$ ; or (iii)  $S \in \pi_i \cap \pi'_j$  for some  $i, j$ . By symmetry, we  
 1045 can restrict ourselves to the cases (i) and (iii): the argument for (i)  
 1046 will also hold for (ii). In both cases (i) and (ii) we have  $S \in \pi_i$  for  
 1047 some  $i \in \{1, 2\}$ . Without loss of generality, assume  $s \in \pi_2$ .



1047 For case (iii), assume, also without loss of generality, that  $s \in \pi'_1$ . If furthermore  $s \neq W'$ , we can  
 1048 construct the following two paths with no non-trivial intersections:

$$\begin{cases} \gamma_1 = \pi'_1|_S \circ \pi_2|_s \circ \alpha : V \dashrightarrow U'_1 \\ \gamma_2 = \pi'_2 \circ \alpha' : V \dashrightarrow U'_2 \end{cases} \quad (14)$$

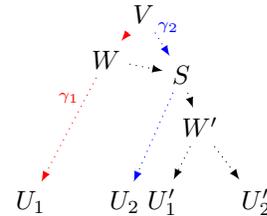
1049 To see that these paths have non non-trivial intersections, start by noticing that, by definition of  
 1050  $S$ , there is no intersection between  $\pi_2$  and  $\pi'_2$  at nodes  $A \triangleleft S$ , so that  $\pi_2|_S \cap \pi'_2 = \emptyset$ . And since  
 1051  $\pi'_1 \cap \pi'_2 = \{W'\}$  and  $S \neq W'$ , we have  $\pi'_1|_S \cap \pi'_2 = \emptyset$ . Finally,  $\pi_2 \cap \alpha' = \pi'_2 \cap \alpha = \pi_1 \cap \alpha' = \emptyset$ ,  
 1052 since otherwise there would be elements of  $S$  which are ancestors of  $s$ . Notice that this argument  
 1053 still holds if  $S = W$ , in which case  $\gamma_1$  reduces to  $\pi'_1|_W \circ \alpha$ . This shows that  $V \in \Lambda(\mathbf{U}, \mathbf{U})$  for case  
 1054 (iii), in case  $S \neq W'$ . If instead  $S = W'$ , we can simply choose paths similar to those for the case  
 1055  $S = W$  (just changing the numbers and the prime) as follows:  $\gamma_1 = \pi_1 \circ \alpha$  and  $\gamma_2 = \pi_2|^{W'} \circ \alpha'$ .



1067 We now turn to case (i), where  $S \in \pi_2 \cap \alpha' \setminus \{W\}$ . Construct the  
 1068 paths:

$$\begin{cases} \gamma_1 = \pi_1 \circ \alpha : V \dashrightarrow U_1 \\ \gamma_2 = \pi_2|_S \circ \alpha'|_s : V \dashrightarrow U_2 \end{cases} \quad (15)$$

1069 Notice that  $\alpha \cap \pi_2|_S = \emptyset$ , otherwise there would be a cycle in  
 1070 the DAG. Also,  $\pi_1 \cap \alpha'|_S = \emptyset$  by definition of  $S$ . And trivially  
 1071  $\pi_1 \cap \pi_2|_S = \emptyset$  and  $\alpha \cap \alpha' = \{V\}$ .

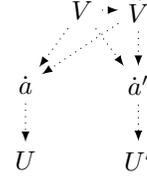


1072 It follows that  $\gamma_1$  and  $\gamma_2$  intersect only trivially, so that  $(v, \gamma_1, \gamma_2)$  forms a  $\Lambda$ -structure over  $(\mathbf{U}, \mathbf{U})$ .

1073  $\square_{\subseteq}$

1074  
 1075  
 1076  
 1077  
 1078  
 1079 <sup>15</sup>Such  $W, W'$  must exist by the definition of  $\mathcal{L}^\infty(\mathbf{U})$  whenever  $\mathcal{L}^\infty(\mathbf{U}) \supseteq \mathbf{U}$ .

1080 Proof of  $\supseteq$ : Let  $V \in \Lambda(\mathbf{U}, \mathbf{U})$ . Then, there is a pair of nodes  $U, U' \in \mathbf{U}$  over which  $V$   
1081 forms  $\Lambda$ -structures. We are going to show that  $V \in \mathcal{L}^\infty(\{U, U'\})$ . Let  $\mathbf{L}$  be the set of all the  
1082  $\Lambda$ -structures  $\lambda_i = (V, \pi_i: V \dashrightarrow U, \pi'_i: V \dashrightarrow U')$  over  $(U, U')$ . Let  $A$  be the set of nodes  
1083 in  $\mathcal{L}^\infty(\{U, U'\})$  which belong to some  $\pi_i$ . Formally,  $A = \{a \in \mathcal{L}^\infty(\{U, U'\}): \exists i \text{ s.t. } \lambda_i \in$   
1084  $\mathbf{L}, a \in \pi_i\} \setminus \{V\}$ . Let  $\dot{a}$  be the first element of  $A$  with respect to a chosen topological  
1085 order  $\preceq$ . Denote by  $\Pi'(\dot{a})$  the set of paths  $\pi'_i$  belonging to some  $\Lambda$ -structure  $(V, \pi'_i, \pi_i)$  in  $\mathbf{L}$   
1086 such that  $\pi_i$  contains  $\dot{a}$ . Let  $A'(\dot{a}) = \{a' \in \mathcal{L}^\infty(\{U, U'\}): \exists \pi'_i \in \Pi'(\dot{a}) \text{ s.t. } a' \in \pi'_i\}$ .  
1087 Furthermore, let  $\dot{a}'$  be the first element of  $A'(\dot{a})$  with respect to  
1088  $\preceq$ . Denote by  $(V, \hat{\pi}, \hat{\pi}')$  a  $\Lambda$ -structure of  $\mathbf{L}$  such that  $a \in \hat{\pi}$  and  
1089  $a' \in \hat{\pi}'$ . Notice that  $\dot{a} \neq \dot{a}'$  and  $\hat{\pi}|_{\dot{a}} \cap \hat{\pi}'|_{\dot{a}'} = \{V\}$ , by definition of  
1090  $\Lambda$ -structure. In particular,  $v \in \text{SCA}(\dot{a}, \dot{a}')$ . Suppose, for the sake of  
1091 contradiction, that there is  $\tilde{V} \in \text{SCA}(\dot{a}, \dot{a}')$  such that  $\tilde{V}$  is reachable  
1092 from  $V$ . Then there is  $\lambda = (V, \gamma, \gamma')$  in  $\mathbf{L}$  such that  $\tilde{V}, \dot{a} \in \gamma$ . But  
1093  $\tilde{V} \preceq \dot{a}$ , contradicting minimality of  $\dot{a}$ . Hence  $V \in \text{LSCA}(\dot{a}, \dot{a}')$ .  
1094 Finally, since  $\dot{a}, \dot{a}' \in \mathcal{L}^\infty(\mathbf{U})$ , it follows that  $V \in \mathcal{L}^\infty(\mathbf{U})$ .



1095  $\square_{\supseteq}$

1096  $\square$

### 1098 F.3 THE LSCA CLOSURE IS THE MGISS

1100 **Lemma 37.** *Let  $B \in \mathbf{V}$ . Assume there are nodes  $Z, W$  which are reachable from  $B$  with paths  
1101 whose interiors do not intersect  $\mathcal{L}^\infty(\{Z, W\})$ . Then  $B \in \mathcal{L}^\infty(\{Z, W\})$ .*

1103 *Proof.* Let  $\mathbf{B}$  be the intersection of the two paths. Note that all elements of  $\mathbf{B}$  are comparable in the  
1104 ancestor partial order. Take  $B'$  to be the least element of  $\mathbf{B}$ . Then  $B'$  forms a Lambda structure over  
1105  $\{Z, W\}$ , so by Theorem 12,  $B'$  is in  $\mathcal{L}^\infty(\{Z, W\})$ . Now suppose  $B' \neq B$ : then this contradicts that  
1106 the interiors of the paths are not in  $\mathcal{L}^\infty(\{Z, W\})$ . So  $B = B' \in \mathcal{L}^\infty(\{Z, W\})$ .  $\square$

1107 **Lemma 38.** *Let  $B \in \text{An}(Y)$  and  $B \notin \mathcal{L}^\infty(\text{Pa}(Y))$ . Then there is exactly one node  $Z \in \mathcal{L}^\infty(\text{Pa}(Y))$   
1108 reachable from  $B$  by paths whose interiors do not contain elements from  $\mathcal{L}^\infty(\text{Pa}(Y))$ .*

1110 *Proof.* There must be at least one node in  $\mathcal{L}^\infty(\text{Pa}(Y))$  reachable from  $B$  by paths not containing  
1111 interior elements from  $\mathcal{L}^\infty(\text{Pa}(Y))$ : since  $\text{Pa}(Y) \subseteq \mathcal{L}^\infty(\text{Pa}(Y))$  and  $B \in \text{An}(Y)$ , there are paths  
1112 from  $B$  to  $Y$  crossing  $\mathcal{L}^\infty(\text{Pa}(Y))$  (in fact, paths must at least intersect  $\text{Pa}(Y)$ ). Choose one such  
1113 path  $\pi: B \dashrightarrow Y$ . Let  $Z$  be the first element of  $\mathcal{L}^\infty(\text{Pa}(Y))$  in  $\pi$ . Then the path  $\pi|_Z: B \dashrightarrow Z$   
1114 obtained from  $\pi$  by truncating it at  $Z$  has no interior nodes in the closure  $\mathcal{L}^\infty(\text{Pa}(Y))$ . Furthermore,  
1115 if there would be a second path from  $B$  to  $W \in \mathcal{L}^\infty(\text{Pa}(Y)) \setminus \{Z\}$  containing no interior nodes  
1116 from the closure, then, by Lemma 37,  $B$  would be in  $\mathcal{L}^\infty(\{Z, W\})$  and thus in  $\mathcal{L}^\infty(\text{Pa}(Y))$  —  
1117 contradiction. This establishes uniqueness.  $\square$

1118 **Corollary 39.** *Under the assumptions of Lemma 38, all directed paths from  $B$  to  $Y$  must go through  
1119  $Z$ .*

1121 *Proof.* If there was a path from  $B$  to  $Y$  not containing  $Z$ , it would have to go through a parent  $A$  of  $Y$ .  
1122 But the first element of  $\mathcal{L}^\infty(\text{Pa}(Y))$  in this path (perhaps  $A$  itself) would contradict the uniqueness  
1123 of  $Z$  from Lemma 38.  $\square$

1125 To facilitate the exposition of the proof in Theorem 13, we introduce the concept of superiority *in a*  
1126 *given SCM*. It is simply a version of Definition 2 with no universal quantifier over  $\mathcal{C}$ . This definition  
1127 can be extended for sets of nodes in the obvious way (in an identical manner to how Definition 3  
1128 extended Definition 1 and Definition 2).

1129 **Definition 40** (Superiority in an SCM). *Let  $\mathcal{C}$  be an SCM with causal graph  $G$ . Let  $X, W, Y$  be  
1130 nodes in  $G$ .  $X$  is (deterministically) atomic-intervention superior to  $W$  relative to  $Y$  in  $G$ , denoted  
1131  $X \succeq_{G,Y}^{\text{det},a} W$ , if for every unit  $\mathbf{n}$  there is  $x \in R_X$  such that no atomic intervention on  $W$  results in a  
1132*

1133 <sup>16</sup>Notice that  $A'$  (and  $A$ ) are not empty (at least one of  $\{U, U'\}$  is in  $A'$  (and  $A$ )).

larger  $Y$  than the value of  $Y$  resulting from setting  $X = x$ . That is, Equation (2) holds (in  $\mathfrak{C}$ ) for all  $\mathbf{n} \in R_{\mathbb{N}}$ .

**Theorem 13** (Superiority of the LSCA Closure). *Let  $G$  be a causal graph and  $Y$  a node of  $G$  with at least one parent. Then, the LSCA closure  $\mathcal{L}^\infty(\text{Pa}(Y))$  of the parents of  $Y$  is the minimal globally interventionally superior set  $\text{mGISS}_Y(G)$  of  $G$  relative to  $Y$ .*

*Proof.* We need to prove two results:

- (i)  $\mathcal{L}^\infty(\text{Pa}(Y))$  is globally interventionally superior with respect to  $Y$ . That is:  $\mathcal{L}^\infty(\text{Pa}(Y)) \succeq_Y^{\text{det},a} \mathbf{V} \setminus (\mathcal{L}^\infty(\text{Pa}(Y)) \cup \{Y\})$ .
- (ii) Furthermore, this is the minimal set with this property. Namely, any proper subset  $\mathbf{I}$  of  $\mathcal{L}^\infty(\text{Pa}(Y))$  would not be interventionally superior to  $\mathbf{V} \setminus (\mathbf{I} \cup \{Y\})$ .

*Proof of (i):* Let  $B \in \mathbf{V} \setminus \mathcal{L}^\infty(\text{Pa}(Y))$  and  $B \neq Y$ . We want to show that there is  $A$  in the closure  $\mathcal{L}^\infty(\text{Pa}(Y))$  such that  $A \succeq_Y^{\text{det},a} B$ .

If  $B$  is not an ancestor of  $Y$ , then trivially  $\bar{f}_Y^{\text{do}(B=b)}(\mathbf{n}) = \bar{f}_Y(\mathbf{n})$  for all  $\mathbf{n} \in R_{\mathbb{N}}$  and for all  $b \in R_B$ , so that in particular  $\max_{b \in R_B} \bar{f}_Y^{\text{do}(B=b)}(\mathbf{n}) = \bar{f}_Y(\mathbf{n})$ . Now, let  $A$  be a parent of  $Y$ , and  $a^* = \bar{f}_A(\mathbf{n})$  (i.e.  $a^*$  is the value that  $A$  would attain if no intervention was performed). Then, from the definition of unrolled assignment and atomic intervention,  $\bar{f}_Y^{\text{do}(A=a^*)}(\mathbf{n}) = \bar{f}_Y(\mathbf{n})$ . Thus,  $\max_{a \in R_A} \bar{f}_Y^{\text{do}(A=a)}(\mathbf{n}) \geq \bar{f}_Y(\mathbf{n}) = \max_{b \in R_B} \bar{f}_Y^{\text{do}(B=b)}(\mathbf{n})$ . That is,  $A \succeq_Y^{\text{det},a} B$ .

Assume from now on that  $B$  is an ancestor of  $Y$ . From Lemma 38 there is one and only one node  $Z \in \mathcal{L}^\infty(\text{Pa}(Y))$  reachable from  $B$  by paths not containing intermediate elements from  $\mathcal{L}^\infty(\text{Pa}(Y))$ . Let  $z^* \in \arg \max_{z \in R_Z} [\bar{f}_Y^{\text{do}(Z=z)}(\mathbf{n})]$ . Further, let  $b \in R_B$ . From Lemma 22, we have that  $\bar{f}_Y[B](b, \mathbf{n}) = \bar{f}_Y[Z](\bar{f}_Y[B](b, \mathbf{n}), \mathbf{n})$ , which of course is at most  $\bar{f}_Y[Z](z^*, \mathbf{n})$ . Finally, Lemma 21 allows us to relate this to a post-intervention unrolled assignment as  $\bar{f}_Y[Z](z^*, \mathbf{n}) = \bar{f}_Y^{\text{do}(Z=z^*)}(\mathbf{n})$ . This shows that  $\max_{b \in R_B} \bar{f}_Y^{\text{do}(B=b)}(\mathbf{n}) \leq \max_{z \in R_Z} \bar{f}_Y^{\text{do}(Z=z)}(\mathbf{n})$ , so that  $Z \succeq_Y^{\text{det},a} B$ .

□<sub>(i)</sub>

*Proof of (ii):* We want to show that, for any causal graph  $G$  and node  $Y$  from  $G$ , for any set  $\mathbf{I} \subsetneq \mathcal{L}^\infty(\text{Pa}(Y))$  there is an SCM  $\mathfrak{C}$  (with causal graph  $G$ ) such that  $\mathbf{I}$  is not interventionally superior to  $\bar{\mathbf{I}} \setminus \{Y\}$  in  $\mathfrak{C}$ , i.e.  $\mathbf{I} \not\succeq_{\mathfrak{C}, Y}^{\text{det},a} \bar{\mathbf{I}} \setminus \{Y\}$ . In other words, we want to prove that:

$$\begin{aligned} \forall \text{ DAG } G = (\mathbf{V}, E), \forall Y \in \mathbf{V}, \forall \mathbf{I} \subsetneq \mathcal{L}^\infty(\text{Pa}(Y)), \\ \exists \text{ SCM } \mathfrak{C} \text{ s.t. } G^{\mathfrak{C}} = G \text{ and } \mathbf{I} \not\succeq_{\mathfrak{C}, Y}^{\text{det},a} \bar{\mathbf{I}} \setminus \{Y\}. \end{aligned} \quad (16)$$

Let  $G$  be a DAG,  $Y$  be a node of  $G$  and  $\mathbf{I}$  a proper subset of  $\mathcal{L}^\infty(\text{Pa}(Y))$ . Take  $B$  a minimal element of  $\mathcal{L}^\infty(\text{Pa}(Y)) \setminus \mathbf{I}$ . In particular,  $B \in \bar{\mathbf{I}} \setminus \{Y\}$ . We will show that there is an SCM  $\mathfrak{C}$  in which no element of  $\mathbf{I}$  is interventionally superior to  $B$ , thus proving that  $\mathbf{I} \not\succeq_{\mathfrak{C}, Y}^{\text{det},a} \bar{\mathbf{I}} \setminus \{Y\}$ , and hence  $\mathbf{I} \not\succeq_Y^{\text{det},a} \bar{\mathbf{I}} \setminus \{Y\}$ . We will divide the proof in two cases:  $B \in \text{Pa}(Y)$  and  $B \in \mathcal{L}^\infty(\text{Pa}(Y)) \setminus \text{Pa}(Y)$ . Assume  $B \in \text{Pa}(Y)$ . We can construct an SCM with causal graph  $G$  as follows:

$$\begin{cases} f_Y(\text{Pa}(Y), N_Y) = 2B + \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(Y) \setminus \{B\}} W \right) + N_Y \\ f_B(\text{Pa}(B), N_B) = N_B \cdot \left( 1 - \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(B)} W \right) \right) \\ f_{V \neq Y, B}(\text{Pa}(V), N_V) = \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(V)} W \right) + N_V \\ N_{V \neq B} \sim \delta(0) \\ N_B \sim \text{Ber}(1/2) \end{cases}, \quad (17)$$

where all endogenous variables are binary except for  $Y$  (whose range is  $\mathbb{N}$ ), and all exogenous variables are simply zero except for  $N_B$ , which is also binary. The idea is that  $B$  has a stronger influence on  $Y$  than all the other parents of  $Y$  combined, and there are values of  $\mathbf{n}$  (namely whenever

1188  $N_B = 0)$  for which  $B$  is not influenced by other variables. We need to show that, for all  $X \in \mathbf{I}$ , there  
 1189 is  $\mathbf{n} \in R_{\mathbf{N}}$  such that

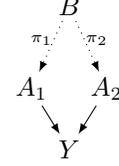
$$1190 \max_{x \in R_X} \bar{f}_Y^{do(X=x)}(\mathbf{n}) < \max_{b \in R_B} \bar{f}_Y^{do(B=b)}(\mathbf{n}). \quad (18)$$

1192 Notice that  $R_{\mathbf{N}} = \{\mathbf{0}, \mathbf{e}_{N_B}\}$ , where  $\mathbf{e}_{N_B}$  is zero everywhere except for the  $N_B$  element, which  
 1193 is 1. Let  $X \in \mathbf{I}$  and choose  $\mathbf{n} = \mathbf{0}$ . We have  $\max_{b \in \{0,1\}} \bar{f}_Y^{do(B=b)}(\mathbf{0}) = \bar{f}_Y^{do(B=1)}(\mathbf{0}) = 2 +$   
 1194  $\mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(B)} W \right) \geq 2$ . Furthermore:

$$1196 \begin{aligned} 1197 \max_{x \in \{0,1\}} \bar{f}_Y^{do(X=x)}(\mathbf{0}) &= \max_{x \in \{0,1\}} \left( 2n_B \cdot \left( 1 - \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(B)} W \right) \right) + \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(Y) \setminus \{B\}} W \right) \right) \\ 1198 &= \max_{x \in \{0,1\}} \left( 0 + \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(Y) \setminus \{B\}} W \right) \right) \\ 1200 &\leq 1 < 2 \leq \max_{b \in \{0,1\}} \bar{f}_Y^{do(B=b)}(\mathbf{0}). \end{aligned} \quad (19)$$

1207 This proves the result for  $B \in \text{Pa}(Y)$ .

1208 Assume now that  $B \in \mathcal{L}^\infty(\text{Pa}(Y)) \setminus \text{Pa}(Y)$ . From Theorem 12, there are  
 1209 nodes  $A_1, A_2 \in \text{Pa}(Y)$  which are reachable from  $B$  by paths  $\pi_1, \pi_2$  which  
 1210 only intersect at  $B$ . Denote by  $\text{pr}_i$  the operator which, given a node  $A$  in  
 1211 the path  $\pi_i$  different from  $B$ , outputs the previous node in that path. We  
 1212 construct an SCM with causal graph  $G$  as follows:



$$1213 \begin{cases} 1214 f_Y(A_1, A_2, \text{Pa}_Y \setminus \{A_1, A_2\}, N_Y) = 2A_1 \cdot A_2 + \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(Y) \setminus \{A_1, A_2\}} W \right) + N_Y \\ 1215 f_B(\text{Pa}(B), N_B) = N_B \cdot \left( 1 - \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(B)} W \right) \right) \\ 1216 f_{A \in \pi_i \setminus \{B\}}(\text{pr}_i(A), \text{Pa}(A) \setminus \{\text{pr}_i(A)\}, N_A) = \text{pr}_i(A) + N_A \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(A) \setminus \{\text{pr}_i(A)\}} W \right) \\ 1217 f_{V \notin \pi_1 \cup \pi_2 \cup \{Y\}}(\text{Pa}(V), N_V) = \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(V)} W \right) + N_V \\ 1218 N_{V \neq B} \sim \delta(0) \\ 1219 N_B, N_{A \in \pi_i} \sim \text{Ber}(1/2) \end{cases}, \quad (20)$$

1222 where again all endogenous variables except  $Y$  are binary, and all exogenous variables are zero  
 1223 except for those of the type  $N_A, A \in \pi_i$ , which is also binary. Let  $X \in \mathbf{I}$ . We again need to show  
 1224 that there is  $\mathbf{n} \in R_{\mathbf{n}}$  such that Equation (18) holds. One again chooses the setting  $\mathbf{n} = \mathbf{0}$ . The  
 1225 intuition behind this SCM is similar to that of Equation (17), with the added property that the  
 1226 elements of the paths  $\pi_i$  are simply noisy copies of  $B$ , and perfect copies when  $\mathbf{N} = \mathbf{0}$ . In particular,  
 1227  $A_i = B, i \in \{1, 2\}$ , or, using the language of unrolled assignments,  $\bar{f}_{A_i}(\mathbf{0}) = \bar{f}_B(\mathbf{0})$ . Notice that  
 1228  $\bar{f}_{A_1}(\mathbf{0}) \cdot \bar{f}_{A_2}(\mathbf{0}) = \bar{f}_B(\mathbf{0})^2 = \bar{f}_B(\mathbf{0})$ , since  $B$  is binary. These equalities still hold in the SCMs  
 1229 resulting from atomically intervening on  $B$ . Hence:

$$1230 \bar{f}_Y^{do(B=b)}(\mathbf{0}) = 2b + \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(Y) \setminus \{A_1, A_2\}} \bar{f}_W(\mathbf{0}) \right) \geq 2b. \quad (21)$$

1234 Now, if  $X = A \in \pi_1 \setminus \{B\}$ , then  $A_1$  is a perfect copy of  $A$  while  $A_2$  is still a perfect copy of  $B$ .  
 1235 Hence:

$$1236 \begin{aligned} 1237 \bar{f}_Y^{do(A=a)}(\mathbf{0}) &= 2 \underbrace{\bar{f}_{A_1}^{do(A=a)}(\mathbf{0})}_a \cdot \underbrace{\bar{f}_{A_2}^{do(A=a)}(\mathbf{0})}_0 + \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(Y) \setminus \{A_1, A_2\}} \bar{f}_W(\mathbf{0}) \right) \\ 1238 &= \mathbf{1}_{>0} \left( \sum_{W \in \text{Pa}(Y) \setminus \{A_1, A_2\}} \bar{f}_W(\mathbf{0}) \right) \leq 1 < 2 \leq \bar{f}_Y^{do(B=1)}(\mathbf{0}) \end{aligned} \quad (22)$$

The same argument holds if  $X = A \in \pi_2 \setminus \{B\}$ .  
 Finally, if instead  $X \notin \pi_1 \cup \pi_2 \cup \{Y\}$ , then  $\bar{f}_Y^{do(X=x)}(\mathbf{0}) = 0 + \mathbf{1}_{>0} \left( \sum_{X \in \text{Pa}(Y) \setminus \{A_1, A_2\}} \bar{f}_X(\mathbf{0}) \right) \leq$   
 $1 < 2 \leq \bar{f}_Y^{do(B=1)}(\mathbf{0})$ , where the first equality holds because, for  $\mathbf{n} = \mathbf{0}$ , intervening on  $X$  does not  
 affect the elements of the  $\pi_i$ , including the  $A_i$ .

□<sub>(ii)</sub>

□

## G C4 PROOFS

**Definition 41.** We define the following additional notation and terminology.

- For any set of nodes  $\mathbf{B}$  and any node  $V'$ , a path  $\pi_{V'}$  that ends in  $V'$  is uninterrupted by  $\mathbf{B}$  iff  $(\pi_{V'} \cap \mathbf{B}) \setminus \{V'\} = \emptyset$ .
- A  $\Lambda$ -structure which consists of a single node is called degenerate.
- For any set of nodes  $\mathbf{B}$ , any  $\Lambda$ -structure over  $(\mathbf{B}, \mathbf{B})$  is referred to as a  $\Lambda_{\mathbf{B}}$ -structure.
- $U \xleftarrow{\pi_U} V \xrightarrow{\pi_W} W$  denotes a  $\Lambda$ -structure  $(V, \pi_U, \pi_W)$  with paths  $\pi_U : V \dashrightarrow U$ ,  $\pi_W : V \dashrightarrow W$ . If the paths' names are not relevant or clear from the context, we write simply  $U \dashleftarrow V \dashrightarrow W$ .

**Lemma 42.** Let  $G = (\mathbf{V}, E)$  be a DAG,  $\mathbf{U} \subseteq \mathbf{V}$ .  $V \in \mathcal{L}^\infty(\mathbf{U})$  iff there exists a  $\Lambda_{\mathcal{L}^\infty(\mathbf{U})}$ -structure  $V' \dashleftarrow V \dashrightarrow V^*$  for some  $V', V^* \in \mathcal{L}^\infty(\mathbf{U})$ .

*Proof.* It is easily seen that  $\mathcal{L}^\infty(\mathcal{L}^\infty(\mathbf{U})) = \mathcal{L}^\infty(\mathbf{U})$ ; therefore, this lemma is a direct corollary of Theorem 12. □

**Lemma 43** (Existence of  $\Lambda$ -substructure). Let  $\mathbf{B}$  be a set of nodes, and let  $B_1, B_2 \in \mathbf{B}$  s.t.  $B_1 \neq B_2$ . Let  $V \notin \{B_1, B_2\}$  be a node and let  $\pi_1 : V \dashrightarrow B_1$ ,  $\pi_2 : V \dashrightarrow B_2$ , s.t.  $B_1 \notin \pi_2$  and  $B_2 \notin \pi_1$  (note that we do not assume  $\pi_1 \cap \pi_2 = \{V\}$ , meaning that other overlaps remain possible). Then, in the subgraph consisting of the two paths (as in, the graph that includes all the nodes and all the edges that are in at least one of the paths), there exists a  $\Lambda_{\mathbf{B}}$ -structure  $B_1 \dashleftarrow V' \dashrightarrow B_2$  where  $V' \in \pi_1 \cap \pi_2$ .

*Proof.* For  $V' \in \arg \min_{\prec} \pi_1 \cap \pi_2$ ,  $B_1 \xleftarrow{\pi_1|_{V'}} V' \xrightarrow{\pi_2|_{V'}} B_2$  is a  $\Lambda_{\mathbf{B}}$ -structure. □

**Lemma 15.** Let  $G = (\mathbf{V}, E)$  be a DAG,  $\mathbf{U} \subseteq \mathbf{V}$ ,  $V \in \text{An}(\mathbf{U})$ .  $c[V]$  is the unique node s.t. a path  $\pi_{c[V]} : V \dashrightarrow c[V]$  exists where  $\pi_{c[V]} \cap \mathcal{L}^\infty(\mathbf{U}) = \{c[V]\}$  (if  $V$  is its own connector, the path is trivial). This is equivalent to: for every node  $X \in \mathcal{L}^\infty(\mathbf{U})$  and path  $\pi_X : V \dashrightarrow X$ ,  $c[V]$  is the maximal element of  $\pi_X \cap \mathcal{L}^\infty(\mathbf{U})$  w.r.t. the ancestor partial order  $\preceq$ .

*Proof.* We prove the claim by induction on a reverse topological order of  $\text{An}(\mathbf{U})$ . As the base case, for  $V \in \mathbf{U}$ , definitionally  $V \in \mathcal{L}^\infty(\mathbf{U})$ , so  $c[V] = V$ , and we can just take the trivial path. Next, let  $V \in \text{An}(\mathbf{U}) \setminus \mathbf{U}$ , and note that the claim holds for all nodes in  $\text{Ch}(V) \cap \text{An}(\mathbf{U})$  by the inductive hypothesis. Let  $\mathbf{C}$  be as defined in Definition 14. There are two cases:

1. Assume that  $|\mathbf{C}| = 1$ . We can write  $\mathbf{C} = \{X\}$  for some  $X \in \mathbf{V}$ . Note that  $c[V] = X$ . By the inductive assumption,  $X$  is the unique element from  $\mathcal{L}^\infty(\mathbf{U})$  reachable from  $V$  via a non-trivial path uninterrupted by  $\mathcal{L}^\infty(\mathbf{U})$ , as any non-trivial path must go through a child, and we can apply the inductive assumption to each child. However, we still need to rule out the possibility of a trivial path to  $\mathcal{L}^\infty(\mathbf{U})$ , namely to rule out the possibility that  $V \in \mathcal{L}^\infty(\mathbf{U})$ . Since  $V \notin \mathbf{U}$ , then by Theorem 12 it is sufficient to rule out the existence of a non-degenerate  $\Lambda$ -structure from  $V$  to  $\mathbf{U}$ . However, as we noted, any non-trivial path from  $V$  to  $\mathbf{U}$  (and hence to  $\mathcal{L}^\infty(\mathbf{U})$ ) must go through  $X$ , and hence any two paths from  $V$  to  $\mathbf{U}$  must overlap at  $X \neq V$ , meaning that they do not form a  $\Lambda$ -structure.

1296 2. Assume  $|\mathbf{C}| \neq 1$ , so  $c[V] = V$ . Since  $V \in \text{An}(\mathbf{U}) \setminus \mathbf{U}$ , it must have at least one child  
1297 in  $\text{An}(\mathbf{U})$ , and that child has a connector, so  $\mathbf{C} \neq \emptyset$  and thus  $|\mathbf{C}| \geq 2$ . We claim that  
1298  $V \in \mathcal{L}^\infty(\mathbf{U})$  (and hence we can simply take the trivial path to establish the result and  
1299 complete the proof). By Theorem 12, we need to establish the existence of a  $\Lambda$ -structure  
1300 from  $V$  to  $\mathbf{U}$ . Since  $|\mathbf{C}| > 1$ , then there exist  $S_1, S_2 \in \mathbf{C}$  s.t.  $S_1 \neq S_2$ , and there exist  
1301 children  $T_1, T_2$  of  $V$  s.t.  $c[T_1] = S_1$  and  $c[T_2] = S_2$ ; by the inductive assumption,  $S_1$  and  
1302  $S_2$  are in  $\mathcal{L}^\infty(\mathbf{U})$ , and there exist paths  $\pi_1: T_1 \dashrightarrow S_1$  and  $\pi_2: T_2 \dashrightarrow S_2$  uninterrupted  
1303 by  $\mathcal{L}^\infty(\mathbf{U})$ . These paths do not overlap: had they overlapped, then by Lemma 43 they  
1304 would've contained a  $\Lambda$ -substructure  $S_1 \dashleftarrow Z \dashrightarrow S_2$  s.t.  $Z \in \pi_1 \cap \pi_2$ , so by Lemma 42  
1305  $Z \in \mathcal{L}^\infty(\mathbf{U})$ , making neither  $\pi_1$  nor  $\pi_2$  uninterrupted by  $\mathcal{L}^\infty(\mathbf{U})$ . Since  $T_1$  and  $T_2$  are  
1306 children of  $V$ , we may prepend the edges  $V \rightarrow T_1$  and  $V \rightarrow T_2$  to  $\pi_1$  and  $\pi_2$  respectively  
1307 and get paths  $\pi'_1 = V \rightarrow T_1 \dashrightarrow S_1$  and  $\pi'_2 = V \rightarrow T_2 \dashrightarrow S_2$ ; since  $\pi_1$  and  $\pi_2$  do not  
1308 overlap, these two paths yield a  $\Lambda$ -structure from  $V$  to  $\mathcal{L}^\infty(\mathbf{U})$ , which by Lemma 42 implies  
1309  $V \in \mathcal{L}^\infty(\mathbf{U})$ .  
1310 □

1311  
1312 **Theorem 16.** *C4 correctly computes  $\mathcal{L}^\infty(\mathbf{U})$ , and runs in  $O(|\mathbf{V}| + |E|)$  time.*  
1313

1314 *Proof.* Correctness is immediate from Lemma 15, as it implies  $V \in \mathcal{L}^\infty(\mathbf{U}) \Leftrightarrow c[V] = V$ . As for  
1315 the running time, assume that if the graph is not given in adjacency list representation, we convert it to  
1316 this representation in  $O(|\mathbf{V}| + |E|)$  time. Initialization in C4 is trivially  $O(|\mathbf{V}|)$ . Computing  $\text{An}(\mathbf{U})$   
1317 can be done in  $O(|\mathbf{V}| + |E|)$  time using BFS or DFS, and reverse topological sorting can be done in  
1318  $O(|\mathbf{V}| + |E|)$  using Kahn's algorithm. In the loop, for each  $v \in \text{An}(\mathbf{U}) \setminus \mathbf{U}$ , we go over all outgoing  
1319 edges from  $v$  to compute  $C$ , which because of the adjacency list representation takes  $O(|\text{Ch}(v)|)$   
1320 time. In aggregate over the entire operation of the algorithm, computing  $C$  takes  $O(|E|)$  time overall,  
1321 as each edge is inspected at most once. The loop runs  $O(|\mathbf{V}|)$  times, and all operations in it except  
1322 the computation of  $C$  take  $O(1)$  time, so all steps except computing  $C$  take at most  $O(|\mathbf{V}|)$  time  
1323 overall. Thus, the running time of the algorithm is  $O(|\mathbf{V}| + |E|)$ . □

## 1324 H SUPPLEMENTARY MATERIAL FOR EXPERIMENTAL RESULTS

1325 The results of the experiments testing our search space reduction method are presented in Figure 5  
1326 and Figure 6, for randomly generated graphs and real-world datasets, respectively.

1327 All real-world datasets come from the `bnlearn` repository, except for the `railway` dataset, which  
1328 was provided by ProRail, the institution responsible for traffic control in the Dutch railway system.  
1329 The `bnlearn` repository can be found in <https://www.bnlearn.com/bnrepository/>.  
1330 It was created by Marco Scutari as part of the `bnlearn` R package, and it is licensed under the  
1331 Creative Commons Attribution-ShareAlike 3.0 License (CC BY-SA 3.0).  
1332

1333 The `railway` dataset consists of a graph whose nodes represent train delays in a segment of the  
1334 Dutch railway system, measured at specific “points of interest” (such as train stations). Each node is  
1335 labeled with a code identifying the train, an acronym for the point of interest, a letter indicating the  
1336 train's activity at that location—arriving (A), departing (V), or passing through (D)—and the planned  
1337 time for that activity. Arrows are drawn between delay nodes that are known to influence each other.  
1338 For example, arrows connect nodes of the same train at consecutive times, as the delay of a train at  
1339 time  $t$  will influence its delay at  $t + \Delta t$ . Additionally, arrows may connect nodes corresponding to  
1340 train activities sharing the same platform, since a train must wait for the preceding train to vacate the  
1341 platform before using it. This dataset can be found in the code repository which supplements this  
1342 paper.  
1343

1344 The two search space reduction experiments – on both random and real-world graphs – were  
1345 completed in a few minutes on a CPU-based laptop with 16 GB RAM and 512 GB SSD storage.  
1346 The experiments evaluating the impact of our method on conditional intervention bandits using the  
1347 `bnlearn` models `asia`, `sachs` and `child` were also run on the same laptop. However, the most  
1348 complex model, `pathfinder`, exceeded the laptop's memory capacity due to the large number  
1349 of possible contexts generated from combinations of values assigned to ancestor nodes. Thus, the  
`pathfinder` experiment was run on an internal server with 1 TB RAM (350 GB allocated for the

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

job) and 23 TB SSD storage (our codebase used 1.5 GB). All experiments used CPU-only compute workers. The *asia*, *sachs*, and *child* models completed in approximately 12 hours on the laptop using 2 CPU cores in parallel, while the *pathfinder* experiment took about 16 hours using 27 CPU cores in parallel on the server.

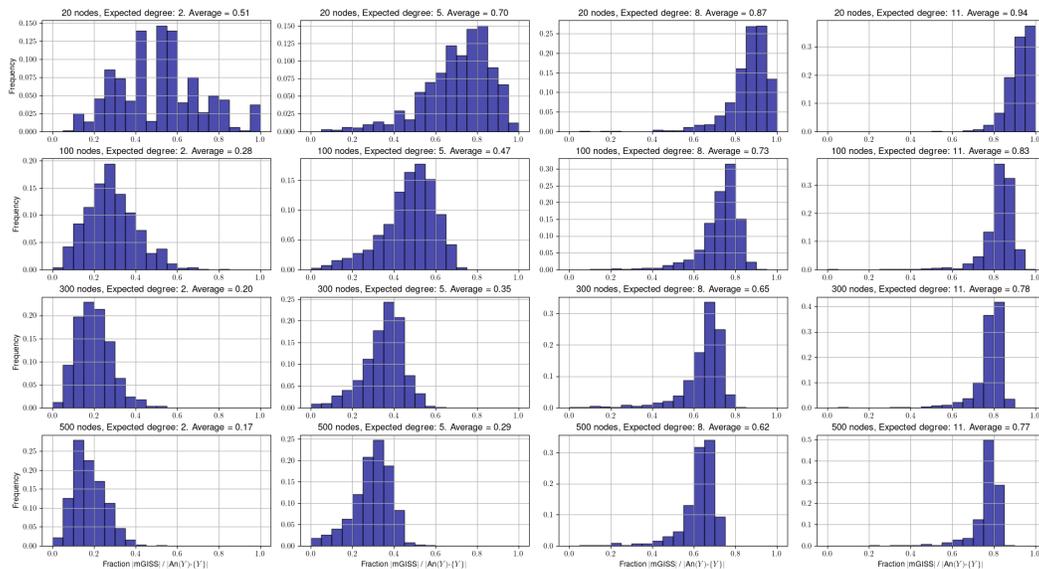


Figure 5: Fraction of nodes remaining after applying our search space filtering procedure, on random graphs. 1000 graphs were generated for each pair (number of nodes, expected degree). The impact of our method decreases with the expected degree, and increases with the number of nodes.

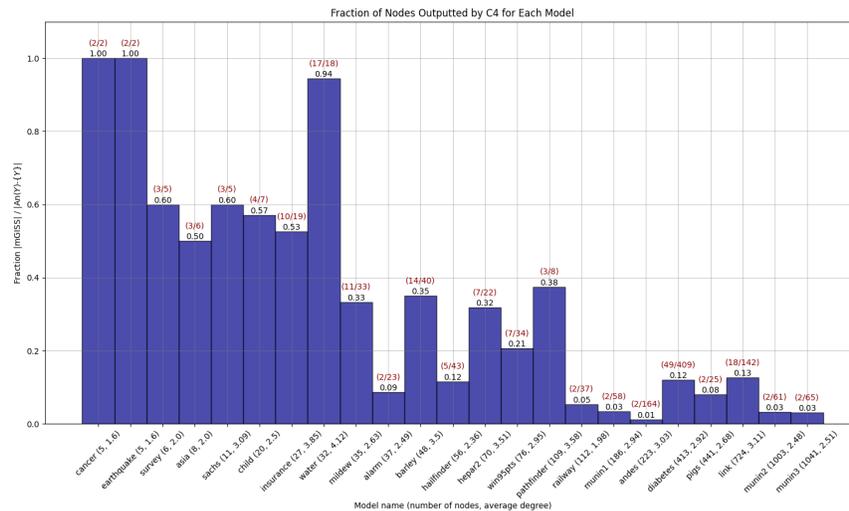


Figure 6: Fraction of nodes remaining after applying our search space filtering procedure, on real-world graphs. All models come from the *bnlearn* repository except for the railway model. The models are sorted by their *total* number of nodes. On top of each bar one can read the fraction value (in black) and the exact numbers (number of nodes in mGISS / number of proper ancestors of  $Y$ ) in red. Notice that models with larger numbers of nodes tend to benefit more from our method.

---

1404 I LLM USAGE STATEMENT  
1405

1406 We have made limited use of LLMs in preparing this paper. Specifically, they were employed to  
1407 rephrase certain complex sentences and suggest alternative wordings. No other use of LLMs was  
1408 made.  
1409

1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457