
COUNTERFACTUAL REALIZABILITY

Anonymous authors

Paper under double-blind review

ABSTRACT

It is commonly believed that, in a real-world environment, samples can only be drawn from observational and interventional distributions, corresponding to Layers 1 and 2 of the *Pearl Causal Hierarchy*. Layer 3, representing counterfactual distributions, is believed to be inaccessible by definition. However, Bareinboim, Forney, and Pearl (2015) introduced a procedure that allows an agent to sample directly from a counterfactual distribution, leaving open the question of what other counterfactual quantities can be estimated directly via physical experimentation. We resolve this by introducing a formal definition of *realizability*, the ability to draw samples from a distribution, and then developing a complete algorithm to determine whether an arbitrary counterfactual distribution is realizable given fundamental physical constraints, such as the inability to go back in time and subject the same unit to a different experimental condition. We illustrate the implications of this new framework for counterfactual data collection using motivating examples from causal fairness and causal reinforcement learning. While the baseline approach in these motivating settings typically follows an interventional or observational strategy, we show that a counterfactual strategy provably dominates both.

1 INTRODUCTION

The *Pearl Causal Hierarchy*, or PCH, is an important recent milestone in our understanding of causality (Pearl & Mackenzie, 2018; Bareinboim et al., 2022). The three layers of the PCH represent the distinct regimes of *seeing*, *doing*, and *imagining*, with regard to an environment. Consider an environment involving a decision variable X and an outcome Y . Layer 1 (\mathcal{L}_1) represents *observational* distributions, such as $P(Y | x)$. Layer 2 (\mathcal{L}_2) represents *interventional* distributions, such as $P(Y; do(x))$, using the $do()$ operator. Layer 3 (\mathcal{L}_3) represents *counterfactual* distributions **dealing with conflicting realities**, such as $P(Y_x | x', y')$: the distribution of Y had X been fixed as x , given that X, Y were in fact naturally observed to be x', y' . Higher layers subsume lower ones, but are underdetermined by them (Ibeling & Icard, 2020; Bareinboim et al., 2022).

Reasoning about \mathcal{L}_3 -quantities plays a vital role in personalized decision-making (Mueller & Pearl, 2023), analysing a causal effect into direct and indirect pathways (Pearl, 2005; Rubin, 2004), and constructing explanations for decisions, among other topics, in applications such as healthcare (Mueller & Pearl, forthcoming), economics (Li & Pearl, 2019), epidemiology (Robins & Greenland, 1992) etc. Suppose an economist were interested in estimating $P(y_x | x')$, an important \mathcal{L}_3 -quantity called the effect of the treatment on the treated, or ETT (Heckman & Robb Jr., 1985; 1986). One approach to computing such quantities is through *identification* (Pearl, 2000, §3.2.4): leveraging causal knowledge about the environment, typically a causal graph or parametric assumptions, to infer the higher-layer quantity using lower-layer data. This approach fails when the quantity is nonidentifiable, e.g. ETT in the general setting (Shpitser & Pearl, 2009; Correa et al., 2021).

However, another approach uses physical experimentation to attempt to directly draw samples from the relevant distribution, $P(Y_x, X)$ in the case of ETT, and then uses statistical methods to estimate $P(Y_x = y, X = x')$. This approach is only possible if there is some sequence of physical actions by which an agent can measure these random variables simultaneously for a single unit. It is generally believed to be feasible to draw samples only from \mathcal{L}_1 - and \mathcal{L}_2 -distributions, the latter by interventions like randomized controlled trials (RCT), à la Fisher (Fisher, 1935), and the former by simply observing the natural behaviour of the system. \mathcal{L}_3 -distributions like $P(Y_x, X)$ are deemed non-realizable in general, **as the potential response Y_x and natural decision X belong to different "worlds"**. **Once a unit naturally adopts decision $X = x'$, Y_x cannot be evaluated in the $do(x)$ regime**

054 for the same unit.¹ However, Bareinboim, Forney & Pearl have shown it is feasible to draw samples
 055 from the ETT distribution $P(Y_x, X)$ through a *counterfactual randomization* procedure (Bareinboim
 056 et al., 2015; Forney et al., 2017). This leaves open the possibility that other \mathcal{L}_3 -distributions, say
 057 perhaps $P(Y_x, X, Y)$, are also realizable through clever experimental setups, allowing one to estimate
 058 important quantities like the probability of sufficiency, $P(y_x | y', x')$ (Pearl, 1999).

059 This brings us to the central question motivating this work: *from which \mathcal{L}_3 -distributions is it possible*
 060 *to draw samples given fundamental physical constraints like the inability to travel back in time and*
 061 *subject the original unit to a different experimental condition?* We resolve this open question with a
 062 rigorous formal treatment of the *realizability of an \mathcal{L}_3 -distribution* (Def. 3.4).

063 Our main contributions in this work are as follows:
 064

- 065 • In Sec. 2 we introduce a physical procedure called *counterfactual randomization* (Def. 2.3) by
 066 which an agent can gather counterfactual data, subsuming previous similar notions.
- 067 • In Sec. 3 we develop the **CTF-REALIZE** algorithm (Algo. 1) to determine whether an \mathcal{L}_3 -
 068 distribution is physically realizable. We prove the algorithm is complete (Thm. 3.5), and derive
 069 important corollaries characterizing realizable distributions (Cors. 3.7,3.8). For instance, we show
 070 that our main result generalizes an influential notion in the causal inference literature, known as
 071 the *fundamental problem of causal inference* (Holland, 1986).
- 072 • In Sec. 4 we discuss important practical implications of counterfactual realizability. The tradi-
 073 tional route of computing \mathcal{L}_3 -quantities through identification often fails. Our work suggests
 074 opportunities for novel experiment-design ideas to directly estimate these quantities, as illustrated
 075 through Examples 1,2 and 3. More concretely,
 - 076 – In Sec. 4.1, we describe an application in causal fairness, where the naive approach of
 077 constraining a classifier using an interventional (\mathcal{L}_2) fairness metric fails to prevent disparities
 078 in outcomes across groups, but where a counterfactual (\mathcal{L}_3) approach works.
 - 079 – In Sec. 4.2, we show how counterfactual randomization can be used to improve RL algo-
 080 rithms. The baseline approach in a multi-arm bandit setting is to use allocation procedures
 081 (e.g., UCB, EXP3, Thompson Sampling) to discover which arm x optimizes the expected
 082 outcome $\mathbb{E}[Y; do(x)]$, which is an interventional (\mathcal{L}_2) strategy (Sutton & Barto, 1998; Latti-
 083 more & Szepesvári, 2020). It turns out there are provably superior strategies (w.r.t expected
 084 outcome) based on directly optimizing counterfactual (\mathcal{L}_3) objectives, as we demonstrate in
 085 Example 3. We prove optimality of our proposed strategy in a bandit setting with a generic
 086 causal template (Thm. F.2, Cor. F.3).

087 Finally, Sec. 5 discusses important themes, future directions, and the limitations of our work. Proofs
 088 and details of simulations are included in Appendices.
 089

090 **Preliminaries.** We denote variables by capital letters, X , and values by small letters, x . Bold
 091 letters, \mathbf{X} , are sets of variables and \mathbf{x} sets of values. $P(\mathbf{x})$ is shorthand for $P(\mathbf{X} = \mathbf{x})$. $\mathbb{1}[\cdot]$ is the
 092 indicator function. We use *Structural Causal Models* (SCM) to describe the generative process
 093 for a system of interest (Bareinboim et al., 2022, Def. 1)(Pearl, 2000). An SCM \mathcal{M} is a tuple
 094 $\langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$. \mathbf{V} is the set of observable variables. \mathbf{U} is the set of unobservable variables
 095 exogenous to the system, distributed according to $P^{\mathcal{M}}(\mathbf{U})$. $\mathcal{F} = \{f_V\}$ is a set of functions s.t.
 096 each f_V causally generates the value of $V \in \mathbf{V}$ as $V \leftarrow f_V(\mathbf{U}_V, \mathbf{Pa}_V)$, where $\mathbf{U}_V \subseteq \mathbf{U}$ and
 097 $\mathbf{Pa}_V \in \mathbf{V} \setminus V$. Each \mathcal{M} induces a *causal diagram* \mathcal{G} (Bareinboim et al., 2022, Def. 13), which
 098 is a graph containing a vertex for each $V \in \mathbf{V}$, a directed edge from each node in \mathbf{Pa}_V to V ,
 099 and a bidirected edge between V, V' if $\mathbf{U}_V, \mathbf{U}_{V'}$ are not independent. Given a graph \mathcal{G} , $\mathcal{G}_{\overline{\mathbf{X}}\mathbf{W}}$ is
 100 the result of removing edges coming into variables in \mathbf{X} , and edges coming out of \mathbf{W} . We use
 101 standard terminology like parents, descendants of a node (see App. A). Our treatment is limited to
 102 *recursive* SCMs, which implies acyclic diagrams, with finite discrete domains over \mathbf{V} . The *do*(\mathbf{x})
 103 operator indexes a sub-model $\mathcal{M}_{\mathbf{x}}$ where the functions generating variables \mathbf{X} are replaced with
 104 constant values \mathbf{x} . A variable $Y \notin \mathbf{X}$ evaluated in this regime is called a *potential response*, denoted
 105 $Y_{\mathbf{x}}$. $(\mathbf{W}_{\star} = \mathbf{w})$ denotes an arbitrary counterfactual event, e.g. $(Y_x = y \wedge Y_{x'} = y' \wedge X = x'')$.

106 ¹E.g., "The problem with counterfactuals like $[P(Y_x | x')]$ is [that] ... we simply cannot perform an
 107 experiment where the same person is both given and not given treatment." (Shpitser & Pearl, 2007) Also,
 "By definition, one can never observe [counterfactuals], nor assess empirically the validity of any modeling
 assumptions made about them..." (Dawid, 2000)

The probability of such an event is given by the \mathcal{L}_3 -valuation (Bareinboim et al., 2022, Def. 7):

$$P^{\mathcal{M}}(\mathbf{W}_* = \mathbf{w}) = \sum_{\mathbf{u}} \left(\prod_{W_t \in \mathbf{W}_*} \mathbb{1}[W_t(\mathbf{u}) = w] \right) P^{\mathcal{M}}(\mathbf{u}),$$
 with w taken from \mathbf{w} .

2 DATA-COLLECTION PROCEDURES

In this section, we define a procedure, *counterfactual randomization*, that extends the scope of traditional *Fisherian* experimentation (discussed below). Consider a system of interest modeled by unknown SCM \mathcal{M} . Interventions and counterfactual events are typically defined in terms of *symbolic* operations on \mathcal{M} . To conceptually separate this from the *physical* constraints experienced by an agent (natural or artificial), we define the following physical actions that an agent can perform in the system. These are simply the physical counterparts to symbolic procedures.

We call each discrete episode of the system’s behaviour a *unit*. Examples of units are patients in a clinical trial, neighbourhoods in a social science experiment, rounds played on a slot machine etc. We index units WLOG by $i = 1, 2, 3, \dots$, which constitute a target population in the system.

Definition 2.1 (Physical actions). (1) **SELECT**⁽ⁱ⁾: randomly choosing, without replacement, a unit i from the target population, to observe in the system; (2) **READ**(V)⁽ⁱ⁾: **measuring the realized feature $V^{(i)}$ of unit i , produced by a causal mechanism $f_V \in \mathcal{F}$ operating on i** ; (3) **RAND**(X)⁽ⁱ⁾: erasing and replacing i ’s natural mechanism f_X for a decision variable X with an enforced value drawn from a randomizing device having support over $\text{Domain}(X)$. ■

READ(V)⁽ⁱ⁾ = v and **RAND**(X)⁽ⁱ⁾ = x are also overloaded to refer to the values read and enforced, respectively. **RAND**(X)⁽ⁱ⁾ is the standard Fisherian randomization of a decision variable X , corresponding to the symbolic procedure of a *stochastic* intervention on X (Correa & Bareinboim, 2020).² As **RAND**(X)⁽ⁱ⁾ erases the unit i ’s natural decision, **READ**(X)⁽ⁱ⁾ will yield the value randomly assigned to unit i . The discovery of this procedure marked an important achievement in the history of science and experiment-design (Fisher, 1925; 1935). Since the use of a randomizing device eliminates by design any confounding between the assigned decision and the unit’s latent attributes $\mathbf{U}^{(i)}$, it allows researchers to estimate causal effects.

It is evident that the actions in Def. 2.1 are sufficient for an agent to physically draw samples from any \mathcal{L}_1 - or \mathcal{L}_2 -distribution, as discussed in App. C.3. Until recently, it was generally presumed these were the only physical actions possible on units in a system. However, we discuss some important extensions of experimental capabilities next.

Counterfactual data-collection procedures. In an early work from the causal reinforcement learning literature, Bareinboim, Forney & Pearl describe an experimental setting in which it is possible to both randomize a unit’s actual decision, and also record the natural decision the unit *would have normally* taken (Bareinboim et al., 2015; Forney et al., 2017). Subsequently, this procedure has been used to establish benchmarks in counterfactual decision making (Zhang & Bareinboim, 2022). These settings involve an agent introspecting to gauge their natural choice, or otherwise revealing their natural choice by some indication, e.g. physical gestures prior to decision-time. Importantly, this form of randomization does not erase the unit’s natural choice of decision variable X , as schematically illustrated in Fig. 1.

Building on the idea, we formalize this into a more general extension of the agent’s capabilities: the ability to intervene on a variable X ’s value *as perceived* by its causal children. To illustrate this, consider the \mathcal{L}_3 -quantity known as natural direct effect, or NDE, which is used in mediation analysis to measure the effect of X on Y via a ”direct” path, as opposed to an ”indirect” path via a mediator Z (Pearl, 2001) – highly relevant in several fields, as discussed in Sec. 1. The NDE is generally considered as identifiable from experimental data only under certain conditions (Pearl, 2005; Correa et al., 2021). The following example details an experiment design where it is possible to compute the NDE even when these identification conditions are not met, by randomizing the *perception* of X .

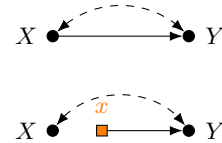


Figure 1: (Top) Causal diagram with decision variable X ; (Bottom) Procedure of randomizing the actual decision without erasing the unit’s natural decision.

²If the device used for enforcing the value of X is a constant function, this action simply becomes **WRITE**($X : x$)⁽ⁱ⁾, corresponding to the atomic intervention $do(x)$. See Preliminaries in Sec. 1.

Example 1 (Mediation analysis). A computer vision company’s tool is being evaluated for an automated speeding ticket system that uses footage from traffic cameras. But the government’s audit team has a concern: it is possible the model is trained on footage with a strong correlation between the color of the car and speeding (perhaps due to color preference of different socioeconomic neighbourhoods), and unfairly penalizes certain car colors.

This amounts to a hypothesis that X (car’s color) affects Y (AI decision to issue a ticket) via a direct path as opposed to the indirect path via Z (speeding). The indirect path describes the causal effects of, say, how pedestrians and other drivers react to a red car and affect its speeding. This hypothesis is true iff NDE is measured to be non-0, where NDE is defined as the following expression: $\text{NDE}_{x,x'}(y) = P(y_{x'Z_x}) - P(y_x)$ (Pearl, 2001). The second term, $P(y_x)$, can be estimated from a Fisherian randomization of X (say, an experiment recruiting drivers and assigning them random cars). Inconveniently, the first term, $P(y_{x'Z_x})$, is nonidentifiable for Fig. 2(a), even using RCT data. So it is unclear how to make progress with this hypothesis test.

However, the audit team recognizes there exists a special mediator, viz. the features W in the video which reveal the car’s color to the model (say, RGB values of pixels in the video frames). They use standard video-editing tools to randomly swap the color of the car in the footage. By randomly assigning a particular car $W \leftarrow \text{red}$, they are able to affect the mechanism f_Y ’s perception of X :

$$P(Y_{W=\text{red}} \mid X = \text{blue}) \quad \text{est. from } \mathcal{L}_2 \text{ data} \quad (1)$$

$$= P(Y_{W=\text{red}, Z} \mid X = \text{blue}) \quad Z : \text{natural value} \quad (2)$$

$$= P(Y_{W=\text{red}, Z_{X=\text{blue}}} \mid X = \text{blue}) \quad \text{consistency property} \quad (3)$$

$$= P(Y_{X=\text{red}, Z_{X=\text{blue}}} \mid X = \text{blue}) \quad \text{Def. 2.2, } X \equiv W \quad (4)$$

$$= P(Y_{X=\text{red}, Z_{X=\text{blue}}}) \quad \text{d-separation} \quad (5)$$

Eq. 4 is justified because W controls Y ’s perception of X given a fixed z (formalized in Lemma D.4). Thus, they are able to directly sample from the \mathcal{L}_3 -distribution $P(Y_{x'Z_x}, X)$ via a physical procedure, and use identification rules to obtain $P(y_{x'Z_x})$. Using the formula for NDE, they can evaluate whether a car’s color has a direct effect on the odds of getting a speeding ticket. ■

Here, one is able to randomize X as perceived by one of its children, by leveraging the variable W (RGB values) that fully encodes information about X (color) and mediates its effect on Y . We capture this intuition with the following (informal) definition.

Definition 2.2 (Counterfactual mediator (informal)). We call W a *counterfactual mediator* of X w.r.t $Y \in Ch(X)$ if the value of X can be retrieved from W by the mechanism generating Y . ■

Other examples of interventions on perceived attributes via counterfactual mediators include changing details on a job application (name, pronouns, keywords) to simulate a perceived alternate demographic identity (Bertrand & Mullainathan, 2003), or editing specific portions of text input to a language model (Feder et al., 2022). Randomizing perception has been discussed in Pearl et al. (2016, §4.4.4). For a detailed discussion of the causal semantics of intervening on perceptions, and the related literature, see (Plecko & Bareinboim, 2024, App. D.1). For the interested reader, we provide a rigorous treatment in App. D, including a formal Def. D.2 of a counterfactual mediator.

This important extension to experimental capabilities is captured in the following definition of a new physical action that an agent might be able to perform in an environment.

Definition 2.3 (Counterfactual (ctf-) randomization). $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$: fixing the value of X as an input to the mechanisms generating $\mathbf{C} \subseteq Ch(X)_{\mathcal{G}}$ using a randomizing device having support over $\text{Domain}(X)$, for unit i , given causal diagram \mathcal{G} . ■

The key differences between the Fisherian $\text{RAND}(X)^{(i)}$ and $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$ are (1) CTF-RAND does not erase the unit i ’s natural decision $X^{(i)}$; and (2) while RAND affects all children of X , CTF-RAND does not affect $Ch(X) \setminus \mathbf{C}$. CTF-RAND can only be enacted under certain structural conditions, viz., either in environments which permit the measurement of a unit’s natural decision while simultaneously randomizing the actual decision (Bareinboim et al., 2015), or where counterfactual mediators can be used to alter X as perceived by a subset of children. Whether the agent is indeed able to perform this action thus depends on the specific experimental setting.

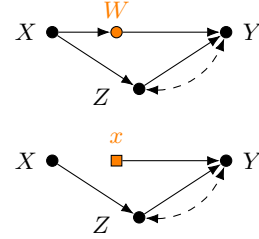


Figure 2: (a) ”Expanded” diagram for Example 1, where W is *counterfactual mediator* for X ; (b) Randomizing the value of X as perceived by Y .

Note: Def. 2.3 implies that it is possible to physically perform multiple randomizations involving the same variable X on a single unit i , with each intervention affecting a different subset of children. Further, CTF-RAND may only be performed w.r.t a graphical child variable; it is not possible to bypass a child and directly affect a descendant’s perception of X .

3 COUNTERFACTUAL REALIZABILITY

Given the possibility of performing ctf-randomization (Def. 2.3), we are interested in knowing which \mathcal{L}_3 -distributions can be accessed directly by experimentation. In this section, we discuss the constraints imposed by nature on an agent. We then formally define *realizability* and develop a complete algorithm to determine whether an \mathcal{L}_3 -distribution is realizable.

The most basic constraints experienced by the agent (natural or artificial) are physical. Each mechanism $f_V \in \mathcal{F}$ represents some physical process that transforms a unit i according to the laws of nature. For instance, taking a drug, X , produces a side effect in the patient, Y , by a biochemical reaction $f_Y(X, U_Y)$, which depends on the drug and the patient’s latent health condition, U_Y . Once patient i has been subjected to mechanism f_Y under $X = x$, there appears to be no way to go back in time and subject the same patient to mechanism f_Y under $X = x'$. Even if technologically feasible to reverse the process (e.g., by taking an antidote to the drug), the latent factors $\mathbf{U} = \mathbf{u}$ might have changed after the experiment (e.g., the patient could have developed tolerance to the drug). Repeating the experiment on this patient is tantamount to testing a *new* unit with unknown latent features $\mathbf{U} = \mathbf{u}'$.³ This observation is made more formal through the following assumption.

Assumption 3.1 (Fundamental constraint of experimentation (FCE)). A unit i in the target population can physically undergo a causal mechanism $f_V \in \mathcal{F}$ at most once. ■

Remark 3.2. The FCE assumption entails that a unit i can only be submitted to a particular mechanism $f_V(\mathbf{Pa}_V, \mathbf{U}_V)$ under a single set of experimental conditions, received as input to f_V . By implication, the physical actions in Defs. 2.1, 2.3 can only be performed at most once per unit i . ■

Once unit i has been subjected to f_V , it is not possible to re-run f_V with differently fixed inputs. $\text{READ}(V)^{(i)}$ thus only yields one value for i . Although ctf-randomization permits multiple interventions involving the same variable X , each such intervention can only be performed once, since it impacts different child mechanisms that can each only occur once for unit i . We also assume that the agent can only perform the physical actions in Defs. 2.1, 2.3, up to isomorphism.

Definition 3.3 (I.i.d sample). Given an \mathcal{L}_3 -distribution $Q = P(\mathbf{W}_*)$ and a sequence of physical actions $\mathcal{A}^{(i)}$ performed on unit i in an environment modeled by SCM \mathcal{M} , producing a vector of realized values $\mathbf{W}_*^{(i)} = \mathbf{w}$ for the variables in \mathbf{W}_* , the vector is said to be an *i.i.d sample* from Q if $P^{\mathbb{C}}(\mathbf{W}_*^{(i)} = \mathbf{w} \mid \mathcal{A}^{(i)}) = P^{\mathcal{M}}(\mathbf{W}_* = \mathbf{w}), \forall \mathbf{w}$, where $P^{\mathbb{C}}$ is the probability measure over the beliefs of the acting agent \mathbb{C} , and the l.h.s is the probability of physical actions $\mathcal{A}^{(i)}$ producing the vector \mathbf{w} when performed on some unit i . ■

Definition 3.4 (Realizability). Given a causal diagram \mathcal{G} and the set of physical actions \mathbb{A} , an \mathcal{L}_3 -distribution $P(\mathbf{W}_*)$ is *realizable given \mathbb{A} and \mathcal{G}* iff there exists a sequence of actions \mathcal{A} from \mathbb{A} by which an agent can draw an i.i.d sample (Def. 3.3) from $P^{\mathcal{M}}(\mathbf{W}_*)$, for any $\mathcal{M} \in M(\mathcal{G})$, the class of SCMs compatible with \mathcal{G} . ■

We emphasize the *distinction between realizability and identifiability*. Identifiability (Pearl, 2000, Def. 3.2.3) from \mathcal{G} states that a distribution (say, $P(\mathbf{v}; \text{do}(x))$) can be uniquely computed from the available data (say, $P(\mathbf{v})$) for any SCM compatible with the assumptions in \mathcal{G} . Realizability of a distribution states that it is physically possible for an agent to actually gather data samples according to this distribution.

We next develop an algorithm to decide whether a distribution is realizable. As an intuition pump, suppose that an agent is able to perform $\text{CTF-RAND}(V \rightarrow C), \forall V, C \in Ch(V)$, w.r.t an input causal diagram, and wants to obtain samples from $P(Z_x, W_t)$. Consider the diagram \mathcal{G}_2 in Fig. 3. By performing $\text{CTF-RAND}(T \rightarrow W)$ and $\text{CTF-RAND}(X \rightarrow Z)$, the distribution is realizable. However, suppose the input diagram is \mathcal{G}_1 . A necessary condition to measure Z_x for a unit is for mechanism f_A

³In the philosophy of science literature, similar ideas have been discussed under the topic of the temporal asymmetry of causation (Reichenbach, 1956, §III-IV).

Algorithm 1 CTF-REALIZE

```

270 1: Input:  $\mathcal{L}_3$ -distribution  $Q = P(\mathbf{W}_*)$ ; causal
271    diagram  $\mathcal{G}$ ; action set  $\mathbb{A}$ 
272
273 2: Output: I.i.d sample  $\mathbf{W}_*^{(i)}$  from  $Q$ ; FAIL if
274     $Q$  is not realizable given  $\mathcal{G}, \mathbb{A}$ 
275
276 3: Fix a topological ordering  $\text{Top}(\mathcal{G})$ 
277
278 4: SELECT(i) for a new unit  $i$ 
279
280 5: for  $V$  in order  $\text{Top}(\mathcal{G})$  do
281   6:  $\text{INT}_V \leftarrow \emptyset$  {Interventions for  $V$ }
282   7:  $\text{OUTPUT}_V \leftarrow \emptyset$  {Index in output vector}
283   8: for each term  $W_t$  in expression  $\mathbf{W}_*$  do
284     9: if  $V \in \text{An}(W)_{\mathcal{G}_{\mathbb{T}}}$  and  $V \neq W$  then
285       10: Call COMPATIBLE( $V, W_t$ ) 2
286     11: end if
287     12: if  $V = W$  then
288       13: Add  $\{W_t\}$  to  $\text{OUTPUT}_V$ 
289     14: end if
290   15: end for
291
292   16: for each  $\{\text{action} : \text{tag}\} \in \text{INT}_V$  do
293     17: Perform the randomization on unit  $i$ 
294     18: If the random-generated value  $\neq$  tag,
295         discard the unit and return to Line 4
296   19: end for
297   20: for each  $W_t \in \text{OUTPUT}_V$  do
298     21: if  $\{\text{RAND}(V) : \cdot\} \in \text{INT}_V$  then
299       22: Return FAIL
300     23: else
301       24: Perform  $\text{READ}(V)^{(i)} = v'$ 
302       25: Assign  $v'$  to each index  $W_t^{(i)}$  in out-
303           put vector  $\mathbf{W}_*^{(i)} = \mathbf{w}$ 
304     26: end if
305   27: end for
306
307 28: end for
308
309 29: Return i.i.d sample  $\mathbf{W}_*^{(i)} = \mathbf{w}$ 

```

to receive the natural value of T , illustrated in green. While a necessary condition to simultaneously measure W_t is for f_W to receive A_t , which in turn requires f_A to receive a fixed t , shown in red. This conflict in necessary conditions renders the query non-realizable.⁴

This "edge-coloring" intuition is formalized in Algo. 1. The algorithm **CTF-REALIZE** takes as input an \mathcal{L}_3 -distribution $P(\mathbf{W}_*)$, a graph \mathcal{G} , and a set of physical actions \mathbb{A} the agent is able to perform in the environment (viz., the **RAND** and **CTF-RAND** actions which are possible in the environment). It returns an i.i.d sample if the distribution is realizable, and **FAIL** otherwise.

The algorithm works as follows (a more detailed walk-through is presented in App. B.2): going over each node V in topological order, the inner loops gather the necessary and sufficient conditions needed w.r.t V for realizing each W_t in the input query \mathbf{W}_* . If there is a conflict in the necessary conditions for evaluating two terms (as we saw for $P(Z_x, W_t)$ in Fig. 3, \mathcal{G}_1), the query is non-realizable. The algorithm is fully general and does not make assumptions about the ability to perform any particular interventions. If the agent cannot perform any counterfactual randomization, the algorithm returns **FAIL** for non- \mathcal{L}_2 queries. If the agent cannot perform any interventions at all, the algorithm returns **FAIL** for non- \mathcal{L}_1 queries (we assume the ability to **READ** all variables). Details about the time and space complexity of Algo. 1 are provided in App. B.3, for the interested reader.

Theorem 3.5 (Correctness and Completeness). *An \mathcal{L}_3 -distribution $Q = P(\mathbf{W}_*)$ is realizable given action set \mathbb{A} and causal diagram \mathcal{G} iff the algorithm **CTF-REALIZE**($Q, \mathcal{G}, \mathbb{A}$) returns a sample. ■*

A further question we may ask is which \mathcal{L}_3 -distributions are realizable if we assume maximum experimental capabilities, notably, the ability to perform separate ctf-randomization for each child of each variable. Given a causal diagram \mathcal{G} , we define the *maximal feasible action set* $\mathbb{A}^\dagger(\mathcal{G})$ as the set containing all of the following actions: $\text{SELECT}^{(i)}$, $\text{READ}(V)^{(i)}$, $\forall V$, and $\text{CTF-RAND}(X \rightarrow C)^{(i)}$, $\forall X$ and $C \in \text{Ch}(X)$. $\mathbb{A}^\dagger(\mathcal{G})$ thus gives the agent the most granular interventional capabilities.

Definition 3.6 (Ancestors of a counterfactual (Correa et al., 2021)). Given a causal diagram \mathcal{G} and a potential response Y_x , the set of (counterfactual) ancestors of Y_x , denoted $\text{An}(Y_x)$, consists of each

⁴To be clear, the input to the algorithm is a graph and an accurate set of actions the agent can perform in the environment. If the graph is per \mathcal{G}_1 in Fig. 3, then $\text{CTF-RAND}(T \rightarrow Z)$ is not possible in this environment. Marginalizing out A and providing graph \mathcal{G}_2 as input does not help.

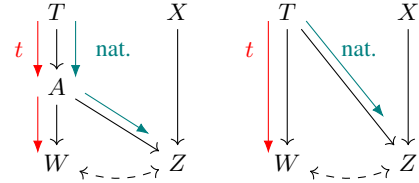


Figure 3: Testing realizability of $P(Z_x, W_t)$ for \mathcal{G}_1 (left) and \mathcal{G}_2 (right). \mathcal{G}_1 yields conflicting requirements.

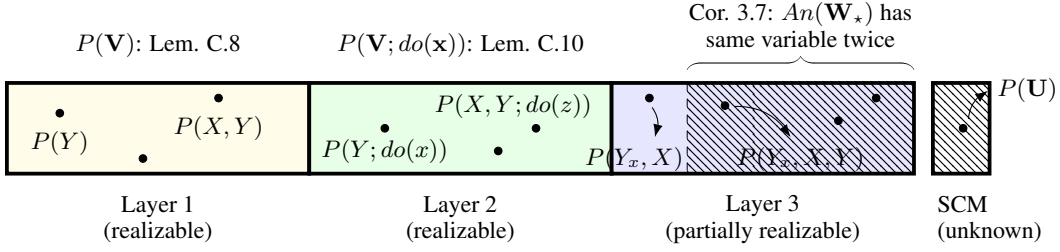


Figure 4: Pearl Causal Hierarchy (PCH) induced by an unknown SCM \mathcal{M} . An \mathcal{L}_3 -distribution is realizable given a graph \mathcal{G} and the maximal feasible action set $\mathbb{A}^\dagger(\mathcal{G})$ iff the ancestor set $An(\mathbf{W}_*)$ does not contain the same variable under different regimes.

$W_{\mathbf{z}}$ s.t. $W \in An(Y)_{\mathcal{G}_{\mathbf{x}}}$, and $\mathbf{z} = \mathbf{x} \cap An(W)_{\mathcal{G}_{\bar{\mathbf{x}}}}$. For a set \mathbf{W}_* , $An(\mathbf{W}_*)$ is defined to be the union of the ancestors of each potential response in the set. ■

Corollary 3.7. *An \mathcal{L}_3 -distribution $Q = P(\mathbf{W}_*)$ is realizable given causal diagram \mathcal{G} and action set $\mathbb{A}^\dagger(\mathcal{G})$ iff the ancestor set $An(\mathbf{W}_*)$ does not contain a pair of potential responses W_t, W_s of the same variable W under different regimes.* ■

For instance, if $\mathbf{W}_* = \{Z_x, W_t\}$ w.r.t graph \mathcal{G}_1 in Fig. 3, then $An(\mathbf{W}_*) = \{Z_x, A, T, W_t, A_t\}$, which contains both A, A_t . Thus, $P(\mathbf{W}_*)$ is not realizable even with maximal experimentation capabilities. In App. B.4, we provide further examples of using the **CTF-REALIZE** algorithm, and the graphical criterion, to demonstrate the realizability of the ETT distribution $P(Y_x, X)$, the non-realizability of the probability of sufficiency distribution $P(Y_x, X, Y)$.

We believe this is an important contribution to causal inference. Cor. 3.7 provides a graphical criterion to delineate how far up the PCH an agent can go via experimental methods, in principle. Often, counterfactuals have been criticized as being hypothetical, untestable, or unscientific assumptions. Our analysis counters this claim, as summarized in Fig. 4.

Corollary 3.8 (Fundamental problem of causal inference (FPCI) (Holland, 1986)). *The distribution $Q = P(Y_x, Y_{x'})$ is not realizable given maximal feasible action set $\mathbb{A}^\dagger(\mathcal{G})$, for any causal diagram \mathcal{G} , and any variables $X, Y \in Desc(X)$.* ■

The FPCI is an influential notion in the literature, and is often taken as a primitive, or in an axiomatic fashion. We show that it is rather a specific consequence of the more general FCE assumption 3.1, and follows from Thm. 3.5 and Cor. 3.7. By itself, the FPCI does not translate to an operational criterion for determining which \mathcal{L}_3 -distributions are realizable (Def. 3.4). For instance, it does not clarify that a distribution with potential responses under conflicting regimes like $P(Y_x, Z_{x'})$ may indeed be realizable via counterfactual randomization, as we show in Example 2. It also does not tell us that $P(Z_x, W_t)$ may be realizable given causal diagram \mathcal{G}_2 in Fig. 3, but not realizable given \mathcal{G}_1 .

4 APPLICATIONS: COUNTERFACTUAL DECISION-MAKING AND FAIRNESS

Next, we highlight the practical relevance of our results with some concrete use-cases. We already discussed in Example 1 how realizability can be used to design experiments for performing **mediation analysis** of direct/indirect effects, an important task in several fields. We now discuss applications in **causal fairness analysis** and **causal reinforcement learning (RL)**. Our goal is to underscore that the standard/baseline approaches in these areas, even among approaches that incorporate counterfactual reasoning, typically use observational (\mathcal{L}_1) or interventional (\mathcal{L}_2) data only, whereas a counterfactual (\mathcal{L}_3) data-collection approach can lead to demonstrably better results. Due to space constraints, we include in App. E the full specification of SCMs used and algorithms implemented.

4.1 CAUSAL FAIRNESS - USING COUNTERFACTUAL DATA FOR FAIRER DECISIONS

Causal fairness analysis is a burgeoning field and a full survey is beyond the scope of this paper (see, e.g., Plecko & Bareinboim (2024) for a review of related works). We limit our discussion to an example where counterfactual realizability is directly relevant.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

A common concern is that models trained to make automated decisions often reveal problematic biases (Angwin et al., 2016; Kodiyam, 2019, e.g.). The causal approach to address this is typically to constrain a classifier to obey some causally-sensitive *fairness measure*, μ (Plecko & Bareinboim, 2024, Def. 3.3). Some measures in the literature involve \mathcal{L}_3 -quantities, and thus face the familiar issue of nonidentifiability (Kusner et al., 2017; Imai & Jiang, 2023). Other approaches acknowledge this limitation and try to construct interventional fairness measures that solely use \mathcal{L}_2 -quantities (Salimi et al., 2019). We present next an example where relying only on \mathcal{L}_2 -data can misleadingly approve a classifier as fair, but where a realizable \mathcal{L}_3 fairness measure actually ensures fairness. This scenario is inspired by a classic experiment in labor economics (Bertrand & Mullainathan, 2003).

Example 2 (Causal fairness). A college is developing an automated system to screen candidates in the first round of college applications, receiving as input a standardized CV per candidate. The system contains two models: model 1 outputs Y and model 2 outputs Z , which are binary decisions of whether the applicant cleared the first review stage for admission and for financial scholarship, respectively. The two models are respectively trained using data from previous years where an admissions team and a separate scholarship team reviewed applications manually. The college wants to ensure fairness w.r.t X , a candidate’s race (a binary variable, for simplicity). In particular, they want to ensure equitable financial access to education for all qualified candidates: a candidate of race $X = 1$ who cleared the admissions screening ($Y = 1$) but was rejected for financial aid ($Z = 0$) should still receive $Z = 0$ had they been of race $X = 0$. The causal diagram is in Fig. 5(a), where the models’ decisions Y, Z might reflect the unconscious race bias of the two committees in previous years (including possibly shared biases, represented by the latent confounder).

The \mathcal{L}_3 fairness measure they ought to minimize is thus

$$\mu_{ctf} = |P(Y_{x_1} = 1, Z_{x_1} = 0) - P(Y_{x_1} = 1, Z_{x_0} = 0)| \quad (6)$$

But the second term $P(y_x, z'_{x'})$ is nonidentifiable from the causal diagram in 5(a). So the college instead uses the following \mathcal{L}_2 measures, as an approximation for the fairness condition:

$$\mu_{int1} = |P(Y = 1; do(x_1)).P(Z = 0; do(x_1)) - P(Y = 1; do(x_1)).P(Z = 0; do(x_0))| \quad (7)$$

$$\mu_{int2} = |P(Y = 1, Z = 0; do(x_1)) - P(Y = 1, Z = 0; do(x_0))| \quad (8)$$

They train the models, adding $\mu_{int1} + \mu_{int2}$ as a penalty in the objective. μ_{int1}, μ_{int2} are estimated using a holdout set of fake CVs, with the intervention $do(x)$ being enacted by randomly choosing an applicant name from an equivalence class which stereotypically indicates one unique race group $X = x$, e.g. names like Lakisha and Jamal for Blacks, or last names like Nguyen or Xi for Asians (cf. Bertrand & Mullainathan (2003)). Since the holdout set’s CV body is independent of X , any effect of X on Y and Z is solely via the perception of race from the candidate name. We show in 5(c) simulations of such an optimization. In blue is the distribution of the true score μ_{ctf} , when the models are trained using μ_{int1}, μ_{int2} . Out of 1000 simulations of classifiers, we see $\mu_{ctf} > 5\%$ for nearly half the \mathcal{L}_2 simulations, indicating statistically significant discrimination roughly 50% the time.

However, the distribution $P(Y_x, Z_{x'})$ is indeed realizable (Def. 3.4) via the interventions $CTF-RAND(X \rightarrow Y), CTF-RAND(X \rightarrow Z)$. The data science team notices that they can *separately and simultaneously* randomize the candidate name as an input to the respective models, and enact these interventions, as shown in 5(b). Thus, they are able to directly use the counterfactual measure μ_{ctf} as a fairness constraint in training. Results from 1000 simulations show that the classifiers trained directly using μ_{ctf} (shown in orange) nearly always meet the fairness requirement.

Details of the implementation are in App. E.2. Note: as in the original experiment, this example requires the structural assumption of race being revealed at the screening stage only by candidate name, which may be more defensible in highly standardized and controlled application processes. ■

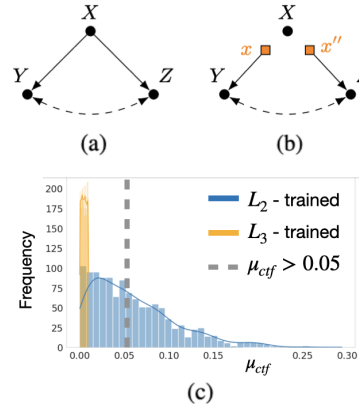


Figure 5: (a) Causal diagram for Example 2; (b) $P(Y_x, Z_{x'})$ is realizable using the interventions $CTF-RAND(X \rightarrow Y)$ and $CTF-RAND(X \rightarrow Z)$; (c) Histogram of 1000 classifiers trained on \mathcal{L}_2 (blue) and \mathcal{L}_3 (orange) fairness measures. \mathcal{L}_2 classifiers show statistically significant discrimination ($\mu_{ctf} > 0.05$).

4.2 CAUSAL RL - COUNTERFACTUAL POLICIES FOR OPTIMAL DECISION-MAKING

Consider a multi-arm bandit problem in which X represents the choice of bandit arm and Y the outcome. The default online learning approach is for the agent to adopt an algorithm like EXP3, UCB or Thompson Sampling to converge to some arm $x^* := \arg \max_x \mathbb{E}[Y; do(x)]$ (Lattimore & Szepesvári, 2020; Sutton & Barto, 1998). Even in methods that explicitly incorporate causal knowledge, the common approach is to use a combination of offline (\mathcal{L}_1) and online (\mathcal{L}_2) data to converge more efficiently to the \mathcal{L}_2 optimization target $\arg \max_x \mathbb{E}[Y; do(x)]$ (Zhang & Bareinboim, 2017, e.g.). It was already shown in (Bareinboim et al., 2015; Forney et al., 2017) that it is possible to perform better by deploying a counterfactual strategy based on sampling each unit’s natural choice $X = x'$ and randomizing actual choice in the *same* round, thus seeking to converge to $\arg \max_x \mathbb{E}[Y_x | x'], \forall x'$, as we discussed in Sec. 2. We call this the ETT baseline strategy, as it relies on drawing samples from the \mathcal{L}_3 ETT distribution, $P(Y_x, X)$, mentioned in Sec. 1.

We improve on this baseline by showing how an agent can leverage the realizability (Def. 3.4) of more nuanced counterfactuals like $P(Y_x, X, D_{x''})$ to construct superior counterfactual strategies. The following scenario involves an agent faced with adversarial latent confounding.

Example 3 (Counterfactual bandit policies). Consider a user of a social media platform which uses surveillance and predictions to increase user engagement through addictive notifications and recommendations (Zuboff, 2018). The user chooses every evening whether to use the platform via desktop ($X = 0$) or mobile ($X = 1$). Y is a binary indicator of whether she stays within her self-determined social media usage limit per day. She also notices that she receives ads when she logs in each evening as D (0: streaming service, 1: food delivery ads). The usage type X affects D, Y , as shown in Fig. 6(a).

On average, the user experiences $\mathbb{E}[Y] = 0.65$ from the observational (\mathcal{L}_1) policy of following her natural inclination each day. She suspects that the company could be tracking and exploiting her latent preferences, so she decides to randomize her daily choice and pick the best “arm”. Sure enough, this naive \mathcal{L}_2 strategy breaks the adversarial confounding, and incurs a better avg. performance of $\mathbb{E}[Y; do(x)] = 0.7, \forall x$. She then decides to test the ETT-based strategy (\mathcal{L}_3) described earlier, by recording what she naturally feels like doing each day ($X = x'$), and subsequently randomizing her actual choice on the same day to optimize $\mathbb{E}[Y_x | x']$, getting an avg. performance of 0.75. However, at this point, she notices that she can do even better. The \mathcal{L}_3 -distribution $P(Y_x, X, D_{x''})$ is realizable (Def. 3.4), since she can perform *another* counterfactual randomization, by sampling her natural choice ($X = x'$), randomly logging in to just see what ads she gets ($D_{x''} = d$), and again randomizing how she actually uses the platform that day to get Y_x . This strategy seeks an optimal $x^* = \arg \max_x \mathbb{E}[Y_x | x', d_{x''}]$, which performs best as shown in Table 1. Details of the SCM, latent confounders, and the optimal \mathcal{L}_3 -strategy are in App. E.3.

Simulations in the online setting corroborate this finding. Fig. 6(c,d) shows the cumulative regret (CR) and optimal arm probability (OAP) over 2000 iterations averaged over 200 epochs (CI=95%). We adapt Thompson Sampling to implement the strategies in Table 1. Details of implementation are in App. E.3.1. The optimal \mathcal{L}_3 strategy (purple) performs best, improving on the performance of the baseline ETT-based strategy (red). Naive randomizations, the standard \mathcal{L}_2 bandit strategy, are shown in yellow and green. All other algorithms fail to improve in OAP after 2000 iterations. ■

Table 1: Performance of different strategies in Example 3.

Strategy	$\mathbb{E}[Y]$
Behavioral policy (\mathcal{L}_1)	0.65
Naive randomization (\mathcal{L}_2)	0.7
ETT baseline strategy (\mathcal{L}_3)	0.75
Optimal \mathcal{L}_3 strategy (this work)	0.80

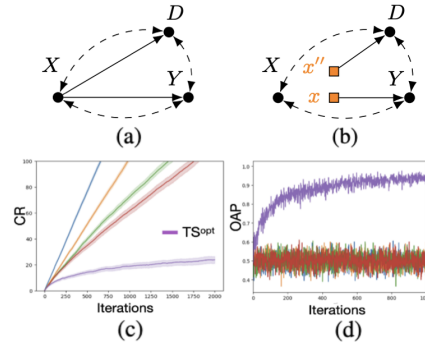


Figure 6: (a) Causal diagram for Example 3; (b) $P(Y_x, X, D_{x''})$ is realizable using the interventions $CTF\text{-}RAND(X \rightarrow Y)$ and $CTF\text{-}RAND(X \rightarrow D)$; (c) Cumulative Regret (CR) for \mathcal{L}_1 strategy (blue) and Thompson Sampling algorithms implementing naive \mathcal{L}_2 (yellow, green), ETT baseline (red), and optimal \mathcal{L}_3 strategy (purple); (d) Optimal Arm Probability (OAP) for all algorithms.

We make two remarks. First, the optimal counterfactual strategy is not simply a contextual Thompson Sampling, where $X, D_{x''}$ are used "merely" as extra context variables per round; indeed, treating this merely as a contextual bandit problem is one of the naive \mathcal{L}_2 -strategies that we test (green plot in Fig. 6(c-d)), which ignores the counterfactual relationship between these variables and incurs dramatically higher regret, as we discuss in App. E.3.1.

Second, an interesting follow-up is whether we can guarantee that our strategy based on maximizing $\mathbb{E}[Y_x | x', d_{x''}]$ is optimal in this problem. Perhaps there are more refined \mathcal{L}_3 -distributions like $P(Y_x, X, D_x, D_{x''})$ etc. that could yield better algorithms? It turns that this is indeed optimal, since most other \mathcal{L}_3 -distributions are not realizable (Def. 3.4). As a bonus, we prove this claim for all bandit problems that fit a specific causal template in App. F. Thereby, we avoid having to conduct an intractable search over the space of all possible \mathcal{L}_3 -strategies, trying to assess their realizability.

5 DISCUSSION

Finally, we discuss some important implications, future directions, and limitations of our work.

Identification and bounding. Much work has been done in the area of \mathcal{L}_3 identification and estimation (Shpitser, 2008; Correa et al., 2021; Geneletti & Dawid, 2011). A natural extension to our work is to investigate the relationship between realizability and identification: which *additional* \mathcal{L}_3 -quantities now become identifiable if the environment permits even some counterfactual randomization? This warrants an update to existing identification algorithms to allow (some) \mathcal{L}_3 -data as input. [Another fascinating research question involves "partial identification", where an input query is tightly bounded within a range that can be computed from available data \(Zhang et al., 2022\): how would the new \$\mathcal{L}_3\$ -data further tighten the bounds for nonidentifiable \$\mathcal{L}_3\$ -quantities?](#)

Experiment design. One of the goals of this paper is to instigate new experiment design ideas that leverage ctf-randomization (Def. 2.3) and go beyond the standard RCT methodology, as in Examples 1-2. For instance, the increasingly automated HR pipeline in companies suggests opportunities for targeted interventions to randomize demographic details in virtual interviews, in standardized aptitude tests, or in performance-evaluation systems for remote workers, to track fairness metrics.

Causal reinforcement learning (CRL). While counterfactual strategies have been studied in CRL, the literature currently focuses on ETT-related strategies based on optimizing $\mathbb{E}[Y_x | x']$ (Bareinboim et al., 2015; Forney et al., 2017; Zhang & Bareinboim, 2022)(Richardson & Robins, 2013, §5.1). We presented an important extension by formalizing ctf-randomization (Def. 2.3) via counterfactual mediators (Def. 2.2), subsuming the previous approach. An ETT-based approach only allows one randomization of a variable X , affecting all downstream mechanisms. Our approach recognizes the possibility of isolating specific causal pathways and randomizing X multiple times per unit, demonstrably surpassing the ETT baseline in Example 3. We proved in App. F an optimality guarantee for our proposed strategy in bandit problems following a causal template. Generalizing this to sequential decision-making settings with arbitrary graphs is an important, non-trivial extension.

Limitations. The first obvious limitation of our framework is that it requires causal knowledge in the form of a graph (or equivalent). This is a standard assumption, needed to make progress in several areas of causal machine learning. Subsequent work could accommodate partial knowledge or model misspecification. Second, it may not always be feasible to perform counterfactual randomization (Def. 2.3) in a given setting. This is why Algo. 1 and Thm. 3.5 are general and do not assume this capability a priori. But where it is possible, even in principle, our work pinpoints opportunities for novel experiment design, as discussed above.

6 CONCLUSION

In this paper, we tackle the open question of which counterfactual distributions are directly accessible by experimental methods - what we define as the *realizability* of a distribution. Countering prevalent belief, we provide a complete algorithm and a graphical criterion for when a counterfactual can indeed be physically sampled from (Fig. 4). We demonstrate the practical relevance of this new framework with examples from causal fairness and causal RL, highlighting that ignoring this possibility could lead to poor outcomes. We believe that switching from an interventional to a counterfactual mindset could help researchers spot opportunities for *counterfactual randomization* that permit exciting new types of experiments, and improved, more personalized decisions.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

REFERENCES

- Tara V. Anand, Adele H. Ribeiro, Jin Tian, and Elias Bareinboim. Causal effect identification in cluster dags. In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'23/IAAI'23/EAAI'23*. AAAI Press, 2023. ISBN 978-1-57735-880-0. doi: 10.1609/aaai.v37i10.26435. URL <https://doi.org/10.1609/aaai.v37i10.26435>.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine Bias: There's software used across the country to predict future criminals. And it's biased against blacks. 2016. URL <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- Elias Bareinboim, Andrew Forney, and Judea Pearl. Bandits with unobserved confounders: A causal approach. In *Advances in Neural Information Processing Systems*, pp. 1342–1350, 2015.
- Elias Bareinboim, Juan D. Correa, Duligur Ibeling, and Thomas Icard. On Pearl's Hierarchy and the foundations of causal inference. In *Probabilistic and Causal Inference: The Works of Judea Pearl*, pp. 507–556. Association for Computing Machinery, New York, NY, USA, 1st edition, 2022.
- Marianne Bertrand and Sendhil Mullainathan. Are Emily and Greg more employable than Lakisha and Jamal? A field experiment on labor market discrimination. Working Paper 9873, National Bureau of Economic Research, July 2003. URL <http://www.nber.org/papers/w9873>.
- J. Correa and E. Bareinboim. A calculus for stochastic interventions: Causal effect identification and surrogate experiments. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, New York, NY, 2020. AAAI Press.
- Juan Correa, Sanghack Lee, and Elias Bareinboim. Nested counterfactual identification from arbitrary surrogate experiments. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 6856–6867. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/36bedb6eb7152f39b16328448942822b-Paper.pdf.
- A Philip Dawid. Causal Inference Without Counterfactuals (with Comments and Rejoinder). *Journal of the American Statistical Association*, 95(450):407–448, 2000.
- Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E. Roberts, Brandon M. Stewart, Victor Veitch, and Diyi Yang. Causal inference in natural language processing: Estimation, prediction, interpretation and beyond. *Transactions of the Association for Computational Linguistics*, 10:1138–1158, 2022. doi: 10.1162/tacl.a.00511. URL <https://aclanthology.org/2022.tacl-1.66>.
- Ronald A. Fisher. *Statistical Methods for Research Workers*. Oliver and Boyd, Edinburgh, 1925.
- Ronald A. Fisher. *The Design of Experiments*. Oliver and Boyd, Edinburgh, 1935.
- Andrew Forney, Judea Pearl, and Elias Bareinboim. Counterfactual Data-Fusion for Online Reinforcement Learners. In *Proceedings of the 34th International Conference on Machine Learning*, 2017. ISBN 9781510855144. doi: <http://dx.doi.org/10.1037/a0022750>.
- Sara Geneletti and A. Philip Dawid. 72834 defining and identifying the effect of treatment on the treated. In *Causality in the Sciences*. Oxford University Press, 03 2011. ISBN 9780199574131. doi: 10.1093/acprof:oso/9780199574131.003.0034. URL <https://doi.org/10.1093/acprof:oso/9780199574131.003.0034>.
- James J. Heckman and Richard Robb Jr. Alternative Methods for Evaluating the Impact of Interventions. In J J Heckman and B Singer (eds.), *Longitudinal Analysis of Labor Market Data*. Cambridge University Press, New York, NY, 1985.

-
- 594 James J. Heckman and Richard Robb Jr. Alternative Methods for Solving the Problem of Selection
595 Bias in Evaluating the Impact of Treatments on Outcomes. In H. Wainer (ed.), *Drawing Inference*
596 *From Self Selected Samples*, pp. 63–107. Springer-Verlag, New York, NY, 1986.
- 597 P W Holland. Statistics and Causal Inference. *Journal of the American Statistical Association*, 81
598 (396):945–960, 12 1986.
- 600 Duligur Ibeling and Thomas Icard. Probabilistic reasoning across the causal hierarchy. In *Proceedings*
601 *of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 10170–10177, 2020.
- 602 Kosuke Imai and Zhichao Jiang. Principal Fairness for Human and Algorithmic Decision-Making.
603 *Statistical Science*, 38(2):317 – 328, 2023. doi: 10.1214/22-STS872. URL [https://doi.](https://doi.org/10.1214/22-STS872)
604 [org/10.1214/22-STS872](https://doi.org/10.1214/22-STS872).
- 606 Akhil Alfons Kodiyan. An overview of ethical issues in using ai systems in hiring with a case study
607 of amazon’s ai based hiring tool. *ResearchGate preprint*, pp. 1–19, 11 2019.
- 608 Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In
609 I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett
610 (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2017/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf)
611 [2017/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf).
- 612 Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- 613
614
615 Ang Li and Judea Pearl. Unit selection based on counterfactual logic. In *Proceedings of the*
616 *Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 1793–1799.
617 International Joint Conferences on Artificial Intelligence Organization, 7 2019. doi: 10.24963/
618 [ijcai.2019/248](https://doi.org/10.24963/ijcai.2019/248). URL <https://doi.org/10.24963/ijcai.2019/248>.
- 619
620 Scott Mueller and Judea Pearl. Personalized decision making – a conceptual introduction. *Journal*
621 *of Causal Inference*, 11(1):20220050, 2023. doi: doi:10.1515/jci-2022-0050. URL [https:](https://doi.org/10.1515/jci-2022-0050)
622 [//doi.org/10.1515/jci-2022-0050](https://doi.org/10.1515/jci-2022-0050).
- 623 Scott Mueller and Judea Pearl. Perspective on ‘harm’ in personalized medicine – an alternative
624 perspective. forthcoming.
- 625
626 Judea Pearl. Probabilities of causation: Three counterfactual interpretations and their identification.
627 *Synthese*, 121(1–2):93–149, 11 1999.
- 628 Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York,
629 NY, USA, 2nd edition, 2000. ISBN 978-0-521-89560-6.
- 630
631 Judea Pearl. Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty*
632 *in Artificial Intelligence*, pp. 411–420. Morgan Kaufmann, San Francisco, CA, 2001.
- 633
634 Judea Pearl. Direct and indirect effects. In *Proceedings of the American Statistical Association, Joint*
635 *Statistical Meetings*, pp. 1572–1581. {MIRA} Digital Publishing, Minn., MN, 2005.
- 636
637 Judea Pearl and Dana Mackenzie. *The Book of Why*. Basic Books, New York, 2018. ISBN
638 978-0-465-09760-9.
- 639
640 Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. *Causal inference in statistics: A Primer*.
641 John Wiley & Sons, 2016. ISBN 9781119186847.
- 642
643 Drago Plecko and Elias Bareinboim. Causal fairness analysis: A causal toolkit for fair machine
644 learning. *Foundations and Trends in Machine Learning*, 17(3):304–589, Jan 2024. ISSN 1935-8237.
645 doi: 10.1561/2200000106. URL <https://doi.org/10.1561/2200000106>.
- 646
647 Hans Reichenbach. *The Direction of Time*. University of California Press, Berkeley, 1956.
- Thomas S. Richardson and James M. Robins. Single world intervention graphs (SWIGs): A
unification of the counterfactual and graphical approaches to causality. Working Paper 128, Center
for the Statistics and the Social Sciences, 2013.

648 James M Robins and Sander Greenland. Identifiability and Exchangeability for Direct and Indirect
649 Effects. *Epidemiology*, 3(2):143–155, 1992.

650 Donald B Rubin. Direct and indirect causal effects via potential outcomes. *Scandinavian Journal of*
651 *Statistics*, 31:161–170, 2004.

652 Babak Salimi, Luke Rodriguez, Bill Howe, and Dan Suciu. Interventional fairness: Causal database
653 repair for algorithmic fairness. SIGMOD ’19, pp. 793–810, New York, NY, USA, 2019. Association
654 for Computing Machinery. ISBN 9781450356435. doi: 10.1145/3299869.3319901. URL
655 <https://doi.org/10.1145/3299869.3319901>.

656 Ilya Shpitser. *Complete Identification Methods for Causal Inference*. PhD thesis, Computer Science
657 Department, University of California, Los Angeles, CA, 4 2008.

658 Ilya Shpitser and Judea Pearl. What Counterfactuals Can Be Tested. In *Proceedings of the Twenty-*
659 *Third Conference on Uncertainty in Artificial Intelligence*, pp. 352–359. AUAI Press, Vancouver,
660 BC, Canada, 2007.

661 Ilya Shpitser and Judea Pearl. Effects of Treatment on the Treated: Identification and Generalization.
662 In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, Montreal,
663 Quebec, 2009. AUAI Press.

664 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 1998.

665 Kevin Xia and Elias Bareinboim. Neural causal abstractions. *Proceedings of the AAAI Conference*
666 *on Artificial Intelligence*, 38(18):20585–20595, Mar. 2024. doi: 10.1609/aaai.v38i18.30044. URL
667 <https://ojs.aaai.org/index.php/AAAI/article/view/30044>.

668 Junzhe Zhang and Elias Bareinboim. Transfer learning in multi-armed bandits: A Causal approach.
669 In *IJCAI International Joint Conference on Artificial Intelligence*, 2017. ISBN 9780999241103.

670 Junzhe Zhang and Elias Bareinboim. Can humans be out of the loop? In Bernhard Schölkopf,
671 Caroline Uhler, and Kun Zhang (eds.), *Proceedings of the First Conference on Causal Learning*
672 *and Reasoning*, volume 177 of *Proceedings of Machine Learning Research*, pp. 1010–1025. PMLR,
673 11–13 Apr 2022. URL <https://proceedings.mlr.press/v177/zhang22a.html>.

674 Junzhe Zhang, Jin Tian, and Elias Bareinboim. Partial counterfactual identification from observational
675 and experimental data. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari,
676 Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine*
677 *Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 26548–26558. PMLR,
678 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/zhang22ab.html>.

679 Shoshana Zuboff. *The Age of Surveillance Capitalism: The Fight for a Human Future at the New*
680 *Frontier of Power*. 1st edition, 2018. ISBN 1610395697.

681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

APPENDICES

Appendix A: Graphical terminology

Appendix B: Details on the **CTF-REALIZE** algorithm

Appendix C: Assumptions and realizability proofs

Appendix D: Details on counterfactual randomization

Appendix E: Details on examples

Appendix F: Optimality result and proof

A GRAPHICAL TERMINOLOGY

Structural Causal Models (SCM) and causal diagrams are described in the preliminaries in Sec. 1. See (Bareinboim et al., 2022) for full treatment. We use the following graphical kinship nomenclature w.r.t causal diagram \mathcal{G} :

- Parent(s) of V , denoted \mathbf{Pa}_V : the set of variables $\{V'\}$ s.t. there is a direct edge $V' \rightarrow V$ in \mathcal{G} . \mathbf{Pa}_V does not include V .
- Children of V , denoted $Ch(V)$: the set of variables $\{V'\}$ s.t. there is a direct edge $V \rightarrow V'$ in \mathcal{G} . $Ch(V)$ does not include V .
- Ancestors of V , denoted $An(V)$: the set of variables $\{V'\}$ s.t. there is a path (possibly length 0) from V' to V consisting only of edges pointing toward V , $V' \rightarrow \dots \rightarrow V$. $An(V)$ is defined to include V .
- Descendants of V , denoted $Desc(V)$: the set of variables $\{V'\}$ s.t. there is a path (possibly length 0) from V to V' consisting only of edges pointing toward V' , $V \rightarrow \dots \rightarrow V'$. $Desc(V)$ is defined to include V .
- Non-descendants of V , denoted $NDesc(V)$: the set $\mathbf{V} \setminus Desc(V)$. $NDesc(V)$ does not include V .

Given a graph \mathcal{G} , $\mathcal{G}_{\overline{\mathbf{X}}\mathbf{W}}$ is the result of removing edges coming into variables in \mathbf{X} , and edges coming out of \mathbf{W} .

756 B DETAILS ON THE CTF-REALIZE ALGORITHM

757 B.1 SUB-ROUTINE OF CTF-REALIZE ALGORITHM (ALGO. 1)

760 Algorithm 2 COMPATIBLE (sub-routine)

```

762 1: Input:  $V \in \mathbf{V}$  of  $\mathcal{G}$ ;  $W_t \in \mathbf{W}_*$  of  $Q$            20:   end if
763 2: for each  $C \in Ch(V)$  do                               21:   end if
764 3:   if  $C \in An(W)$  then                                22:   if  $V \notin \mathbf{T}$  then
765 4:     if  $V \in \mathbf{T}$  then                                   23:     for each  $\mathbf{C} \ni C$  s.t.
766 5:       Let  $v :=$  value of  $V$  in subscript  $\mathbf{t}$            24:       CTF-RAND( $V \rightarrow \mathbf{C}$ )  $\in \mathbb{A}$  do
767 6:       Find smallest  $\mathbf{C} \ni C$  s.t.                   25:       if {CTF-RAND( $V \rightarrow \mathbf{C}$ ) :  $\cdot$ }  $\in$ 
768 7:       CTF-RAND( $V \rightarrow \mathbf{C}$ )  $\in \mathbb{A}$                 INT $_V$  and its label is not "Natural"
769 8:       if {CTF-RAND( $V \rightarrow \mathbf{C}$ ) :  $\cdot$ }  $\in$          then
770 9:       INT $_V$  and its label is not " $v$ " then           26:         Return FAIL
771 10:      Return FAIL                                       27:       else
772 11:     else                                               28:         Add {CTF-RAND( $V \rightarrow \mathbf{C}$ ) :
773 12:       Add {CTF-RAND( $V \rightarrow \mathbf{C}$ ) :  $v$ }           Natural} to INT $_V$ , with the la-
774 13:       to INT $_V$ , with the label " $v$ "                 bel "Natural"
775 14:     end if                                           29:       end if
776 15:     if no such  $\mathbf{C} \ni C$  s.t.                       30:     end for
777 16:     CTF-RAND( $V \rightarrow \mathbf{C}$ )  $\in \mathbb{A}$  then         31:     if {RAND( $V$ ) :  $\cdot$ }  $\in$  INT $_V$  and its
778 17:     if {RAND( $V$ ) :  $\cdot$ }  $\in$  INT $_V$  and its             label is not "Natural" then
779 18:     label is not " $v$ " then                             32:       Return FAIL
780 19:     Return FAIL                                       33:     else if RAND( $V$ )  $\in \mathbb{A}$  then
781 20:     else if RAND( $V$ )  $\notin \mathbb{A}$  then                 34:       Add {RAND( $V$ ) : Natural} to
782 21:     Return FAIL                                       35:       INT $_V$ , with the label "Natural"
783 22:     else                                               36:     end if
784 23:       Add {RAND( $V$ ) :  $v$ } to INT $_V$ ,                 37:     end if
785 24:       with the label " $v$ "                               38:   end if
786 25:     end if                                           39: end for

```

787 B.2 WALK-THROUGH OF THE CTF-REALIZE ALGORITHM (ALGO. 1)

788 **General strategy.** For each variable V in the input graph, we check what are the necessary and
789 sufficient interventions (or lack of interventions) we need to perform w.r.t each term W_t in the input
790 query \mathbf{W}_* . This is what the inner loops and subroutine **COMPATIBLE** are doing - accumulating
791 correct and complete conditions in topological order. If there is no conflict across these conditions
792 collectively, and if the feasible action set contains the necessary actions, the query is realizable.
793 Otherwise not. [We use induction to prove that it is enough to check for conflicting conditions with](#)
794 [regard to prior loops, in topological order. Proof of completeness results in Appendix C.](#)

795 Walk-through for Algo. 1 **CTF-REALIZE**:

- 798 i. **Lines 5-7:** we go over each node V in the input graph, in topological order, and maintain a
799 tracker of interventions (and lack of interventions) needed w.r.t V to realize \mathbf{W}_* ; we also
800 check whether V needs to be added to the final output vector.
- 801 ii. **Line 10:** for each W_t in the input query \mathbf{W}_* , we check if there is any conflict in necessary
802 and sufficient conditions w.r.t V for realizing W_t , by calling **COMPATIBLE**(V, W_t). This
803 only needs to be done if $V \in An(W)$; otherwise V has no effect on W_t .
- 804 iii. **Line 13:** if $V = W$, the measured value of V needs to be added to the output vector \mathbf{w} .
- 805 iv. **Lines 17-18:** we perform all the interventions (if any) that are needed for V . [Step \[ii\] has](#)
806 [already checked whether these actions can be performed, and if there are any conflicts.](#)
807 – **Note on rejection sampling:** since we framed our actions as randomizations, in order
808 to enact an intervention like $do(x)$, we draw a random value and reject if the draw is not
809 x . This is for clarity of presentation, aligned with the rest of the paper. We could easily

810 introduce deterministic actions, or add some concentration guarantees of drawing x
811 within finite samples etc. but this is well-understood and would be a distraction.
812
813 v. **Lines 21-25:** if V itself is part of the output, the set of necessary actions cannot involve a
814 Fisherian RAND of V (because unlike CTF-RAND, Fisherian RAND overrides the mecha-
815 nism f_V generating V).
816 vi. **Line 29:** if the above steps have been completed w.r.t V for each W_t , there will be no further
817 conflicts arising w.r.t V and all nodes topologically prior to V , regardless of conditions
818 needed w.r.t subsequent nodes (by an induction argument). If this is can be completed for
819 all V , then the query is realizable and we output the vector w .

820 Walk-through for sub-routine Algo. 2 **COMPATIBLE**(V, W_t) called in Step [ii]:
821

822 ii.a. **Lines 2-3:** the necessary and sufficient conditions w.r.t V for W_t involve how the children
823 $Ch(V)$ receive the value of V as an input. We only care about the children that belong in
824 $An(W)$ for this sub-routine call; if a child is not in $An(W)$ it wouldn't affect W_t .
825
826 ii.b. **Lines 4-19:** if $V \in \mathbf{T}$, this means the potential response W_t involves an intervention on
827 V . We find the minimal interventions needed to achieve this (CTF-RAND for the smallest
828 subset of children possible, and failing this, a Fisherian RAND); and we update our tracker
829 of necessary actions for V .
830
831 ii.b.1. The necessary and sufficient conditions to measure W_t are that \mathbf{T} should be fixed
832 as t (by intervention) as an input to all children $C \in Ch(\mathbf{T}) \cap An(W)$; and that all
833 other ancestors of W are received naturally by their respective children. These two
834 conditions ensure that W is evaluated in the $do(t)$ regime, as defined in the SCM.
835
836 ii.b.2. In particular, for $V \in \mathbf{T}$ in this sub-routine call, this means that we need to fix $v \in \mathbf{t}$
837 (**Line 5**) as input to each relevant child $C \in Ch(V) \cap An(W)$. So for each such child
838 C , we find the smallest subset $C' \subseteq Ch(V)$ s.t. $C \in C'$ and CTF-RAND($V \rightarrow C'$) is
839 in the input action set \mathbb{A} (**Line 6**), and we add this to the required actions, along with the
840 tagged value " v " that needs to be fixed for V by this action (**Line 10**). If CTF-RAND is
841 not available, we add Fisherian RAND(V) to the required action list (**Line 18**). We call
842 this chosen action the "minimal action".
843
844 ii.b.3. If no such randomization action is available, return FAIL (**Line 16**). If the minimal
845 action is already being used in a previous loop to enforce a different value $v' \neq v$ that
846 is not compatible with t , return FAIL (**Lines 7-8, 13-14**). Or if the minimal action
847 has already been tagged with the value "Natural" in a previous loop, return FAIL -
848 this means a previous loop already recorded that this action must *necessarily not* be
849 performed (**Lines 7-8, 13-14**). Such a conflict means that W_t cannot be realised.
850
851 ii.c. **Lines 22-34:** if $V \notin \mathbf{T}$, this means the potential response W_t requires that V be received
852 without intervention (i.e. "naturally") by the relevant child nodes; we also add this necessary
853 condition to our tracker.
854
855 ii.c.1. Again, the necessary and sufficient conditions to measure W_t are that each ancestor
856 $A \in An(W)_{\mathcal{G}_{\mathbf{T}}}$, with $A \notin \mathbf{T} \cup \{W\}$, needs to be received "naturally" (i.e., without
857 intervention) by its children $C \in Ch(A) \cap An(W)$; and that \mathbf{T} needs to be fixed by
858 intervention, as discussed earlier in [ii.b.1]. These two conditions ensure that W is
859 evaluated in the $do(t)$ regime, as defined in the SCM.
860
861 ii.c.2. In particular, for $V \notin \mathbf{T}$ in this sub-routine, this means that we should *necessarily not*
862 intervene on the value of V that is received as input by each relevant child $C \in Ch(V)$
863 which is also an ancestor of W .
864
865 ii.c.3. We track this requirement by adding that for every possible randomization
866 CTF-RAND($V \rightarrow C'$) and RAND(V), where C' contains an ancestor of W , that
867 this action *cannot* be performed, by using the tag "Natural" (**Lines 27, 33**).
868
869 ii.c.4. If any such action identified in [ii.c.3] has already been tagged with some value v' ,
870 this means that a previous loop has recorded that this action *necessarily* needs to be
871 performed in order to fix the value as v' - we return FAIL (**Lines 24-25, 30-31**). Such a
872 conflict means that W_t cannot be realized.

864 B.3 COMPLEXITY ANALYSIS OF **CTF-REALIZE** ALGORITHM (ALGO. 1)

865 The *time complexity* of Algorithm 1 is $\mathcal{O}(kn^2)$, where k is the number of terms in the input query Q ,
 866 and n is the number of variables in the input graph \mathcal{G} .

867 k depends on the domain size of the variables. That is, $k \leq n \cdot \prod_V \left(|Domain(V)| + 1 \right) = \mathcal{O}(m^n)$,
 868 where m is the domain size of the variable with the most possible categorical values.
 869

870 The *space complexity* is the same, as the algorithm needs to store up to all intermediate steps before
 871 terminating.

874 B.4 EXAMPLES USING THE **CTF-REALIZE** ALGORITHM

875 **Example B.1.** (ETT realizability)

876 Query, $Q = P(Y_x, X)$

877 Graph, \mathcal{G} : Fig. 7



882 Figure 7: Graph for Example B.1

883

884 Suppose action set $\mathbb{A} = \mathbb{A}^\dagger(\mathcal{G}) := \{\text{CTF-RAND}(X \rightarrow Y)\}$

885 **CTF-REALIZE**($Q, \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G})$) trace:

- 886
- 887 • Start with X (first in topological order)
 - 888 • For the first term in \mathbf{W}_* : Y_x
 - 889 – Since $X \in An(Y)$, call Algo. 2 **COMPATIBLE**(X, Y_x)
 - 890 * $Y \in Ch(X)$ and $Y \in An(Y)$
 - 891 * $X \in \text{subscript of } Y_x$
 - 892 * $\text{CTF-RAND}(X \rightarrow Y) \in \mathbb{A}^\dagger(\mathcal{G})$
 - 893 * $\text{INT}_X \leftarrow \{\text{CTF-RAND}(X \rightarrow Y) : x\}$
 - 894 • For the second term in \mathbf{W}_* : X
 - 895 – $\text{OUTPUT}_X \leftarrow \{X\}$
 - 896 • Moving to Y (next in topological order)
 - 897 • For the first term in \mathbf{W}_* : Y_x
 - 898 – $\text{OUTPUT}_Y \leftarrow \{Y_x\}$
 - 899 • Perform interventions in INT_X , followed by READ, and assign output vector based on
 900 $\text{OUTPUT}_X, \text{OUTPUT}_Y$
 - 901 • Return i.i.d sample

902 For simplicity, we don't show the steps $\text{SELECT}^{(i)}$ and the rejection sampling involving in the
 903 randomization procedure (steps 17-18 of Algo. 1).

904 Thus, Q is realizable given $\mathcal{G}, \mathbb{A}^\dagger$. This is validated by the ancestor set $An(Y_x, X)_\mathcal{G} = \{Y_x, X\}$,
 905 which doesn't repeat any variables. This is also illustrated in Fig. 8.

906 However, suppose the agent's action set is

907 $\mathbb{A} = \{\text{RAND}(X)\}$, i.e., does not permit any counterfactual randomization procedures.

908 In this case,

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

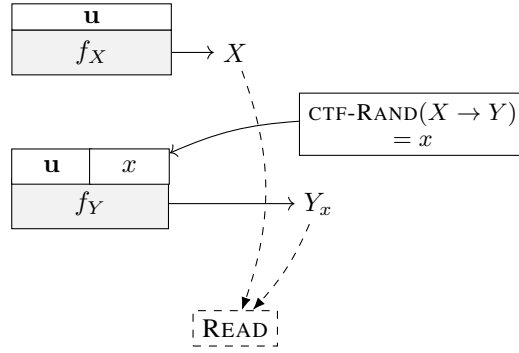


Figure 8: $P(Y_x, X)$ is realizable given the graph in Fig. 7 and $\mathbb{A}^\dagger(\mathcal{G})$.

CTF-REALIZE($Q, \mathcal{G}, \mathbb{A}$) trace:

- Start with X (first in topological order)
- For the first term in \mathbf{W}_* : Y_x
 - Since $X \in An(Y)$, call Algo. 2 **COMPATIBLE**(X, Y_x)
 - * $Y \in Ch(X)$ and $Y \in An(Y)$
 - * $X \in$ subscript of Y_x
 - * $RAND(X) \in \mathbb{A}$, and no other ctf-randomization procedure
 - * $INT_X \leftarrow \{RAND(X) : x\}$
- For the second term in \mathbf{W}_* : X
 - $OUTPUT_X \leftarrow \{X\}$
- Moving to Y (next in topological order)
- For the first term in \mathbf{W}_* : Y_x
 - $OUTPUT_Y \leftarrow \{Y_x\}$
- $OUTPUT_X$ contains X , but the intervention set INT_X contains $RAND(X)$
- **FAIL** (Line 22 of Algo. 1)

■

Example B.2. (Probability of sufficiency (PS) realizability)

Query, $Q = P(Y_x, X, Y)$

Graph, \mathcal{G} : Fig. 9

$$X \longrightarrow Y$$

Figure 9: Graph for Example B.2

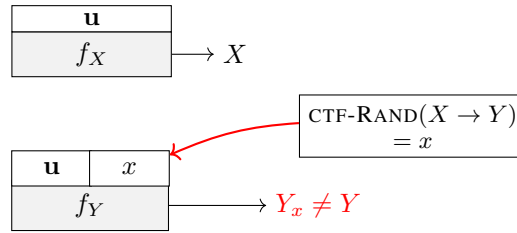
Suppose action set $\mathbb{A} = \mathbb{A}^\dagger(\mathcal{G}) := \{\text{CTF-RAND}(X \rightarrow Y)\}$

CTF-REALIZE($Q, \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G})$) trace:

- Start with X (first in topological order)
- For the first term in \mathbf{W}_* : Y_x
 - Since $X \in An(Y)$, call Algo. 2 **COMPATIBLE**(X, Y_x)

- 972 * $Y \in Ch(X)$ and $Y \in An(Y)$
- 973 * $X \in \text{subscript of } Y_x$
- 974 * $\text{CTF-RAND}(X \rightarrow Y) \in \mathbb{A}^\dagger(\mathcal{G})$
- 975 * $\text{INT}_X \leftarrow \{\text{CTF-RAND}(X \rightarrow Y) : x\}$
- 976
- 977 • For the second term in \mathbf{W}_* : X
- 978 – $\text{OUTPUT}_X \leftarrow \{X\}$
- 979
- 980 • For the third term in \mathbf{W}_* : Y
- 981 – Since $X \in An(Y)$, call Algo. 2 **COMPATIBLE**(X, Y)
- 982 * $Y \in Ch(X)$ and $Y \in An(Y)$
- 983 * $X \notin \text{subscript of } Y$; X needs to be received naturally
- 984 * But INT_X already contains $\{\text{CTF-RAND}(X \rightarrow Y) : x\}$ with label $x \neq \text{"Natural"}$
- 985 * **FAIL** (Line 25 of Algo. 2)
- 986
- 987

988 Thus, Q is not realizable given $\mathcal{G}, \mathbb{A}^\dagger$. This is validated by the ancestor set $An(Y_x, X, Y)_{\mathcal{G}} =$
 989 $\{Y_x, X, Y\}$, which contains both Y_x, Y . This is also illustrated in Fig. 10.



990 Figure 10: $P(Y_x, X, Y)$ is not realizable given the graph in Fig. 9 and $\mathbb{A}^\dagger(\mathcal{G})$.

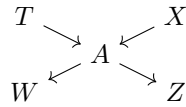
1000

1001

1002

1003 **Example B.3.** Query, $Q = P(W_{xt}, Z_{x'})$

1004 Graph, \mathcal{G} : Fig. 11



1005

1006

1007

1008

1009

1010

1011 Figure 11: Graph for Example B.3

1012 Suppose action set $\mathbb{A} = \mathbb{A}^\dagger(\mathcal{G}) := \{\text{CTF-RAND}(T \rightarrow A), \text{CTF-RAND}(X \rightarrow A), \text{CTF-RAND}(A \rightarrow$
 1013 $W), \text{CTF-RAND}(A \rightarrow Z)\}$

1014 **CTF-REALIZE**($Q, \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G})$) trace:

- 1015
- 1016 • Start with X (first in topological order)
- 1017
- 1018 • For the first term in \mathbf{W}_* : W_{xt}
- 1019 – Since $X \in An(W)$, call Algo. 2 **COMPATIBLE**(X, W_{xt})
- 1020 * $A \in Ch(X)$ and $A \in An(W)$
- 1021 * $X \in \text{subscript of } W_{xt}$
- 1022 * $\text{CTF-RAND}(X \rightarrow A) \in \mathbb{A}^\dagger(\mathcal{G})$
- 1023 * $\text{INT}_X \leftarrow \{\text{CTF-RAND}(X \rightarrow A) : x\}$
- 1024
- 1025 • For the second term in \mathbf{W}_* : $Z_{x'}$

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

- Since $X \in An(Z)$, call Algo. 2 **COMPATIBLE**($X, Z_{x'}$)
 - * $A \in Ch(X)$ and $A \in An(Z)$
 - * $X \in$ subscript of $Z_{x'}$; X needs to be fixed as x'
 - * But INT_X already contains $\{CTF-RAND(X \rightarrow A) : x\}$ with label $x \neq x'$
 - * **FAIL** (Line 8 of Algo. 2)

Thus, Q is not realizable given $\mathcal{G}, \mathbb{A}^\dagger$. This is validated by the ancestor set $An(W_{xt}, Z_{x'})_{\mathcal{G}} = \{W_{xt}, A_{xt}, Z_{x'}, A_{x'}\}$, which contains both $A_{xt}, A_{x'}$. This is also illustrated in Fig. 12.

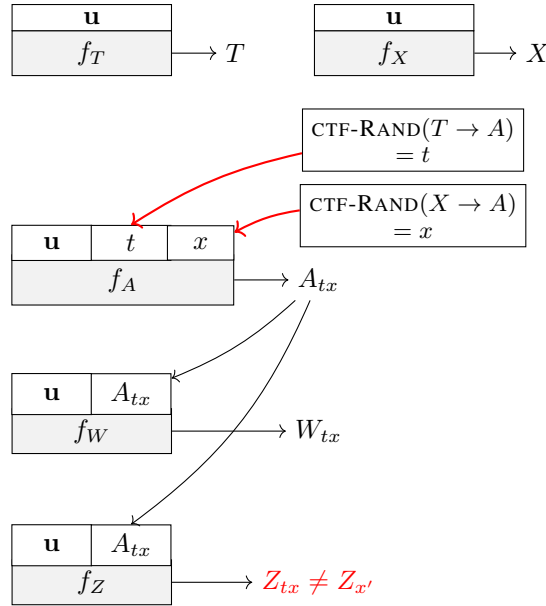


Figure 12: $P(W_{xt}, Z_{x'})$ is not realizable given the graph in Fig. 11 and $\mathbb{A}^\dagger(\mathcal{G})$.

Example B.4. Query, $Q = P(Y_x, Z_{x'}, W_{x''})$

Graph, \mathcal{G} : Fig. 13

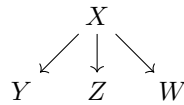


Figure 13: Graph for Example B.4

Suppose action set $\mathbb{A} = \{\text{RAND}(X), \text{CTF-RAND}(X \rightarrow \{Z, W\})\}$

CTF-REALIZE($Q, \mathcal{G}, \mathbb{A}$) trace:

- Start with X (first in topological order)
- For the first term in \mathbf{W}_* : Y_x
 - Since $X \in An(Y)$, call Algo. 2 **COMPATIBLE**(X, Y_x)
 - * $Y \in Ch(X)$ and $Y \in An(Y)$
 - * $X \in$ subscript of Y_x
 - * $\text{RAND}(X) \in \mathbb{A}$; no other ctf-randomization procedures affecting Y
 - * $INT_X \leftarrow \{\text{RAND}(X) : x\}$

-
- 1080 • For the second term in \mathbf{W}_* : $Z_{x'}$
1081
1082 – Since $X \in An(Z)$, call Algo. 2 **COMPATIBLE**($X, Z_{x'}$)
1083 * $Z \in Ch(X)$ and $Z \in An(Z)$
1084 * $X \in$ subscript of $Z_{x'}$
1085 * $\text{CTF-RAND}(X \rightarrow \{Z, W\}) \in \mathbb{A}$; no smaller ctf-randomization procedures affect-
1086 ing Z
1087 * $\text{INT}_X \leftarrow \text{INT}_X \cup \{\text{CTF-RAND}(X \rightarrow \{Z, W\}) : x'\}$
1088
1089 • For the third term in \mathbf{W}_* : $W_{x''}$
1090 – Since $X \in An(W)$, call Algo. 2 **COMPATIBLE**($X, W_{x''}$)
1091 * $W \in Ch(X)$ and $W \in An(W)$
1092 * $X \in$ subscript of $W_{x''}$; X needs to be fixed as x''
1093 * But INT_X already contains $\{\text{CTF-RAND}(X \rightarrow \{Z, W\}) : x'\}$ with label $x' \neq x''$
1094 * No smaller ctf-randomization procedures affecting W
1095 * **FAIL** (Line 8 of Algo. 2)
1096

1097 Thus, Q is not realizable given \mathcal{G}, \mathbb{A} .

1098 However, suppose instead that action set $\mathbb{A}' = \{\text{RAND}(X), \text{CTF-RAND}(X \rightarrow$
1099 $\{Z, W\}), \text{CTF-RAND}(X \rightarrow Z)\}$

1100 **CTF-REALIZE**($Q, \mathcal{G}, \mathbb{A}'$) trace:

- 1101
1102 • Start with X (first in topological order)
1103
1104 • For the first term in \mathbf{W}_* : Y_x
1105 – Since $X \in An(Y)$, call Algo. 2 **COMPATIBLE**(X, Y_x)
1106 * $Y \in Ch(X)$ and $Y \in An(Y)$
1107 * $X \in$ subscript of Y_x
1108 * $\text{RAND}(X) \in \mathbb{A}$; no other ctf-randomization procedures affecting Y
1109 * $\text{INT}_X \leftarrow \{\text{RAND}(X) : x\}$
1110
1111 • For the second term in \mathbf{W}_* : $Z_{x'}$
1112 – Since $X \in An(Z)$, call Algo. 2 **COMPATIBLE**($X, Z_{x'}$)
1113 * $Z \in Ch(X)$ and $Z \in An(Z)$
1114 * $X \in$ subscript of $Z_{x'}$
1115 * $\text{CTF-RAND}(X \rightarrow Z) \in \mathbb{A}$
1116 * $\text{INT}_X \leftarrow \text{INT}_X \cup \{\text{CTF-RAND}(X \rightarrow Z) : x'\}$
1117
1118 • For the third term in \mathbf{W}_* : $W_{x''}$
1119 – Since $X \in An(W)$, call Algo. 2 **COMPATIBLE**($X, W_{x''}$)
1120 * $W \in Ch(X)$ and $W \in An(W)$
1121 * $X \in$ subscript of $W_{x''}$
1122 * $\text{CTF-RAND}(X \rightarrow \{Z, W\}) \in \mathbb{A}$
1123 * $\text{INT}_X \leftarrow \text{INT}_X \cup \{\text{CTF-RAND}(X \rightarrow \{Z, W\}) : x''\}$
1124
1125 • Moving to Y (next in topological order)
1126 – $\text{OUTPUT}_Y \leftarrow \{Y_x\}$
1127
1128 • Moving to Z (next in topological order)
1129 – $\text{OUTPUT}_Z \leftarrow \{Z_{x'}\}$
1130
1131 • Moving to W (next in topological order)
1132 – $\text{OUTPUT}_W \leftarrow \{W_{x''}\}$
1133

- Perform interventions in INT_X , followed by READ, and assign output vector based on $\text{OUTPUT}_Y, \text{OUTPUT}_Z, \text{OUTPUT}_W$
- Return i.i.d sample

Thus, Q is realizable given \mathcal{G}, \mathbb{A}' .

Lastly, it is evident that Q is realizable given $\mathcal{G}, \mathbb{A}^\dagger$. This is validated by the ancestor set $An(Y_x, Z_{x'}, W_{x''})_{\mathcal{G}} = \{Y_x, Z_{x'}, W_{x''}\}$, which does not repeat any variables. This is also illustrated in Fig. 14.

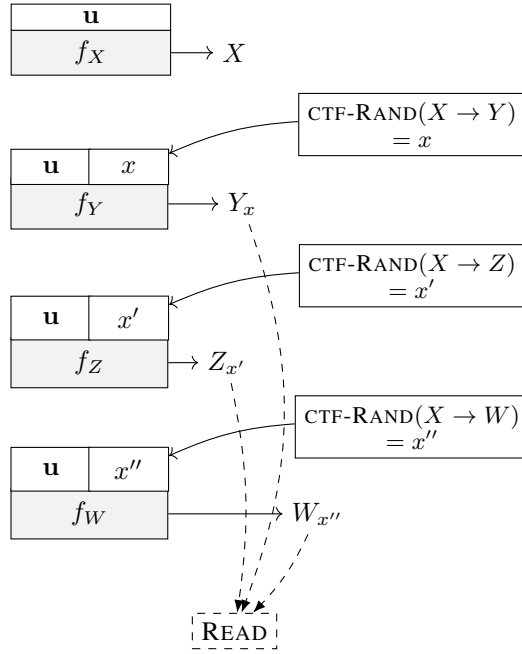


Figure 14: $P(Y_x, Z_{x'}, W_{x''})$ is realizable given the graph in Fig. 13 and $\mathbb{A}^\dagger(\mathcal{G})$.

■

1188 C ASSUMPTIONS AND REALIZABILITY PROOFS

1189 C.1 ASSUMPTIONS

1190 In this section, we gather together all the structural assumptions we make in this paper, for ease of
1191 reference. We also include related remarks.

1192 **Assumption C.1** (Unobservability). An agent deployed in the environment does not know the
1193 underlying SCM \mathcal{M} of the environment, and does not know the latent features $\mathbf{U}^{(i)}$ of any unit i in
1194 the target population. ■

1195 **Assumption C.2** (Feasible actions). Given causal diagram \mathcal{G} , the physical actions that an agent can
1196 perform on any unit i in the target population are limited to: $\text{SELECT}^{(i)}$, $\text{READ}(V)^{(i)}$, $\text{RAND}(X)^{(i)}$,
1197 and $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$, for some $V, X \in \mathbf{V}$ and $\mathbf{C} \subseteq \text{Ch}(X)_{\mathcal{G}}$, per Defs. 2.1,2.3. ■

1200 **Assumption 3.1** (Fundamental constraint of experimentation (FCE)). A unit i in the target population
1201 can physically undergo a causal mechanism $f_V \in \mathcal{F}$ at most once. ■

1202 We define the probability measure $P^{\mathbb{C}}(\cdot)$ from the perspective of an exogenous agent (i.e., an agent
1203 external to the system) \mathbb{C} 's beliefs about the environment, distinguished by superscript from $P^{\mathcal{M}}(\cdot)$,
1204 the true unknown distribution.

1205 *Remark C.3.* Let $\mathcal{A}^{(i)}$ be a sequence of actions taken by agent \mathbb{C} on unit i that is not conditional on
1206 any data gathered regarding i . The assumption of \mathbb{C} behaving exogenously means that $P^{\mathbb{C}}(\mathbf{U}^{(i)} =$
1207 $\mathbf{u} \mid \mathcal{A}^{(i)}) = P^{\mathcal{M}}(\mathbf{u})$. ■

1208 Regarding the structural conditions involving counterfactual randomization (Def. 2.3), we make the
1209 following assumption, mainly as a simplifying step for use in the proofs.

1210 **Assumption D.3** (Tree structure). Given a variable X , causal diagram \mathcal{G} , and an "expanded" diagram
1211 \mathcal{G}^+ (Def. D.1) including the set of all the counterfactual mediators \mathbf{W} (Def. D.2) of X in the
1212 environment, each $W \in \mathbf{W}$ has only one parent in \mathcal{G}^+ , and each $C \in \text{Ch}(X)_{\mathcal{G}}$ has at most one
1213 $W \in \mathbf{W}$ as a parent in \mathcal{G}^+ . ■

1214 From this assumption, and from the definition of a counterfactual mediator (Def. D.2), we can derive
1215 the following observations:

1216 *Remark C.4* (No bypassing children). Given causal diagram \mathcal{G} , the procedure $\text{CTF-RAND}(X \rightarrow \mathbf{C})$,
1217 either by eliciting a unit's natural decision or via a counterfactual mediator, can only be performed
1218 w.r.t $\mathbf{C} \subseteq \text{Ch}(X)_{\mathcal{G}}$. It cannot by-pass child mechanisms and directly affect a descendant. This is
1219 elaborated in App. D.3, and specifically in Lemma D.7. ■

1220 *Remark D.6* (Procedure containment). Assumption D.3 implies that if an agent is capable of performing
1221 both $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$ and $\text{CTF-RAND}(X \rightarrow \mathbf{C}')^{(i)}$ s.t. $\mathbf{C} \neq \mathbf{C}'$ and $\mathbf{C} \cap \mathbf{C}' \neq \emptyset$, then
1222 either $\mathbf{C} \subseteq \mathbf{C}'$ or $\mathbf{C}' \subseteq \mathbf{C}$. ■

1223 *Remark D.5* (Superseding action). Given a decision variable X , the action $\text{CTF-RAND}(X \rightarrow \mathbf{C}')^{(i)}$
1224 can *supersede* the action $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$ if $\mathbf{C}' \subsetneq \mathbf{C}$, where *supersede* means that the
1225 former action $\text{CTF-RAND}(X \rightarrow \mathbf{C}')^{(i)}$ blocks any effect that the latter action has on the variables \mathbf{C}' .
1226 Additionally, the action $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$ *supersedes* the action $\text{RAND}(X)^{(i)}$. ■

1227 Counterfactual randomization permits multiple randomizations for the same variable X for a single
1228 unit i . But some randomizations block the effects of others. See App. D.2.

1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

1242 C.2 PROOFS FOR SECTION 3

1243
1244 Recall, $P^{\mathbb{C}}(\cdot)$ is the probability measure from the perspective of an exogenous agent (i.e., an agent
1245 external to the system) \mathbb{C} 's beliefs about the environment, distinguished by superscript from $P^{\mathcal{M}}(\cdot)$,
1246 the true unknown distribution.

1247 Since unit selection is randomized, $\text{SELECT}^{(i)}$ yields an unbiased sample of a unit with latent features
1248 distributed according to the target population frequency $P(\mathbf{u})$. I.e., $P^{\mathbb{C}}(\mathbf{U}^{(i)} = \mathbf{u} \mid \text{SELECT}^{(i)}) =$
1249 $P^{\mathcal{M}}(\mathbf{u})$.

1250 **Lemma C.5** (I.i.d requirement). *Consider a sequence of actions $\mathcal{A}^{(i)}$ performed on unit i in the target*
1251 *population, that yields a vector of realized values $\mathbf{W}_*^{(i)}$. $\mathbf{W}_*^{(i)}$ is an i.i.d sample from $P^{\mathcal{M}}(\mathbf{W}_*)$, for*
1252 *arbitrary \mathcal{M} iff*

- 1254 i. $P^{\mathbb{C}}(\mathbf{U}^{(i)} = \mathbf{u} \mid \mathcal{A}^{(i)}) = P^{\mathcal{M}}(\mathbf{U} = \mathbf{u})$; and
1255
1256 ii. $\mathbb{1}[\mathbf{W}_*^{(i)} = \mathbf{w} \mid \mathcal{A}^{(i)}, \mathbf{U}^{(i)} = \mathbf{u}] = \mathbb{1}[\mathbf{W}_*(\mathbf{u}) = \mathbf{w}]$.

1257
1258 *Proof.* Recall from Def. 3.3 that $\mathbf{W}_*^{(i)}$ being an i.i.d sample from $P^{\mathcal{M}}(\mathbf{W}_*)$ means that

$$1259 P^{\mathbb{C}}(\mathbf{W}_*^{(i)} = \mathbf{w} \mid \mathcal{A}^{(i)}) = P^{\mathcal{M}}(\mathbf{W}_* = \mathbf{w}), \forall \mathbf{w} \quad (9)$$

1260
1261
1262 Reverse direction:

1263 We simply multiply respective l.h.s and r.h.s of conditions [i] and [ii] and sum over all \mathbf{u} to get

$$1264 \sum_{\mathbf{u}} P^{\mathbb{C}}(\mathbf{U}^{(i)} = \mathbf{u} \mid \mathcal{A}^{(i)}) \cdot \mathbb{1}[\mathbf{W}_*^{(i)} = \mathbf{w} \mid \mathcal{A}^{(i)}, \mathbf{U}^{(i)} = \mathbf{u}] = \sum_{\mathbf{u}} P^{\mathcal{M}}(\mathbf{U} = \mathbf{u}) \cdot \mathbb{1}[\mathbf{W}_*(\mathbf{u}) = \mathbf{w}]$$

$$1265 \quad \quad \quad (10)$$

$$1266 \quad \quad \quad = P^{\mathcal{M}}(\mathbf{W}_* = \mathbf{w}), \quad (11)$$

1267
1268 which we get from the Layer 3 valuation formula (see preliminaries in Sec. 1). On the l.h.s, we apply
1269 the chain rule to get the result we need.

$$1270 P^{\mathbb{C}}(\mathbf{W}_*^{(i)} = \mathbf{w} \mid \mathcal{A}^{(i)}) = P^{\mathcal{M}}(\mathbf{W}_* = \mathbf{w}) \quad (12)$$

1271
1272 Forward direction:

1273 Assume Eq. 9.

1274 From Remark C.3, we conclude that condition [i] automatically holds. Since the agent acts exoge-
1275 nously to the system, $P^{\mathbb{C}}(\mathbf{U}^{(i)} = \mathbf{u} \mid \mathcal{A}^{(i)}) = P^{\mathcal{M}}(\mathbf{U} = \mathbf{u}), \forall \mathbf{u}$.

1276 Applying the chain rule on both sides of Eq. 9,

$$1277 \sum_{\mathbf{u}} P^{\mathbb{C}}(\mathbf{U}^{(i)} = \mathbf{u} \mid \mathcal{A}^{(i)}) \cdot \mathbb{1}[\mathbf{W}_*^{(i)} = \mathbf{w} \mid \mathcal{A}^{(i)}, \mathbf{U}^{(i)} = \mathbf{u}] = \sum_{\mathbf{u}} P^{\mathcal{M}}(\mathbf{U} = \mathbf{u}) \cdot \mathbb{1}[\mathbf{W}_*(\mathbf{u}) = \mathbf{w}]$$

$$1278 \quad \quad \quad (13)$$

1279 The probability terms are equal, for each \mathbf{u} .

1280 Since the probability terms are free parameters, and we need this equation to hold for any arbitrary
1281 probability simplex, it must be the case that the indicator terms are also equal. Thus begetting
1282 condition [ii].

1283 ■

1284
1285 **Lemma C.6.** *Given a causal diagram \mathcal{G} , for any SCM \mathcal{M} compatible with \mathcal{G} , the jointly necessary and*
1286 *sufficient conditions to measure a potential response $W_{\mathbf{t}}(\mathbf{u})$ are [i] \mathbf{T} is fixed as \mathbf{t} (by intervention) as*
1287 *an input to all children $C \in \text{Ch}(\mathbf{T}) \cap \text{An}(W)$; [ii] each $A \in \text{An}(W)_{\mathcal{G}_{\overline{\mathbf{T}}}}$, $A \notin \{\mathbf{T}, W\}$ is received*
1288 *"naturally" (i.e., without intervention) by its children $C \in \text{Ch}(A) \cap \text{An}(W)$; and [iii] the mechanism*
1289 *f_W is not erased and overwritten (by a Fisherian intervention).*

1296 *Proof.* W_t is the variable W evaluated in the sub-model \mathcal{M}_t , where the equations for \mathbf{T} are replaced
 1297 by constant values in \mathbf{t} .

1298 For any changes to the function for $T \in \mathbf{T}$, the function f_W is only affected by any effect on the
 1299 children of T which are also ancestors of W . Any effect of T on some $C' \in Ch(T)$ s.t. $C' \notin An(W)$
 1300 has no effect on W .

1301 Further, in the submodel \mathcal{M}_t there are no interventions on any other ancestors of W in $\mathcal{G}_{\mathbf{T}}$, besides
 1302 \mathbf{T} . Even if there were interventions involving some $X \notin An(W)_{\mathcal{G}_{\mathbf{T}}}$, this would have no effect on W
 1303 in the sub-model \mathcal{M}_t , by Rule 3 of do-calculus.

1304 It is evident that f_W evaluated according to the sub-model \mathcal{M}_t , and evaluated according to a sub-
 1305 model satisfying conditions [i] and [ii] are identical for each \mathbf{u} , since the sequence of structural
 1306 equations that eventually generate W are the same.

1307 Finally, in order to measure W_t , we need to measure the output of the mechanism f_W in the real
 1308 world. The mechanism cannot not be erased and overwritten, as per condition [iii]. ■

1309
 1310
 1311
 1312 **Lemma C.7.** *Given a set \mathbf{W}_* and graph \mathcal{G} , where each member $W_t \in \mathbf{W}_*$ has its respective*
 1313 *conditions [i-iii] (per Lemma C.6), suppose these conditions introduce conflicts when combined*
 1314 *across \mathbf{W}_* . Removing X from \mathbf{W}_* and from all subscripts in \mathbf{W}_* to get a new set \mathbf{W}'_* does not*
 1315 *introduce new conflicts between the terms, if X is first in a topological ordering of \mathcal{G} .*

1316
 1317 *Proof.* Let us consider the conditions in the necessary-and-sufficient set given in Lemma C.6 for
 1318 each $W_t \in \mathbf{W}_*$.

1319 Condition [i] add requirements for each $t \in \mathbf{t}$ of some $W_t \in \mathbf{W}_*$. Since X is removed from every
 1320 subscript, this no longer applies to any term in \mathbf{W}'_* .

1321 Condition [ii] requires that if X is an ancestor of some $W_t \in \mathbf{W}_*$, and it doesn't appear in \mathbf{T} , then
 1322 X must be received without intervention by mediating children. Since X is removed from every
 1323 subscript, X not being intervened upon at all meets this condition [ii] for every term \mathbf{W}'_* , without
 1324 conflicting with condition [i], which no longer applies.

1325 Importantly, even though X is no longer being intervened upon, for each $W_t \in \mathbf{W}_*$, removing
 1326 X from \mathbf{T} (if it appears) does not add any additional ancestors in $\mathcal{G}_{\mathbf{T}}$ that need to be tracked for
 1327 condition [ii], since X is first in topological order.

1328 Condition [iii] would only apply to X itself, which is not present as a potential response in \mathbf{W}'_* . ■

1329
 1330
 1331
 1332
 1333 **Theorem 3.5 :**

1334 Let $\mathcal{A}^{(i)}$ be a sequence of actions conducted by an exogenous agent to beget a vector of values $\mathbf{W}_*^{(i)}$
 1335 for a unit i .

1336 By Lemma C.5, if an agent wants $\mathbf{W}_*^{(i)}$ to be an i.i.d sample from $P(\mathbf{W}_*)$, then for each possible
 1337 $\mathbf{U}^{(i)} = \mathbf{u}$, the vector $\mathbf{W}_*^{(i)}$ needs to be identical to the $\mathbf{W}_*(\mathbf{u})$ as evaluated according to the SCM.
 1338 Essentially, this says that the agent's actions need to output the same vector $\mathbf{W}_*(\mathbf{u})$ as if it has been
 1339 evaluated according to the SCM.

1340 By Lemma C.6, $\mathbf{W}_*(\mathbf{u})$ can be evaluated if and only if the following three conditions are met for
 1341 each $W_t \in \mathbf{W}_*$ simultaneously:

- 1342 i \mathbf{T} is fixed as \mathbf{t} (by intervention) as an input to all children $C \in Ch(\mathbf{T}) \cap An(W)$;
- 1343 ii Each $A \in An(W)_{\mathcal{G}_{\mathbf{T}}}$, $A \notin \{\mathbf{T}, W\}$ is received "naturally" (i.e., without intervention) by
 1344 its children $C \in Ch(A) \cap An(W)$; and
- 1345 iii The mechanism f_W is not erased and overwritten (by a Fisherian intervention).

1346
 1347
 1348 **Inductive hypothesis (IH):**

1350 **CTF-REALIZE**($P(\mathbf{W}_*), \mathcal{G}, \mathbb{A}$) returns FAIL if and only if conditions [i-iii] are not met simultane-
1351 ously, when combined across all $W_t \in \mathbf{W}_*$ w.r.t a causal diagram \mathcal{G} having $\leq n$ nodes
1352

1353 **Base case:**

1354 Consider an SCM with only one variable $V \in \mathbf{V}$. IH is trivially true, since the conditions are always
1355 met, and since **CTF-REALIZE** will just return the value $\text{READ}(V)$.
1356 Assume IH is true for any SCM with causal diagram having $\leq n$ nodes.
1357

1358 **n+1 case:**

1359 Consider an SCM whose causal diagram \mathcal{G} has $n + 1$ nodes. Let X be the first in some topological
1360 ordering of \mathcal{G} . Consider an action set \mathbb{A} that the agent can perform in the environment, and an
1361 arbitrary distribution $P(\mathbf{W}_*)$.
1362 WLOG, we can begin the outer loop of **CTF-REALIZE**($P(\mathbf{W}_*), \mathcal{G}, \mathbb{A}$) with X (first in topological
1363 order).
1364

- 1365 • The inner loop calls **COMPATIBLE**(X, W_t) for each $W \in \text{Desc}(X), X \neq W$.
- 1366 • It maintains a tracker INT_X of the smallest counterfactual interventions needed to satisfy
1367 condition [i] for each W_t , resorting to Fisherian intervention if needed. Note (per Remark
1368 D.6), interventions follow a tree-like structure, so conflicts can be tracked by tagging the
1369 smallest available intervention that is needed for each W_t w.r.t each child of X .
- 1370 • Note also (per Remark D.5) that if there are two simultaneous interventions added to INT_X ,
1371 $\text{CTF-RAND}(X \rightarrow C)$, $\text{CTF-RAND}(X \rightarrow C')$, where $C' \subseteq C$, then the set C' is unaffected
1372 by the first procedure.
1373
- 1374 • This inner loop exactly checks if there are any conflicts in conditions [i-ii] among \mathbf{W}_*
1375 w.r.t X , by "tagging" each procedure with the fixed value x needed for that intervention
1376 (including the requirement of no intervention).
- 1377 • Finally the outer loop in Line 20 of Algo. 1 checks if X appears as a potential response
1378 anywhere in \mathbf{W}_* . If so, INT_X cannot contain the requirement of Fisherian $\text{RAND}(X)$, since
1379 this violates condition [iii] w.r.t X .
- 1380 • X does not appear anywhere else in subsequent algorithm iterations.

1381

1382 Thus, we conclude that **CTF-REALIZE**($P(\mathbf{W}_*), \mathcal{G}, \mathbb{A}$) does not return FAIL on the outer loops
1383 evaluated for X , if and only if there are no conflicts in the conditions [i-iii] for \mathbf{W}_* w.r.t X . *In other*
1384 *words, all conflicts w.r.t X , in the conditions [i-iii] combined across the terms in \mathbf{W}_* , are identified*
1385 *in the algorithm steps that involve X .*
1386

1387 Next, define the new set \mathbf{W}'_* by dropping x from the subscript (if it appears) for each $W_t \in \mathbf{W}_*$, and
1388 dropping X from \mathbf{W}_* (if it appears). Since X is first in topological order of \mathcal{G} , this does not add any
1389 *new* conflicts across conditions [i-iii] induced by each term in \mathbf{W}'_* (by Lemma C.7).

1390 It is also clear that if there are conflicts *not* involving X , that are induced by conditions [i-iii]
1391 across the terms in \mathbf{W}_* , then these conflicts are also induced by \mathbf{W}'_* . Suppose there are two terms
1392 $W_t, Y_h \in \mathbf{W}_*$ s.t. $\mathbf{T} \setminus X$ needs to be received as $t \setminus X$ by mediating children (condition [i]) for W_t ,
1393 and this conflicts with the requirement that $\mathbf{T} \setminus X$ needs to be received as $t' \setminus X$ (or naturally) by the
1394 same mediating children, for Y_h . Removing X does not affect this conflict, since X is topologically
1395 prior.

1396 Next, define the graph \mathcal{G}' as the projection of \mathcal{G} that marginalizes out X (and adds bidi-
1397 rected edges if needed). \mathcal{G}' has $\leq n$ nodes. Therefore, from the IH, we conclude that **CTF-**
1398 **REALIZE**($P(\mathbf{W}'_*), \mathcal{G}', \mathbb{A}$) does not return FAIL if and only if there are no conflicts induced by
1399 conditions [i-iii], combined across terms in \mathbf{W}'_* .

1400 Now, we note that **CTF-REALIZE**($P(\mathbf{W}_*), \mathcal{G}, \mathbb{A}$) is merely **CTF-REALIZE**($P(\mathbf{W}'_*), \mathcal{G}', \mathbb{A}$), plus
1401 all the steps involving X that we discussed earlier (can be verified from inspecting the algorithm -
1402 the former has an outer loop involving X and then contains the same steps as the latter). *Therefore,*
1403 *all conflicts induced by conditions [i-iii] that involve X and do not involve X are identified in the*
algorithm steps when run on \mathcal{G} and \mathbf{W}_ .*

1404 Thus, we show that $\mathbf{CTF-REALIZE}(P(\mathbf{W}_*), \mathcal{G}, \mathbb{A})$ returns FAIL if and only if conditions [i-iii] are
 1405 not met simultaneously, when combined across all $W_t \in \mathbf{W}_*$ w.r.t a causal diagram having $\leq n + 1$
 1406 nodes. The IH stands proved.

1407 By Lemma C.6, we know that conditions [i-iii] are necessary and sufficient to evaluate each term
 1408 in $\mathbf{W}_*(\mathbf{u})$ simultaneously, for any SCM compatible with \mathcal{G} . By Lemma C.5, we know that this is
 1409 equivalent to drawing an i.i.d sample from $P(\mathbf{W}_*)$. This gives us the proof of the theorem.
 1410

1411 Note: we don't discuss the rejection sampling steps involved steps 17-18 of Algo. 1 as this is trivially
 1412 equivalent to intervening using a fixed value. ■

1413
 1414
 1415

1416 **Corollary 3.7** :

1417 The proof intuition is as follows: given a graph \mathcal{G} and a potential response Y_x , the set of (counterfac-
 1418 tual) ancestors of Y_x (Correa et al., 2021) lists each ancestor of Y and *what regime it must be realized*
 1419 *in*, in order for Y_x to be evaluated. In other words $An(Y_x)$ tracks the regimes necessary and sufficient
 1420 for its ancestors to be evaluated under to beget Y_x .
 1421

1422 For instance, in graph \mathcal{G}_1 in Fig. 3, in order to evaluate W_t , we need A_t to be evaluated in the regime
 1423 \mathcal{M}_t . In order to evaluate Z_x , we need A, T to both be evaluated naturally. This reveals a conflict at
 1424 the bottleneck f_A , which renders the distribution non-realizable.

1425 Thus, Corollary 3.7 provides a *sufficient* condition to conclude that a distribution is non-realizable,
 1426 if $An(\mathbf{W}_*)$ contains two potential responses of the same variable under different regimes. It also
 1427 becomes a *necessary* condition for non-realizability, if the agent can perform $\mathbf{CTF-RAND}(X \rightarrow C)$,
 1428 separately for each $C \in Ch(X)$, for all X . I.e., if the action set is $\mathbb{A}^\dagger(\mathcal{G})$.

1429 The proof steps are similar to Theorem 3.5.

1430 **Inductive Hypothesis (IH):**

1431 Given a graph \mathcal{G} with $\leq n$ nodes, and an arbitrary distribution \mathbf{W}_* , $\mathbf{CTF-}$
 1432 $\mathbf{REALIZE}(P(\mathbf{W}_*), \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G}))$ if and only if $An(\mathbf{W}_*)$ does not contain a pair of potential
 1433 responses W_t, W_s of the same variable W under different regimes.
 1434

1435 **Base case:**

1436 For a graph containing only one variable Y , this is trivially true. $An(Y) = Y$, and the distribution
 1437 $P(Y)$ is realizable.
 1438

1439 Assume IH is true for a graph of $\leq n$ nodes.

1440 **n+1 case:**

1441 Consider an SCM whose causal diagram \mathcal{G} has $n + 1$ nodes. Let X be the first in some topological
 1442 ordering of \mathcal{G} . The agent can perform \mathbb{A}^\dagger in the environment, and the distribution is some arbitrary
 1443 $P(\mathbf{W}_*)$. WLOG, we can begin the outer loop of $\mathbf{CTF-REALIZE}(P(\mathbf{W}_*), \mathcal{G}, \mathbb{A})$ with X (first in
 1444 topological order).
 1445

1446 From Lemmas C.6 and C.5, we know that conditions [i-iii] for each $W_t \in \mathbf{W}_*$, combined across \mathbf{W}_*
 1447 form a necessary and sufficient set to realize $P(\mathbf{W}_*)$.

1448 Note that condition [iii] is always satisfied because the agent need never perform a Fisherian $\mathbf{RAND}(V)$
 1449 for any V . It can get the same effect by performing $\mathbf{CTF-RAND}(V \rightarrow C)$ for each $C \in Ch(V)$. Step
 1450 12 of the sub-routine, Algo. 2 would never be invoked.
 1451

1452 From Theorem 3.5, we know that $\mathbf{CTF-REALIZE}(P(\mathbf{W}_*), \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G}))$ returns FAIL if and only if
 1453 there are conflicts in conditions [i-ii] when combined across all the terms \mathbf{W}_* .

1454 Define the new set \mathbf{W}'_* by dropping x from the subscript (if it appears) for each $W_t \in \mathbf{W}_*$, and
 1455 dropping X from \mathbf{W}_* (if it appears). Since X is first in topological order of \mathcal{G} , this does not add any
 1456 *new* conflicts across conditions [i-ii] induced by each term in \mathbf{W}'_* (by Lemma C.7). It also doesn't
 1457 *remove* any conflicts that are not related to X , as argued in the proof of Theorem 3.5, since X comes
 topologically first.

1458 Define the graph \mathcal{G}' as the projection of \mathcal{G} that marginalizes out X (and adds bidirected edges if
 1459 needed). \mathcal{G}' has $\leq n$ nodes. From the IH, we conclude that $\mathbf{CTF-REALIZE}(P(\mathbf{W}'_\star), \mathcal{G}', \mathbb{A}^\dagger(\mathcal{G}'))$
 1460 does not return FAIL if and only if $An(\mathbf{W}'_\star)$ does not contain two potential responses W_t, W_s of the
 1461 same variable under different regimes.

1462 However, note that (as discussed in the proof of Theorem 3.5, and from inspecting
 1463 the algorithm), the only difference between $\mathbf{CTF-REALIZE}(P(\mathbf{W}_\star), \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G}))$ and $\mathbf{CTF-}$
 1464 $\mathbf{REALIZE}(P(\mathbf{W}'_\star), \mathcal{G}', \mathbb{A}^\dagger(\mathcal{G}'))$ is that in the former, the outer loop of $\mathbf{CTF-REALIZE}$ first checks
 1465 for conflicts in the conditions [i-ii] across \mathbf{W}_\star w.r.t X . After that, the steps for both algorithms are
 1466 the identical.

1467 Therefore, any conflicts detected by $\mathbf{CTF-REALIZE}(P(\mathbf{W}_\star), \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G}))$ that are *not* detected by
 1468 $\mathbf{CTF-REALIZE}(P(\mathbf{W}'_\star), \mathcal{G}', \mathbb{A}^\dagger(\mathcal{G}'))$ must be conflicts w.r.t X . By the IH, these additional conflicts
 1469 (unrelated to X) cannot be because of a pair of conflicting potential responses in $An(\mathbf{W}'_\star)$.

1471 We have already established that removing X to make \mathbf{W}'_\star does not remove or add any conflicting
 1472 potential response pairs that don't involve X . Therefore, our task is to now show that each of
 1473 these additional conflicts (involving X) must correspond to at least one conflicting pair of potential
 1474 responses in $An(\mathbf{W}_\star)$, that are not present in $An(\mathbf{W}'_\star)$. And conversely, we need to show that each
 1475 pair of conflicting potential responses in $An(\mathbf{W}_\star)$ involving X (i.e., that is not present in $An(\mathbf{W}'_\star)$)
 1476 corresponds to at least one conflict detected by $\mathbf{CTF-REALIZE}(P(\mathbf{W}_\star), \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G}))$ in the outer
 1477 loop involving X .

1478 Forward direction:

1479 As discussed in the proof of Theorem 3.5, $\mathbf{CTF-REALIZE}(P(\mathbf{W}_\star), \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G}))$ returns FAIL in the
 1480 outer loop involving X if and only if the input of X to some $C \in Ch(X)$ is required to be some x to
 1481 satisfy condition [i] w.r.t some $W_t \in \mathbf{W}_\star$, but also required to be x' or "natural" to satisfy condition
 1482 [i/ii] w.r.t some $Y_h \in \mathbf{W}_\star$.

1483 Note that the action set is $\mathbb{A}^\dagger(\mathcal{G})$. Therefore step 6 of sub-routine Algo. 2 would always pick only
 1484 the procedure $\mathbf{CTF-RAND}(X \rightarrow C)$ whenever C needs to receive a fixed value. The interventions
 1485 affecting other $C' \in Ch(X)$ would not affect C .

1486 In this case, it is easy to see that the set $An(W_t)$ must contain $C_{x\dots}$ per Def. 3.6, and $An(Y_h)$
 1487 must contain $C_{x'\dots}$ or a potential response of C without X in the subscript. Thus, if $\mathbf{CTF-}$
 1488 $\mathbf{REALIZE}(P(\mathbf{W}_\star), \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G}))$ returns FAIL in the outer loop involving X , there must be a pair of
 1489 conflicting potential responses in $An(\mathbf{W}_\star)$.

1491 Reverse direction:

1492 Assume there exists a conflicting pair of potential responses $A_t, A_s \in An(\mathbf{W}_\star)$ where $x \in t$ and s
 1493 contains some x' or does not contain X at all.

1494 This means there is some $W_h \in \mathbf{W}_\star$ s.t. $A \in An(W)_{\mathcal{G}_X}$ and some $Y_j \in \mathbf{W}_\star$ s.t. $A \in An(Y)_{\mathcal{G}_X}$.
 1495 I.e., A mediates the effect of X on W, Y . Further, from Def. 3.6, it means that A needs to be realized
 1496 in conflicting regimes w.r.t X .

1497 From Lemma C.6, such conflict happens because for A_t , condition [i] requires that X is fixed by
 1498 intervention to be x for all children $C \in Ch(X) \cap An(A)$. Whereas, for A_s , condition [i] or [ii]
 1499 requires that each child $C \in Ch(X) \cap An(A)$ receives X either fixed as x' , or naturally (as the
 1500 case may be, for s). For any such $C \in Ch(X) \cap An(A)$, it is clear from the proof of Theorem 3.5
 1501 that this conflict will trigger a FAIL from $\mathbf{CTF-REALIZE}(P(\mathbf{W}_\star), \mathcal{G}, \mathbb{A}^\dagger(\mathcal{G}))$ in the first outer loop
 1502 involving X .

1504 Thus, we have shown that the IH holds for any \mathbf{W}_\star involving a graph with $n + 1$ nodes.

1505 Since Theorem 3.5 shows $\mathbf{CTF-REALIZE}$ is complete, we have thus proved that \mathbf{W}_\star is realizable
 1506 given \mathcal{G} and $\mathbb{A}^\dagger(\mathcal{G})$ if and only if $An(\mathbf{W}_\star)$ does not contain a pair of conflicting potential responses
 1507 for the same variable under different regimes.
 1508 ■

1509
 1510
 1511

Corollary 3.8 :

1512 This follows from Corollary 3.7. For any causal diagram, the ancestral set of $\{Y_x, Y_{x'}\}$ would include
 1513 both these potential responses.

1514 Thus, the query is not realizable.
 1515

1516 ■

1517
 1518
 1519
 1520

1521 C.3 REALIZABILITY OF \mathcal{L}_1 - AND \mathcal{L}_2 -DISTRIBUTIONS

1522
 1523 It is widely known and acknowledged that it is possible to draw samples from \mathcal{L}_1 - and \mathcal{L}_2 -distributions:
 1524 the former by simply observing a system's natural behaviour, and the latter by intervening in the
 1525 system through interventions like Fisherian randomization.

1526 Still, we find it educational to derive these proofs from first principles. This sub-section is not strictly
 1527 needed to follow the main contributions in Secs. 2 and 3.

1528 We define the probability measure $P^{\mathbb{C}}(\cdot)$ from the perspective of an exogenous agent (i.e., an agent
 1529 external to the system) \mathbb{C} 's beliefs about the environment, distinguished by superscript from $P^{\mathcal{M}}(\cdot)$,
 1530 the true unknown distribution.

1531 Since unit selection is randomized, $\text{SELECT}^{(i)}$ yields an unbiased sample of a unit with latent features
 1532 distributed according to the target population frequency $P(\mathbf{u})$. I.e., $P^{\mathbb{C}}(\mathbf{U}^{(i)} = \mathbf{u} \mid \text{SELECT}^{(i)}) =$
 1533 $P^{\mathcal{M}}(\mathbf{u})$. We also assume that target population size is large enough that $\text{SELECT}^{(i)}$ does not
 1534 significantly change the distribution of the remaining population.

1535 Further, we assume that the actions $\text{READ}(V)^{(i)}$ and $\text{RAND}(V)^{(i)}$ do not disrupt any other mechanism
 1536 $f_{V'}$ for unit i .

1537
 1538
 1539
 1540

1541 **Lemma C.8** (Observational sample). *An agent \mathbb{C} can draw an i.i.d sample distributed according to*
 1542 *the \mathcal{L}_1 query $P(\mathbf{V})$ associated with an SCM \mathcal{M} , by the following actions:*

1543 i. $\text{SELECT}^{(i)}$

1544 ii. $\text{READ}(\mathbf{V})^{(i)} = \mathbf{v} \sim P(\mathbf{V})$

1545 Given N i.i.d samples, the consistent unbiased estimate of $P(\mathbf{v})$ is

1546
 1547
 1548

$$1549 \hat{P}(\mathbf{v}) := \frac{1}{N} \sum_i \prod_{v \in \mathbf{v}} \mathbb{1}[\text{READ}(V)^{(i)} = v] \quad (14)$$

1550
 1551
 1552

1553 *Proof.* This follows directly from the definitions of the actions. $\text{SELECT}^{(i)}$ chooses a unit at random
 1554 from the population. By Remark C.3, $P^{\mathbb{C}}(\mathbf{U}^{(i)} = \mathbf{u} \mid \text{SELECT}^{(i)}) = P^{\mathcal{M}}(\mathbf{u})$. For randomly selected
 1555 unit i ,

$$1556 P^{\mathbb{C}}(\text{READ}(\mathbf{V})^{(i)} = \mathbf{v} \mid \text{SELECT}^{(i)}) \quad (15)$$

$$1557 = \sum_{\mathbf{u}} P^{\mathbb{C}}(\mathbf{U}^{(i)} = \mathbf{u} \mid \text{SELECT}^{(i)}). \quad (16)$$

$$1559 P^{\mathbb{C}}(\text{READ}(\mathbf{V})^{(i)} = \mathbf{v} \mid \mathbf{U}^{(i)} = \mathbf{u}, \text{SELECT}^{(i)}) \quad \text{Chain rule}$$

$$1560 = \sum_{\mathbf{u}} P^{\mathbb{C}}(\mathbf{U}^{(i)} = \mathbf{u} \mid \text{SELECT}^{(i)}) \cdot \mathbb{1}^{\mathcal{M}}[\mathbf{V}(\mathbf{u}) = \mathbf{v}] \quad \text{Def. 2.1(ii)} \quad (17)$$

$$1561 = \sum_{\mathbf{u}} P^{\mathcal{M}}(\mathbf{u}) \cdot \mathbb{1}^{\mathcal{M}}[\mathbf{V}(\mathbf{u}) = \mathbf{v}] \quad \text{Rem. C.3} \quad (18)$$

$$1562 = P^{\mathcal{M}}(\mathbf{v}) \quad \text{Definition} \quad (19)$$

1563
 1564
 1565

1566 I.e., this record is an i.i.d. sample from $P^{\mathcal{M}}(\mathbf{V})$. Now consider the estimator below.

$$1567 \hat{P}(\mathbf{v}) := \frac{1}{N} \sum_n \prod_{v \in \mathbf{v}} \mathbb{1}^{\mathbb{C}}[\text{READ}(V)^{(i)} = v] \quad (20)$$

$$1570 = \frac{1}{N} \sum_n \sum_{\mathbf{u}} \prod_{v \in \mathbf{v}} \mathbb{1}^{\mathcal{M}}[\mathbf{U}^{(i)} = \mathbf{u}] \cdot \mathbb{1}^{\mathcal{M}}[V(\mathbf{u}) = v] \quad (21)$$

1572 Un-biasedness is established by taking expectation on either side, w.r.t the agent \mathbb{C} 's actions (choice
1573 of units to observe):

$$1575 \mathbb{E}_{\mathbb{C}}[\hat{P}(\mathbf{v})] = \mathbb{E}_{\mathbb{C}} \left[\frac{1}{N} \sum_n \sum_{\mathbf{u}} \prod_{v \in \mathbf{v}} \mathbb{1}^{\mathcal{M}}[\mathbf{U}^{(i)} = \mathbf{u}] \cdot \mathbb{1}^{\mathcal{M}}[V(\mathbf{u}) = v] \right] \quad (22)$$

$$1578 = \sum_{\mathbf{u}} \frac{1}{N} \mathbb{E}_{\mathbb{C}} \left[\sum_n \mathbb{1}^{\mathcal{M}}[\mathbf{U}^{(i)} = \mathbf{u}] \prod_{v \in \mathbf{v}} \mathbb{1}^{\mathcal{M}}[V(\mathbf{u}) = v] \right] \quad \text{Linearity of expectation} \quad (23)$$

$$1581 = \sum_{\mathbf{u}} \frac{1}{N} \mathbb{E}_{\mathbb{C}} \left[\sum_n \mathbb{1}^{\mathcal{M}}[\mathbf{U}^{(i)} = \mathbf{u}] \right] \prod_{v \in \mathbf{v}} \mathbb{1}^{\mathcal{M}}[V(\mathbf{u}) = v] \quad V(\mathbf{u}) \text{ constant wrt } \mathbb{C} \quad (24)$$

$$1583 = \sum_{\mathbf{u}} \frac{1}{N} \left[N \cdot P^{\mathcal{M}}(\mathbf{u}) \right] \prod_{v \in \mathbf{v}} \mathbb{1}^{\mathcal{M}}[V(\mathbf{u}) = v] \quad \text{Def. 2.1(i), Rem. C.3} \quad (25)$$

$$1586 = P^{\mathcal{M}}(\mathbf{v}) \quad \text{Definition} \quad (26)$$

1587 Consistency is established by the fact that as \mathcal{N} (target population size) $\rightarrow \infty$, and N (sample size)
1588 $\rightarrow \infty$,

$$1590 \frac{1}{N} \sum_n \mathbb{1}^{\mathcal{M}}[\mathbf{U}^{(i)} = \mathbf{u}] \rightarrow P^{\mathcal{M}}(\mathbf{u}) \quad (27)$$

1594 **Lemma C.9.** *The \mathcal{L}_2 distribution of an atomic intervention is equivalent to the \mathcal{L}_2 distribution of the
1595 corresponding conditional stochastic intervention.*

$$1597 P^{\mathcal{M}}(\mathbf{v}; do(\mathbf{x})) = P^{\mathcal{M}}(\mathbf{v} | \mathbf{x}; \sigma_{\mathbf{X}}) \quad (28)$$

$$1598 = \sum_{\mathbf{u}} \underbrace{\mathbb{1}[\mathbf{V}_{\sigma_{\mathbf{X}}}(\mathbf{u}) = \mathbf{v} | X_{\sigma_{\mathbf{X}}} = \mathbf{x}]}_{\textcircled{1}} \cdot \underbrace{P(\mathbf{u})}_{\textcircled{2}} \quad (29)$$

1601 *Proof.* The step from the r.h.s of Eq. 28 to Eq. 29 is derived as follows: in the submodel $\mathcal{M}_{\sigma_{\mathbf{X}}}$, if we
1602 are given that \mathbf{X} has been randomly assigned \mathbf{x} , then the remaining variables are deterministically
1603 generated as a function of \mathbf{u} and \mathbf{x} via their respective equations. The probability mass is collected
1604 over all the \mathbf{u} which produce the output \mathbf{v} over all these equations.

$$1606 P^{\mathcal{M}}(\mathbf{v} | \mathbf{x}; \sigma_{\mathbf{X}}) = \sum_{\mathbf{u}} \mathbb{1}[\mathbf{V}_{\sigma_{\mathbf{X}}}(\mathbf{u}) = \mathbf{v} | X_{\sigma_{\mathbf{X}}} = \mathbf{x}] \cdot P^{\mathcal{M}}(\mathbf{u}) \quad (30)$$

1608 Notice: if \mathbf{v} is incompatible with \mathbf{x} , the indicator in the r.h.s evaluates to 0. Next, we prove. Eq. 28.

1609 In $\mathcal{M}_{\sigma_{\mathbf{X}}}$, as defined, \mathbf{X} is assigned according to an independent random vector. Notate this vector as
1610 $\mathbf{X}_{\sigma_{\mathbf{X}}}$ and let the distribution of this vector be $P_{\sigma_{\mathbf{X}}}(\mathbf{X})$, defined by the assignment frequency over the
1611 target population.

1612 $\mathcal{M}_{\sigma_{\mathbf{X}}}$ is defined such that the target population is split into groups, each assigned ($X_{\sigma_{\mathbf{X}}} = \mathbf{x}$) for
1613 some \mathbf{x} . Note, the assignment vector $\mathbf{X}_{\sigma_{\mathbf{X}}}$ is independent of the latent features \mathbf{U} across the target
1614 population iff each finite group assigned ($X_{\sigma_{\mathbf{X}}} = \mathbf{x}$) has the same distribution of latent features
1615 $P(\mathbf{U})$ as in the overall target population.

1617 The above discussion handles the finite size of the target population. Starting with the r.h.s of Eq. 28,

$$1618 P^{\mathcal{M}}(\mathbf{v} | \mathbf{x}; \sigma_{\mathbf{X}}) = \frac{P(\mathbf{v}, \mathbf{x}; \sigma_{\mathbf{X}})}{P(\mathbf{x}; \sigma_{\mathbf{X}})} = \begin{cases} P(\mathbf{v}; \sigma_{\mathbf{X}}) / P(\mathbf{x}; \sigma_{\mathbf{X}}) & \text{if } \mathbf{v} \text{ compatible with } \mathbf{x} \\ 0 & \text{otherwise} \end{cases} \quad (31)$$

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Evaluating for when \mathbf{v} is compatible with \mathbf{x} :

$$\frac{P(\mathbf{v}; \sigma_{\mathbf{X}})}{P(\mathbf{x}; \sigma_{\mathbf{X}})} = \frac{P(\mathbf{v}; \sigma_{\mathbf{X}})}{P_{\sigma_{\mathbf{X}}}(\mathbf{x})} \quad (32)$$

$$= \frac{\sum_{\mathbf{u}} \left(P(\mathbf{u}) \prod_{V_i \in \mathbf{V} \setminus \mathbf{X}} P(v_i | \mathbf{pa}_i, \mathbf{u}_i) \cdot P_{\sigma_{\mathbf{X}}}(\mathbf{x}) \right)}{P_{\sigma_{\mathbf{X}}}(\mathbf{x})} \quad \text{Truncated factorization product} \quad (33)$$

$$= \sum_{\mathbf{u}} P(\mathbf{u}) \prod_{V_i \in \mathbf{V} \setminus \mathbf{X}} P(v_i | \mathbf{pa}_i, \mathbf{u}_i) \quad (34)$$

$$= P^{\mathcal{M}}(\mathbf{v}; do(\mathbf{x})) \quad \text{Truncated factorization product} \quad (35)$$

Eq. 33 uses the fact that each sub-group assigned ($X_{\sigma_{\mathbf{X}}} = \mathbf{x}$), by independence, has the same frequency of latent features $P(\mathbf{u})$. ■

Lemma C.10 (Interventional sample). *An agent \mathbb{C} can draw an i.i.d sample distributed according to the \mathcal{L}_2 query $P(\mathbf{V}; do(\mathbf{x}))$ associated with an SCM \mathcal{M} , by the following actions:*

- i. SELECT⁽ⁱ⁾
- ii. RAND(\mathbf{X})⁽ⁱ⁾
- iii. If RAND(\mathbf{X})⁽ⁱ⁾ = \mathbf{x} , then READ(\mathbf{V})⁽ⁱ⁾ = $\mathbf{v} \sim P(\mathbf{V}; do(\mathbf{x}))$, else repeat i-iii.

Given $N_{\mathbf{x}}$ i.i.d samples, the consistent unbiased estimate of Eq. 29 is given by

$$\hat{P}(\mathbf{v}; do(\mathbf{x})) = \underbrace{\frac{1}{N_{\mathbf{x}}}}_{\textcircled{2}} \sum_i \underbrace{\mathbb{1}[\text{READ}(\mathbf{V})^{(i)} = \mathbf{v}, \text{RAND}(\mathbf{X})^{(i)} = \mathbf{x}]}_{\textcircled{1}}, \quad (36)$$

Proof. The proof steps are similar to the ones used for the Observational i.i.d sample case. Note that Remark C.3 still hold since even though the agent is conditioning on the value randomly assigned to a particular unit i , this value is independent of the unit's latent features $\mathbf{U}^{(i)}$. ■

1674 D DETAILS ON COUNTERFACTUAL RANDOMIZATION

1675
1676 In this appendix, we provide a formal account for the procedure we define in Sec. 2, called
1677 CTF-RAND($X \rightarrow \mathbf{C}$) (Def. 2.3). This appendix is intended for the interested reader, and de-
1678 tails some useful remarks to be used in the proofs of our results in Appendix C. These details are not
1679 strictly needed to follow the results presented in the main body of the paper.

- 1680 • In Sec. D.1, we lay out the structural conditions under which it is possible to perform this
1681 procedure, and we provide an algorithm (Algorithm 3) by which an agent can translate the
1682 structural conditions in the environment into a list of CTF-RAND procedures that it is able
1683 to perform in the given setting.
- 1684 • In Sec. D.2, we emphasize that it is possible for an agent to enact multiple randomization
1685 procedures involving the same variable X for a single unit i , and illustrate this with an
1686 example.
- 1687 • In Sec. D.3, we discuss the constraints implied by the assumptions we make. In particular,
1688 we discuss why CTF-RAND($X \rightarrow \mathbf{C}$) can only be performed on some $\mathbf{C} \subseteq Ch(X)$, and
1689 not by-pass children to directly affect some distant descendants of X .

1690 D.1 STRUCTURAL CONDITIONS REQUIRED FOR COUNTERFACTUAL RANDOMIZATION

1691 Counterfactual randomization (Def. 2.3) can be performed under two circumstances:

- 1692 i. CTF-RAND($X \rightarrow Ch(X)$) can be performed by eliciting a unit’s natural decision X , while
1693 simultaneously randomizing its actual enforced decision. Thus, the agent can affect the
1694 value of the decision X as received by all the children of X , whilst also recording the natural
1695 realization of X . As discussed in Sec. 2, this was established in (Bareinboim et al., 2015;
1696 Forney et al., 2017; Zhang & Bareinboim, 2022).
- 1697 ii. CTF-RAND($X \rightarrow \mathbf{C}$) can also be performed for some $\mathbf{C} \subseteq Ch(X)$ if there is a special
1698 *counterfactual mediator* (defined below) by which the mechanisms generating \mathbf{C} perceive
1699 the value of X . This counterfactual mediator then allows the agent to intervene on the value
1700 of X as *perceived* by \mathbf{C} , thus mimicking an actual intervention on X .

1701 **Definition D.1** (Expanded SCM). Given an SCM \mathcal{M} containing observable variables \mathbf{V} , we define
1702 an *expanded SCM* \mathcal{M}^+ of the same environment to be a model containing a bigger set of observable
1703 variables $\mathbf{V}^+ \supset \mathbf{V}$, and which relaxes the positivity requirement. I.e., it is possible that $P^{\mathcal{M}^+}(\mathbf{v}^+) =$
1704 0, for some \mathbf{v}^+ in \mathcal{L}_1 . We call the causal diagram of \mathcal{M}^+ an *expanded causal diagram* \mathcal{G}^+ . ■

1705 **Definition D.2** (Counterfactual mediator (formal)). Given a variable X in a causal diagram \mathcal{G} , we
1706 call any variable $W \notin \mathbf{V}$ a *counterfactual mediator* of X w.r.t $Y \in Ch(X)_{\mathcal{G}}$ if

- 1707 i. In an “expanded” SCM of the environment \mathcal{M}^+ (Def. D.1), W is generated according to an
1708 invertible mechanism $W \leftarrow f_W(X, \mathbf{U}_W)$ with \mathbf{U}_W possibly empty, s.t. $f_W^{-1}(W) = X$;
- 1709 ii. It is physically possible to perform RAND(W)⁽ⁱ⁾ (Def. 2.1); and
- 1710 iii. In \mathcal{M}^+ , Y is generated by the mechanism $Y \leftarrow f_Y(f_W^{-1}(W), \mathbf{A}, \mathbf{U}_Y)$, where \mathbf{A} is the set
1711 $\mathbf{Pa}_Y \setminus X$ in \mathcal{G} . ■

1712 The intuition behind Def. D.2 is that a *counterfactual mediator* is a real variable in the environment
1713 which fully encodes information about the variable X , and which mediates how Y perceives the
1714 value of X via the “direct” causal path. For instance, in Example 1 (*Mediation analysis*), the RGB
1715 values of the video frames W are a counterfactual mediator for the mechanism f_Y (decision to issue
1716 a speeding ticket) to perceive the car’s color X via the “direct” path, not via the actual speeding of
1717 the car).

1718 **Condition [i]** of Def. D.2 divides the domain of W into *equivalence classes* s.t. each value w belongs
1719 to an equivalence class $\{w' : f_W^{-1}(w) = x\}$ for some value x .

1720 **Condition [iii]** of Def. D.2 essentially says that the mechanism f_Y only cares about which equivalence
1721 class W belongs to. I.e., Y only cares about what W reveals about X .

Note: these conditions does not require an agent to have full knowledge of the SCM. They are rather structural assumptions about the underlying mechanisms which can be verified in a given setting. In Example 1, treating W as counterfactual mediator means the assumption that (1) the video features W uniquely map back to the actual color of the car in the footage; and (2) the computer vision system only cares W reveals about X , and is indifferent to any stochasticity *within* some equivalence class $\{w' : f_W^{-1}(w) = x\}$.

Assumption D.3. (Tree structure) Given a variable X , causal diagram \mathcal{G} , and an "expanded" diagram \mathcal{G}^+ (Def. D.1) including the set of all the counterfactual mediators \mathbf{W} (Def. D.2) of X in the environment, each $W \in \mathbf{W}$ has only one parent in \mathcal{G}^+ , and each $C \in Ch(X)_{\mathcal{G}}$ has at most one $W \in \mathbf{W}$ as a parent in \mathcal{G}^+ . ■

Assumption D.3 enforces that each child of X perceives X through at most one proxy pathway. This assumption rules out possible structures like Fig. 15(a) where a child perceives X through multiple proxy pathways.

This assumption is general enough to allow most cases of interest. If X is a construct like gender identity, then it is possible that a child perceives X via a cluster of personal attributes \mathbf{W} which indicate X . In this case, no single attribute solely satisfies Def. D.2 of a counterfactual mediator. However, the cluster of attributes \mathbf{W} could be collapsed into a single variable having domain equal to the cartesian product of the sub-domains (Anand et al., 2023; Xia & Bareinboim, 2024). This single node \mathbf{W} would indeed satisfy the definition of a counterfactual mediator and would comply with the tree structure in Assumption D.3. For a comprehensive discussion of the semantics of interventions on the perception of a compound attribute such as race or gender identity, see (Plecko & Bareinboim, 2024, App. D.1). The following Lemma is the key property that enables path-specific randomization.

Lemma D.4. *Given a causal diagram \mathcal{G} containing variables X and $Y \in Ch(X)_{\mathcal{G}}$. Let W be a counterfactual mediator of X w.r.t Y (Def. D.2). For any value x , we have*

$$Y_{w\mathbf{a}}(\mathbf{u}) = Y_{x\mathbf{a}}(\mathbf{u}), \quad \forall \mathbf{u}, \forall w \in \{w' : f_W^{-1}(w) = x\}, \quad (37)$$

where $\mathbf{A} := \text{Pa}_Y \setminus X$ in \mathcal{G} .

Proof. This follows from Def. D.2. Suppose we are given values (w, x) where $f_W^{-1}(w) = x$. Let $\mathbf{A} := \text{Pa}_Y \setminus X$ in \mathcal{G} .

The variable $W_x(\mathbf{u}) = W_{x\mathbf{a}}(\mathbf{u}) = f_W(x, \mathbf{u})$, in the enhanced submodel $\mathcal{M}_{x\mathbf{a}}^+$. Adding \mathbf{a} to the subscript does not matter - by Assumption D.3 and Lemma D.7, \mathbf{A} cannot be an ancestor of W in \mathcal{M}^+ .

Since f_W is invertible by condition [i] in \mathcal{M}^+ , it is also invertible in submodel submodel $\mathcal{M}_{x\mathbf{a}}^+$. Therefore, we have $f_W^{-1}(W_{x\mathbf{a}}(\mathbf{u})) = x$.

$$\begin{aligned} Y_{w\mathbf{a}}(\mathbf{u}) &= f_Y(f_W^{-1}(w), a, \mathbf{u}) \\ &= f_Y(x, a, \mathbf{u}) \\ Y_{W_{x\mathbf{a}}\mathbf{a}}(\mathbf{u}) &= f_Y(f_W^{-1}(W_{x\mathbf{a}}(\mathbf{u})), a, \mathbf{u}) \\ &= f_Y(x, a, \mathbf{u}) \end{aligned}$$

The r.h.s is identical, giving us $Y_{w\mathbf{a}} = Y_{W_{x\mathbf{a}}\mathbf{a}}$. Finally, we argue that $Y_{W_{x\mathbf{a}}\mathbf{a}} = Y_{x\mathbf{a}}$.

The counterfactual $Y_{x\mathbf{a}}$ is evaluated in a submodel of \mathcal{M}^+ , where f_W receives input x and this value of W_x is an input to f_Y , while \mathbf{A} is fixed to be \mathbf{a} . Structurally, this is identical to how the counterfactual $Y_{W_{x\mathbf{a}}\mathbf{a}} = Y_{W_{x\mathbf{a}}}$ is evaluated. Therefore, it is evident that $Y_{W_{x\mathbf{a}}\mathbf{a}} = Y_{x\mathbf{a}}$. ■

Given a variable X , the way an agent actually performs the action CTF-RAND is as follows:

- i. **Performing CTF-RAND by eliciting natural decision:** The agent can perform CTF-RAND($X \rightarrow Ch(X)$)⁽ⁱ⁾ by randomizing the unit's decision. The agent can further perform READ(X)⁽ⁱ⁾ to elicit the unit's natural decision, which has not been erased. This is described in in (Bareinboim et al., 2015; Forney et al., 2017; Zhang & Bareinboim, 2022).

1782
1783
1784
1785
1786
1787
1788
1789
1790

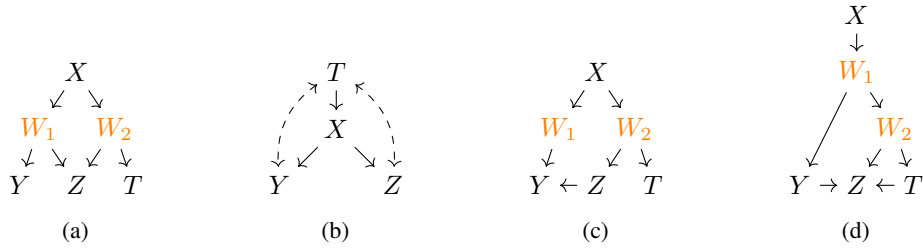


Figure 15: "Expanded" causal diagrams, where counterfactual mediators (labeled W_i) of X have been marked in red. (a) is not permitted by Assumption D.3. Assume the environment in diagram (b) permits the agent to elicit the unit's natural decision X even when randomizing the actual decision.

1791
1792
1793
1794
1795
1796
1797
1798
1799

- ii. **Performing CTF-RAND using counterfactual mediators:** If X has a counterfactual mediator W in the environment, and $C \subseteq Ch(X)$ are the children which perceive X via W , then the agent can perform $CTF-RAND(X \rightarrow C)^{(i)}$ by randomizing W . Each value w mimics randomizing X as perceived by C , per Lemma D.4. The agent can still perform $READ(X)^{(i)}$ by measuring $X^{(i)}$ to get the unit's natural decision, which has not been erased.

1800
1801
1802
1803
1804

Having described the structural conditions that permit counterfactual randomization, we want to abstract away the mediators and succinctly describe the agent's physical actions via the definition of CTF-RAND. Given a variable X , and assumptions/knowledge about X in the environment stated in points [i] and [ii] above, we translate this knowledge into a set of counterfactual randomizations that the agent is physically able to perform in the environment, using Algorithm 3.

1805
1806

Algorithm 3 CTF-PROCEDURES

1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821

- 1: **Input:** Causal diagram \mathcal{G} with decision variable X ; "expanded" diagram \mathcal{G}^+ (Def. D.1) including the counterfactual mediators of X in the environment
 - 2: **Output:** \mathbb{A}_X - the set of CTF-RAND actions that can be performed involving X
 - 3: $\mathbb{A}_X \leftarrow \emptyset$
 - 4: **if** environment allows eliciting natural decision X even when randomizing actual decision **then**
 - 5: **if** X can be randomized **then**
 - 6: $\mathbb{A}_X \leftarrow \mathbb{A}_X \cup \{CTF-RAND(X \rightarrow Ch(X)_{\mathcal{G}})\}$
 - 7: **end if**
 - 8: **end if**
 - 9: **for** each counterfactual mediator W of X **do**
 - 10: Let $C := \{C \mid C \in Ch(X)_{\mathcal{G}} \text{ and perceives } X \text{ via } W\}$
 - 11: $\mathbb{A}_X \leftarrow \mathbb{A}_X \cup \{CTF-RAND(X \rightarrow C)\}$
 - 12: **end for**
 - 13: **Return** \mathbb{A}_X
-

1822
1823
1824
1825
1826

Consider Fig. 15(a-d). (a) is not permitted by Assumption D.3. We assume that in the environment represented by (b) X can be randomized for a unit in the target population without erasing the unit's natural decision, satisfying condition [i] mentioned earlier. Thus, when applying Algorithm 3 to diagrams (b-d), we get the following resulting set of counterfactual randomization procedures which are permitted by the structural assumptions made (unit superscript i is omitted for legibility):

1827
1828
1829
1830
1831

- (b) $\{CTF-RAND(X \rightarrow \{Y, Z\})\}$
- (c) $\{CTF-RAND(X \rightarrow Y), CTF-RAND(X \rightarrow \{Z, T\})\}$
- (d) $\{CTF-RAND(X \rightarrow \{Y, Z, T\}), CTF-RAND(X \rightarrow \{Z, T\})\}$

1832
1833

D.2 MULTIPLE SIMULTANEOUS RANDOMIZATIONS ARE POSSIBLE, FOR A SINGLE UNIT

1834
1835

For a particular decision variable X , there could be multiple randomization procedures which an agent can perform. Consider the example in Fig. 16. The "expanded" diagram on the left shows two counterfactual mediators, W_1, W_2 in a causal structure which permit an agent to perform all

1836 of the following randomization procedures: $\text{RAND}(X)^{(i)}$, $\text{CTF-RAND}(X \rightarrow \{Z, T, B\})^{(i)}$ and
 1837 $\text{CTF-RAND}(X \rightarrow \{T, B\})^{(i)}$ for the same unit i .
 1838

1839 However, if all three actions are performed in parallel, randomizing W_1 to enact $\text{CTF-RAND}(X \rightarrow$
 1840 $\{Z, T, B\})^{(i)}$ will only affect variable Z . This is since the action of randomizing W_2 to further
 1841 enact $\text{CTF-RAND}(X \rightarrow \{T, B\})^{(i)}$ blocks any effect on T, B from the previous action. Similarly,
 1842 $\text{RAND}(X)^{(i)}$ ends up affecting only variable Y , because $\text{CTF-RAND}(X \rightarrow \{Z, T, B\})^{(i)}$ blocks any
 1843 effect from the previous action on Z, T, B . We formalize this observation in Remark D.5.

1844 *Remark D.5* (Superseding action). Given a decision variable X , the action $\text{CTF-RAND}(X \rightarrow \mathbf{C}')^{(i)}$
 1845 can *supersede* the action $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$ if $\mathbf{C}' \subsetneq \mathbf{C}$, where *supersede* means that the
 1846 former action $\text{CTF-RAND}(X \rightarrow \mathbf{C}')^{(i)}$ blocks any effect that the latter action has on the variables \mathbf{C}' .
 1847 Additionally, the action $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$ *supersedes* the action $\text{RAND}(X)^{(i)}$. ■
 1848

1849 Further, Assumption D.3 ensures that all such
 1850 procedures follow a "nested" structure. I.e.,
 1851 given any two randomization procedures involv-
 1852 ing the same variable, the sets of children affec-
 1853 ted by one will be a subset of the set affected
 1854 by the other, as shown in Fig. 16.

1855 *Remark D.6.* (Procedure containment) Assump-
 1856 tion D.3 implies that if an agent is capable of
 1857 performing both $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$ and
 1858 $\text{CTF-RAND}(X \rightarrow \mathbf{C}')^{(i)}$ s.t. $\mathbf{C} \neq \mathbf{C}'$ and
 1859 $\mathbf{C} \cap \mathbf{C}' \neq \emptyset$, then either $\mathbf{C} \subseteq \mathbf{C}'$ or $\mathbf{C}' \subseteq \mathbf{C}$. ■
 1860

1861 D.3 COUNTERFACTUAL 1862 RANDOMIZATION IS ONLY 1863 POSSIBLE FOR DIRECT CHILDREN OF X

1864 Our definition of $\text{CTF-RAND}(X \rightarrow \mathbf{C})^{(i)}$, is
 1865 only valid for some $\mathbf{C} \subseteq \text{Ch}(X)$ in the causal
 1866 diagram (Def. 2.3). This action essentially randomizes the value of decision variable X as an input
 1867 to the mechanisms generating its causal children \mathbf{C} , while leaving open the possibility of measuring
 1868 the unit i 's natural decision (what it would have normally decided in the \mathcal{L}_1 regime), and also the
 1869 possibility of separately and in parallel randomizing the value of X as an input to other causal
 1870 children $\mathbf{C}' = \text{Ch}(X) \setminus \mathbf{C}$.
 1871

1872 However, the notion of "child" is an abstraction w.r.t a specific diagram of the environment under
 1873 study. Consider Fig. 17(a-Left), where \mathcal{G}_1 is the diagram of some environment. Assume there exists
 1874 a counterfactual mediator W_1 of X (Def. D.2) as shown in Fig. 17(a-Middle), which means an agent
 1875 is able to perform the physical action $\text{CTF-RAND}(X \rightarrow A)^{(i)}$, while still being able to measure the
 1876 natural value of X for unit i .

1877 Now consider the diagram \mathcal{G}_2 shown in Fig. 17(b-Left). \mathcal{G}_2 is a valid projection of \mathcal{G}_1 obtained by
 1878 marginalizing out variable A , and is thus also a valid causal diagram of the environment.

1879 Suppose that there exists a counterfactual mediator W_2 as shown in 17(b-Middle). This means that the
 1880 agent can also perform $\text{CTF-RAND}(X \rightarrow Z)^{(i)}$ in the same environment. However, since we are re-
 1881 ferring to the same environment, this means that the agent is able to perform $\text{CTF-RAND}(X \rightarrow Z)^{(i)}$
 1882 w.r.t the diagram \mathcal{G}_1 , where Z is not a child node of X ! This would translate to even greater experimen-
 1883 tal power w.r.t graph \mathcal{G}_1 , where an agent is able to perform counterfactual randomization of X w.r.t
 1884 further descendants like Z and draw i.i.d samples from queries like $P(A_x, Z_{x'})$ by simultaneously
 1885 performing both counterfactual randomizations (i.e. by randomizing W_1, W_2 simultaneously).
 1886

1887 However, this scenario is not possible. Essentially, this would require an "expanded" causal diagram
 1888 (Def. D.1) like shown in Fig. 18, where W_2 is a counterfactual mediator of X w.r.t Z that comes
 1889 after another variable A . If A satisfies positivity w.r.t X , i.e., if $P^{\mathcal{M}}(x, a) > 0, \forall x, a$ in \mathcal{L}_1 , then W_2
 cannot be a counterfactual mediator since it cannot be uniquely mapped back to X .

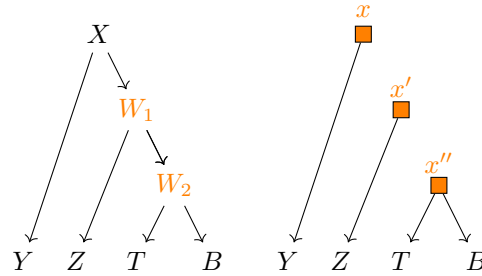


Figure 16: (Left) "Expanded" causal diagram showing counterfactual mediators W_1, W_2 of X ; (Right) Agent performing actions $\text{RAND}(X)^{(i)}$, $\text{CTF-RAND}(X \rightarrow \{Z, T, B\})^{(i)}$ and $\text{CTF-RAND}(X \rightarrow \{T, B\})^{(i)}$ all together on the single unit i .

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

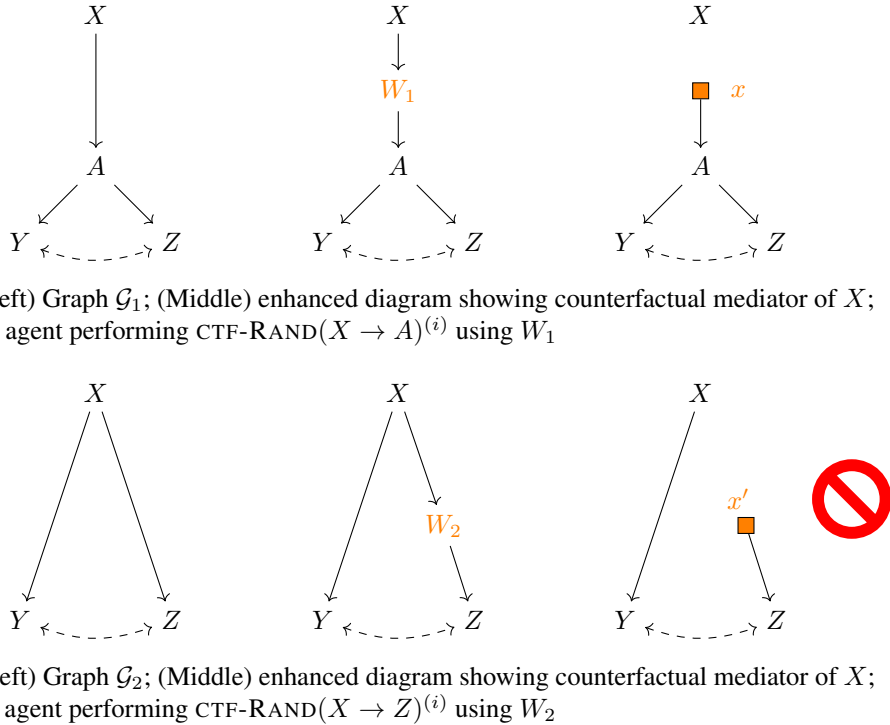


Figure 17: Given a causal diagram \mathcal{G}_1 of the environment, \mathcal{G}_2 is a valid projection of \mathcal{G}_1 (marginalizing A). If A satisfies positivity w.r.t X , then there *cannot* be a counterfactual mediator W_2 as shown in (b-Middle). Which means an agent *cannot* perform $\text{CTF-RAND}(X \rightarrow Z)^{(i)}$ as shown in (b-Right).

Lemma D.7. Given a causal diagram \mathcal{G} of a true SCM \mathcal{M} with a variable X and $A \in \text{Desc}(X)_{\mathcal{G}}$ where $P(x, a) > 0, \forall x, a$. There cannot be a variable W in an "expanded" SCM \mathcal{M}^+ of the environment (Def. D.1) s.t.

- $W \in \text{Desc}(A)_{\mathcal{G}^+}$, where \mathcal{G}^+ is the "expanded" causal diagram of \mathcal{M}^+ ; and
- W is invertible to X , i.e. exists f_W^{-1} s.t. $f_W^{-1}(W) = X$. □

Proof. If A satisfies positivity w.r.t X , then a given value w cannot be mapped back to a unique x , even if we marginalize out A from the SCM.

Note that, by Assumption D.3, a counterfactual mediator has only one parent in the "expanded" causal diagram (Def. D.1). I.e., if it were a descendant of A , its perception of X is fully mediated by A .

If $f'_W(X, \mathbf{U})$ is invertible from W to X , then so is $f_W \circ f_A(X, \mathbf{U})$. It is evident that f'^{-1}_W is well defined iff $f^{-1}_A \circ f^{-1}_W$ is well defined.

f^{-1}_A is not defined. The positivity condition entails that a given value a could have been generated by any value x (when unit is unknown).

Since f'_W is not invertible, W cannot be a counterfactual mediator. ■

Lemma D.7 leads to some important conclusions.

Remark D.8. There cannot be an "expanded" causal diagram (such as in Fig. 18), with a counterfactual mediator that bypasses a child-node and directly fixes a descendent-node's perception of X . I.e., an agent cannot perform $\text{CTF-RAND}(X \rightarrow D)^{(i)}$ for some $D \in \text{Desc}(X) \setminus \text{Ch}(X) \cup \{X\}$. ■

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

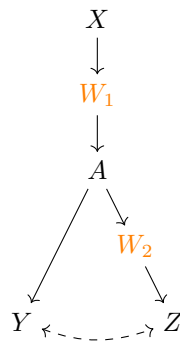


Figure 18: "Expanded" diagram that is needed to sample directly from $P(A_x, Z_{x'})$. This is not possible, per Lemma D.7.

Remark D.9. Conversely, given a graph like \mathcal{G}_2 in Fig. 17(b), if we are told that the agent can perform the action $\text{CTF-RAND}(X \rightarrow Z)^{(i)}$, then \mathcal{G}_2 cannot be a projection of \mathcal{G}_1 (Fig. 17(a)) for the same environment. ■

The upshot of this discussion is that, in general (i.e. without making further assumptions), counterfactual randomization can only be done via counterfactual mediators (Def. D.2) of a decision variable X , and it can only be performed on the children-nodes of X in the general case.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

E DETAILS ON EXAMPLES

E.1 EXAMPLE 1 (MEDIATION ANALYSIS) - FURTHER DISCUSSION

For the interested reader who wants to track the formal definition of a counterfactual mediator (Def. D.2), we provide below a discussion of how these structural assumptions may be verified in a field experiment setting.

This discussion is not strictly needed for following the main contributions in Secs. 2 and 3.

Verification of structural assumptions: In order for W to be a counterfactual mediator of X w.r.t Y , it needs to satisfy 3 conditions in Def. D.2.

Condition [i] says that each w (say, a specific RGB range of pixels) belongs to an equivalence class that maps back uniquely to a car color x . This can be verified using the RCT data. **Condition [ii]** is satisfied since they can do a targeted randomization of W . **Condition [iii]** stipulates that f_Y is not affected by any artefacts introduced by the color-editing tool: this can be verified, for instance, by swapping a car's color from x to x' and then swap it back from x' to x , to ensure that the model's decision Y does not change, thus verifying that the mechanism f_Y only cares about what the color features W reveal about X , and not about any image artefacts that may be introduced by editing.

E.2 EXAMPLE 2 (CAUSAL FAIRNESS)

E.2.1 SIMULATIONS

The details for the simulations shown in Fig. 5(c) are as follows. We first parameterize the space of SCMs that are compatible with the causal diagram in Fig. 5(a) using *canonical parameters* (Zhang et al., 2022, Def. 1, Thm. 1).

These parameters essentially discretize the domain of the latent confounder between Y and Z where each value of the confounder represents a joint mapping in $(X \rightarrow Y) \times (X \rightarrow Z)$. I.e., $2^2 \times 2^2 = 16$ values. We then set up a constrained optimization program to draw samples from the space of all SCMs (i.e., from the space of all trained models in Example 2) which satisfy either

- \mathcal{L}_2 fairness measure $\mu_{int1} + \mu_{int2}$ being penalized; or
- \mathcal{L}_3 fairness measure μ_{ctf} being penalized

Drawing 1000 samples from each gives us the chart in Fig. 5(c).

E.2.2 SCM SPECIFICATION

Just to give intuition for how this disparity can arise, we present below an instantiation of an SCM using canonical parameters, showing how \mathcal{L}_2 measures misleadingly suggest no discrimination, whereas the \mathcal{L}_3 measure actually detects unfairness.

For simplicity, we assume X is binary (0 indicates Race A, 1 indicates Race B). Y, Z are binary outcomes indicating, respectively, the CV passing through the 1st stage of screening for college admission, and for receiving a financial scholarship.

U_X is the random assignment of applicant race in the CV, in the absence of intervention.

The mechanisms f_Y and f_Z represent the decisions of the trained classifiers used by the college data science team, which have been trained using data from previous years' decisions of committees for CV screening and financial aid.

2052
 2053
 2054
 2055
 2056
 2057
 2058
 2059
 2060
 2061
 2062
 2063
 2064
 2065
 2066
 2067
 2068
 2069
 2070
 2071
 2072
 2073
 2074
 2075
 2076
 2077
 2078
 2079
 2080
 2081
 2082
 2083
 2084
 2085
 2086
 2087
 2088
 2089
 2090
 2091
 2092
 2093
 2094
 2095
 2096
 2097
 2098
 2099
 2100
 2101
 2102
 2103
 2104
 2105

SCM \mathcal{M}^*

$$X \leftarrow U_X \sim \text{Bernoulli}(0.5)$$

$$Y \leftarrow \begin{cases} 1, & \text{if Y-type = always-approve} \\ X, & \text{if Y-type = approve iff } x_1 \\ 1 - X, & \text{if Y-type = approve iff } x_0 \\ 0, & \text{if Y-type = always-reject} \end{cases}$$

$$Z \leftarrow \begin{cases} 1, & \text{if Z-type = always-approve} \\ X, & \text{if Z-type = approve iff } x_1 \\ 1 - X, & \text{if Z-type = approve iff } x_0 \\ 0, & \text{if Z-type = always-reject} \end{cases}$$

Y-type	Z-type	$P(U_{YZ})$
Always-approve	Always-approve	0.040
Always-approve	Approve iff x_1	0.175
Always-approve	Approve iff x_0	0.160
Always-approve	Always-reject	0.010
Approve iff x_1	Always-approve	0.040
Approve iff x_1	Approve iff x_1	0.055
Approve iff x_1	Approve iff x_0	0.170
Approve iff x_1	Always-reject	0.010
Approve iff x_0	Always-approve	0.040
Approve iff x_0	Approve iff x_1	0.140
Approve iff x_0	Approve iff x_0	0.025
Approve iff x_0	Always-reject	0.025
Always-reject	Always-approve	0.050
Always-reject	Approve iff x_1	0.010
Always-reject	Approve iff x_0	0.025
Always-reject	Always-reject	0.025

The CV bodies of fake applicants used in the holdout set are divided into "canonical types" such that each type elicits an approval/reject response from the models f_Y and f_Z . Factors influencing this decision could be the prejudice of the committees, the accomplishments listed on the CV etc. (since those preferences went into building the model). U_{YZ} represents of the distribution of these CV types. There are 16 such types, based on the 4 types each per model, as shown above.

For instance, row 1 of the probability table indicates a CV body such that both f_Y and f_Z would approve such a candidate regardless of the perceived race. Row 2 indicates a CV body such that f_Y would always approve such a candidate, but f_Z is biased to only approve such a candidate if they belonged to Race B.

We want to track whether, given a candidate of Race B who passed the CV screening but was rejected for financial aid, this candidate would still be denied financial aid had they been of Race A. In

2106 particular, they care about the fairness metric μ_{ctf} , defined as follows.
 2107

$$2108 \quad \mu_{ctf} := P(Y_{x_1} = 1, Z_{x_1} = 0) - P(Y_{x_1} = 1, Z_{x_0} = 0) = P(y_x, z'_x) - P(y_x, z'_{x'}) = 10\% \quad (38)$$

2109
 2110 The actual value of this measure can be computed from the probabilities in the SCM above. In
 2111 practice, μ_{ctf} can be directly estimated using the counterfactual randomization procedure illustrated
 2112 above in Fig. 5(b).

2113 If the college data scientists instead follow the standard procedure of using only \mathcal{L}_2 -data from a
 2114 Fisherian RCT, they can only estimate \mathcal{L}_2 fairness metrics, such as μ_{int1}, μ_{int2} defined below.

$$2115 \quad \mu_{int1} := P(y_x) \cdot P(z'_x) - P(y_x) \cdot P(z'_{x'}) = 0 \quad (39)$$

$$2116 \quad \mu_{int2} := P(y, z'; do(x)) - P(y, z'; do(x')) = 0 \quad (40)$$

2117
 2118 The \mathcal{L}_2 -metrics μ_{int1}, μ_{int2} show no issues with fairness, and thus fail to capture the insight obtained
 2119 from the \mathcal{L}_3 -metric $\mu_{ctf} = 10\%$. This counterfactual insight helps the college to quantitatively
 2120 characterize the financial hurdles faced by different racial groups in accessing college education, and
 2121 to prevent unfair disparities.
 2122

2123 E.3 EXAMPLE 3 (COUNTERFACTUAL BANDIT POLICIES)

2124
 2125 The SCM used in this hypothetical scenario to generate data is as follows:

$$2126 \quad U_1 \sim \text{Bernoulli}(0.5)$$

$$2127 \quad U_2 \sim \text{Bernoulli}(0.5)$$

$$2128 \quad U_3 \sim \text{Bernoulli}(0.5)$$

2129
 2130

$$2131 \quad X \leftarrow U_1 \oplus U_2 \quad (\oplus \text{ is the XOR function})$$

$$2132 \quad D \leftarrow X \oplus U_3$$

2133

2134 Since Y is a function of X , the average outcome is shown below for different realizations of the
 2135 latents

2136 i. Avg. Y , when $U_3 = 0$

2137

	$U_3 = 0$			
	$U_1 = 0$		$U_1 = 1$	
	$U_2 = 0$	$U_2 = 1$	$U_2 = 0$	$U_2 = 1$
$do(x_0)$	0.6	0.9	0.8	0.5
$do(x_1)$	0.9	0.6	0.5	0.8

2144 Natural choice of X marked **bold**

2145

2146 ii. Avg. Y , when $U_3 = 1$

2147

	$U_3 = 1$			
	$U_1 = 0$		$U_1 = 1$	
	$U_2 = 0$	$U_2 = 1$	$U_2 = 0$	$U_2 = 1$
$do(x_0)$	0.8	0.7	0.6	0.7
$do(x_1)$	0.7	0.8	0.7	0.6

2154 iii. Avg. Y , with U_3 marginalized (consolidating i. and ii.)

2155

	$U_1 = 0$		$U_1 = 1$	
	$U_2 = 0$	$U_2 = 1$	$U_2 = 0$	$U_2 = 1$
$do(x_0)$	0.7	0.8	0.7	0.6
$do(x_1)$	0.8	0.7	0.6	0.7

2159

2160 Let us call the social media user Alice, for ease of reference. U_1, U_2, U_3 are latent attributes affecting
 2161 Alice’s decisions each evening. In particular, U_1 indicates whether she is tired, U_2 indicates whether
 2162 she had a busy day and is distracted, U_3 indicates whether she is hungry, on any given evening.

2163 If Alice is either tired but mentally relaxed ($X = 1 \oplus 0$), or if she is physically energetic but distracted
 2164 ($X = 0 \oplus 1$), Alice decides to take a walk and use social media via mobile app. If Alice is neither
 2165 tired nor distracted, she prefers to continue working on her desktop and uses social media via desktop
 2166 app during breaks ($X = 0 \oplus 0$). If she is both tired and distracted, she also decides to use the social
 2167 media app on her desktop because she has no energy to take a walk ($X = 1 \oplus 1$).

2168 There are so many possible factors affecting her decisions, Alice is unaware that these are the specific
 2169 unconscious causes of her natural choices. However, the social media company’s unscrupulous data
 2170 scientists surveil U_1, U_2, U_3 (perhaps by tracking Alice’s wearable health monitor and calendar)
 2171 and predict her natural choice. The company then uses behavioural insights to ping Alice with the
 2172 precise notifications and content to maximize her time spent on the platform for each realization of
 2173 U_1, U_2, U_3 .

2174 D is the type of ads Alice sees when she logs in to the social media app for the day.

2175 The detailed causal diagram is shown in Fig. 19.

2176

2177

2178

2179

2180

2181

2182

2183

2184

2185

2186

2187

2188

2189

2190

2191

2192

2193

2194

2195

2196

2197

2198

2199

2200

2201

2202

2203

2204

2205

2206

2207

2208

2209

2210

2211

2212

2213

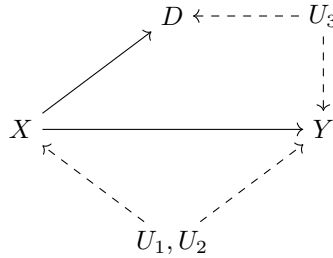


Figure 19: Causal diagram for Example 3. X : app usage type; D : advertisement-type received; Y : compliance with app usage time-limit; U_1 : agent tiredness; U_2 : agent busyness earlier in the day; U_3 : indicator of whether the agent is hungry.

\mathcal{L}_1 -regime: The observational data is contained in Table (iii) in the SCM above, where the bold values correspond to Alice’s natural choices. Given that all combinations of latents happen with equal probability, it is easy to see that the expected reward in the observational regime is $E[Y] = (0.25)(0.7 + 0.7 + 0.6 + 0.6) = 0.65$.

\mathcal{L}_2 -regime: Applying the interventions $do(x_0), do(x_1)$, we can compute the expected outcome from the SCM as shown in Table 2. This is simply the average of all the values in Table (iii) of the SCM above.

	$E[Y; do(x)]$
$do(x_0)$	0.7
$do(x_1)$	0.7

An interventional strategy of randomizing ones actions (or fixing a constant action) outperforms the observational \mathcal{L}_1 regime of allowing one’s actions to be determined by natural inclination.

Table 2: Expected outcome $E[Y_x]$ computed under the interventional regime.

\mathcal{L}_3 -regime - ETT: By counterfactual randomization Alice can sample from the \mathcal{L}_3 distribution $P(Y_x, X)$. She records her natural choice $X = x'$ on a particular evening (what she would have normally done) and randomizes the choice of X that she actually undertakes, during the explore phase. Using this distribution, she then performs the following action, for the natural $X = x'$ that she observes in the exploit phase:

$$do(X = \arg \max_x E[Y_x | x'])$$

We can compute this from Table (iii) of the SCM. Alice simply chooses to do the opposite of what she naturally feels like doing (corresponding to the the non-bold cells of the Table). This ”ETT” \mathcal{L}_3

2214 strategy yields an expected outcome of
 2215

2216
$$\sum_{x'} P(x') \cdot \max_x E[Y_x | x'] = (0.5)[0.7 + 0.8] = 0.75,$$

 2217

2218 outperforming both \mathcal{L}_1 and \mathcal{L}_2 strategies.
 2219

2220 Of course, the explore-exploit phases are combined adaptively in a bandit algorithm like Thompson
 2221 Sampling.

2222 **\mathcal{L}_3 -regime - Optimal:** Finally, Alice leverages her ability to perform *path-specific* randomization
 2223 to sample from the distribution $P(Y_x, X, D_{x'})$. She then adapts a bandit algorithm to performs the
 2224 following actions in the exploit phase:
 2225

2226
$$\text{READ}(X) = x'$$

 2227
 2228
$$\text{CTF-WRITE}(x'' \rightarrow D), \text{ where } x'' = \arg \max_{x''} \left(\max_x E[Y_x | x', D_{x''}] \right); \text{READ}(D) = d$$

 2229
 2230
$$\text{CTF-WRITE}(x \rightarrow Y), \text{ where } x = \arg \max_x E[Y_x | x', d_{x''}]$$

 2231
 2232 ,

2233 where CTF-WRITE is simply the deterministic equivalent of CTF-RAND.

2234 In words, during the explore phase, Alice gathers data on which arm x optimizes $\mathbb{E}[Y_x | x', d_{x''}]$, for
 2235 all x', x'', d . Then, during the exploit phase, Alice first observes $D_{x''} = d$ and $X = x'$, and then
 2236 performs the action x which maximizes her outcome Y_x . Performing another optimization over the
 2237 x'' gives her the best global optimum of $E[Y_x | x', d_{x''}]$.

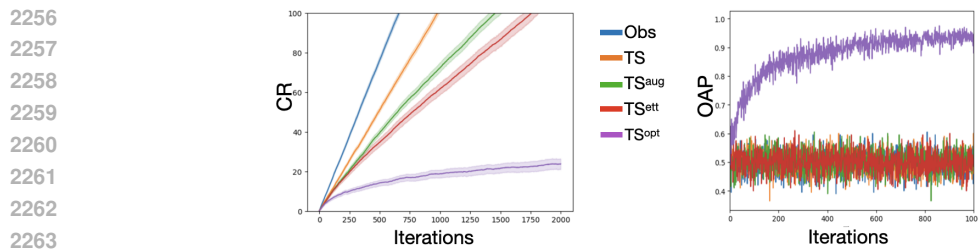
2238 Again, the explore-exploit phases are not separated in a bandit algorithm like Thompson Sampling,
 2239 but incorporated adaptively.

2240 Computing this from the SCM, suppose Alice chooses to record $D_{x_0} = 0 \oplus U_3 = U_3$.

- 2241
- 2242 • When $D_{x_0} = U_3 = 0$, Alice sees according to Table (i) of the SCM that the optimal strategy
 - 2243 is to choose the opposite of what she naturally feels like doing (the values not in bold),
 - 2244 giving the expected outcome $E[Y_x | x', D_{x_0} = 0]$, where $x \neq x'$, as $(0.5)[0.9 + 0.8] = 0.85$
 - 2245
 - 2246 • When $D_{x_0} = U_3 = 1$, Alice sees according to Table (ii) of the SCM that the optimal
 - 2247 strategy is to go with her natural inclination (the values in bold), giving the expected
 - 2248 outcome $E[Y_{x'} | x', D_{x_0} = 1] = (0.5)[0.8 + 0.7] = 0.75$
 - 2249
 - 2250 • Overall, since both values of D_{x_0} are equally likely, this strategy yields an expected outcome
 - 2251 of $0.5[0.85 + 0.75] = 0.8$, which outperforms $\mathcal{L}_1, \mathcal{L}_2$ and \mathcal{L}_3 -ETT strategies.

2252 This walk-through can be repeated identically from the SCM had Alice chosen to measure D_{x_1}
 2253 instead.

2254 E.3.1 SIMULATIONS



2264 Figure 20: Example 3: Cumulative Regret (CR) and Optimal Arm Probability (OAP) for all strategies
 2265 tested via Thompson Sampling.
 2266

2267 The simulation compares the performance of four algorithms

2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

- TS is the conventional Thompson Sampling algorithm that optimizes the \mathcal{L}_2 learning objective $\mathbb{E}[Y; do(x)]$;
- TS^{aug} is a contextual Thompson Sampling algorithm that treats $\{X = x', D_{x''} = d\}$ as merely some extra context variables in each round, and ignoring the \mathcal{L}_3 significance of these variables;
- TS^{ett} is given in Algorithm 5, implementing the ETT baseline strategy described earlier;
- TS^{opt} is given in Algorithm 4, implementing the \mathcal{L}_3 -optimal strategy described earlier.

Importantly, TS^{opt} doesn't treat $\{X = x', D_{x''} = d\}$ merely as extra context variables. Rather, the counterfactual significance of these variables is leveraged via the consistency property

$$E[Y_x | x, d_x] = E[Y | x, d] \tag{41}$$

This means that for several arms being explored, the r.h.s allows us to hot-start the Thompson Sampling using offline (\mathcal{L}_1) data, as implemented in Line 18 of Algorithm 4. This allows for a dramatically faster convergence of the purple vs. green plot in Fig. 20.

Simulations were run for 2000 iterations, 200 epochs (Confidence Interval = 95%).

2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375

Algorithm 4 TS^{opt} : Thompson Sampling OPTIMAL (Bernoulli-Beta case)

```

1: Input: No. of timesteps,  $T$ ; Observational data,  $P(\mathbf{v})$ 
2: for  $x'' \in \text{Domain}(X)$  do
3:   for  $x' \in \text{Domain}(X)$  do
4:      $\alpha_D[x''] [x'] \leftarrow 1$ 
5:      $\beta_D[x''] [x'] \leftarrow 1$  {Initializing  $D$ -priors}
6:   end for
7: end for
8: for  $i \in \text{Domain}(X)$  do
9:   for  $j \in \text{Domain}(X)$  do
10:    for  $d \in \text{Domain}(D)$  do
11:      for  $k \in \text{Domain}(X)$  do
12:         $\alpha_Y[x_i][x_j][d][x_k] \leftarrow 1$ 
13:         $\beta_Y[x_i][x_j][d][x_k] \leftarrow 1$  {Initializing  $Y$ -priors}
14:      end for
15:    end for
16:  end for
17: end for
18:  $t = 1$ 
19: while  $t \leq T$  do
20:   Perform  $\text{READ}(X) = x'$ , for unit
21:   for  $j \in \text{Domain}(X)$  do
22:      $\mu_j^D \sim \text{Beta}(\alpha_D[x''] [x_j], \beta_D[x''] [x_j])$ 
23:   end for
24:   Perform  $\text{CTF-WRITE}(x' \rightarrow D)$  for  $x' = x_j; j := \arg \max_{j'} \mu_{j'}^D$ 
25:   Perform  $\text{READ}(D) = d$ , for unit {Get value of  $D_{x''}$ }
26:   for  $k \in \text{Domain}(X)$  do
27:     if  $x_k = x' = x''$  then
28:        $\mu_k^Y \leftarrow E[Y \mid x'', d]$  {Hot-start using obs. data}
29:     else
30:        $\mu_k^Y \sim \text{Beta}(\alpha_Y[x''] [x'] [d][x_k], \beta_Y[x''] [x'] [d][x_k])$ 
31:     end if
32:   end for
33:   Perform  $\text{CTF-WRITE}(x \rightarrow Y)$  for  $x = x_k; k := \arg \max_{k'} \mu_{k'}^Y$ 
34:   Perform  $\text{READ}(Y) = y$ , for unit {Get value of  $Y_x$ }
35:    $\alpha_D[x''] [x'] \leftarrow \alpha_D[x''] [x'] + y$ 
36:    $\beta_D[x''] [x'] \leftarrow \beta_D[x''] [x'] + 1 - y$  {Update  $D$ -priors}
37:   if  $\neg(x = x' = x'')$  then
38:      $\alpha_Y[x''] [x'] [d][x] \leftarrow \alpha_Y[x''] [x'] [d][x] + y$ 
39:      $\beta_Y[x''] [x'] [d][x] \leftarrow \beta_Y[x''] [x'] [d][x] + 1 - y$  {Update  $Y$ -priors}
40:   end if
41:    $t \leftarrow t + 1$ 
42: end while

```

2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429

Algorithm 5 TS^{ett} : Thompson Sampling ETT (Bernoulli-Beta case)

```
1: Input: No. of timesteps,  $T$ ; Observational data,  $P(\mathbf{v})$ 
2: for  $z \in \text{Domain}(Z)$  do
3:   for  $x' \in \text{Domain}(X)$  do
4:      $\alpha[z][x'] \leftarrow 1$ 
5:      $\beta[z][x'] \leftarrow 1$  {Initializing priors}
6:   end for
7: end for
8:  $t = 1$ 
9: while  $t \leq T$  do
10:  Perform  $\text{READ}(Z) = z$ , for unit
11:  Perform  $\text{READ}(X) = x'$ , for unit
12:  for  $i \in \text{Domain}(X)$  do
13:    if  $x_i = x'$  then
14:       $\mu_i \leftarrow E[Y \mid x', z]$  {Hot-start using obs. data}
15:    else
16:       $\mu_i \sim \text{Beta}(\alpha[z][x_i], \beta[z][x_i])$ 
17:    end if
18:  end for
19:  Perform  $\text{CTF-WRITE}(x \rightarrow Y)$  where  $x = x_i$  s.t.  $i := \arg \max_{i'} \mu_{i'}$ 
20:  Perform  $\text{READ}(Y) = y$ , for unit {Get value of  $Y_x$ }
21:  if  $x \neq x'$  then
22:     $\alpha[z][x'] \leftarrow \alpha[z][x'] + y$ 
23:     $\beta[z][x'] \leftarrow \beta[z][x'] + 1 - y$  {Update priors}
24:  end if
25:   $t \leftarrow t + 1$ 
26: end while
```

2430 F OPTIMALITY RESULT AND PROOF

2431
2432
2433
2434
2435
2436
2437
2438
2439

In this appendix, we identify a strategy that is provably optimal for decision-making in multi-arm bandit (MAB) problems. To focus the discussion, we define a generic *MAB template* (Fig. 21) that is generally representative of a broad class of bandit problems in the literature (the discussion can also be extended to other settings such as sequential or Markov decision processes in future work). X is the decision variable, Z is a context variable, Y is the reward, and D is a descendant of X confounded with Y .

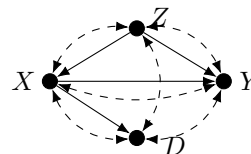


Figure 21: MAB template.

2440 **Definition F.1** (Decision strategy). Given a decision problem following the MAB template (Fig. 21),
2441 a *decision strategy* π is a mapping from a set of variables \mathbf{W}_* (possibly counterfactual) to a set of
2442 actions \mathcal{A} involving decision variable X . The expected reward of following this strategy is notated
2443 $\mu_\pi := \mathbb{E}[Y_{\mathcal{A}} \mid \mathbf{W}_*]$, where $Y_{\mathcal{A}}$ is the potential response of Y under the actions \mathcal{A} .⁵ ■

2444 *Example.* The \mathcal{L}_1 strategy of simply observing the natural behavior of some behavioral agent is
2445 $\pi^{\text{obs}} : \{\} \mapsto \{\}$, which incurs the observational reward of $\mu_{\pi^{\text{obs}}} = \mathbb{E}[Y]$.
2446

2447 *Example.* As discussed in Sec. 4.2, the typical approach in the RL literature is the \mathcal{L}_2 strategy
2448 $\pi^{\text{int}} : \{z\} \mapsto \{\text{WRITE}(X : x^*)\}$, where $x^* := \arg \max_x \mathbb{E}[Y_x \mid z]$. In words, this strategy involves
2449 observing context $Z = z$ for a each round, and then performing the intervention $do(x)$ that maximizes
2450 the \mathcal{L}_2 quantity $\mathbb{E}[Y_x \mid z]$, also known as the *conditional average treatment effect*, or CATE.

2451 *Example.* As discussed in Secs. 2, 4.2, a valid counterfactual strategy would be $\pi^{\text{ett}} : \{x', z\} \mapsto$
2452 $\{\text{CTF-WRITE}(x^* \rightarrow Y)\}$, where $x^* := \arg \max_x \mathbb{E}[Y_x \mid x', z]$.⁶ In words, this strategy involves
2453 observing context $Z = z$ and the unit's natural inclination $X = x'$ for each each round, and then
2454 performing the intervention $do(x)$ that maximizes the \mathcal{L}_3 quantity $\mathbb{E}[Y_x \mid x', z]$, related to the *effect*
2455 *of the treatment on the treated*, or ETT.

2456 In Example 3, we introduced a superior counterfactual strategy

2457
$$\pi^{\text{opt}} : \{X, Z, D_{x''}\} \mapsto \{\text{CTF-WRITE}(x \rightarrow Y), \text{CTF-WRITE}(x'' \rightarrow D)\}, \quad (42)$$

2458 where $x, x'' := \arg \max_{x, x''} \mathbb{E}[Y_x \mid Z, X, D_{x''}]$.
2459

2460 With minor abuse of notation, this is the strategy that (1) observes $Z = z, X = x'$ for a round;
2461 (2) maps from $\{z, x'\} \mapsto x''$, to perform the counterfactual intervention $\text{CTF-WRITE}(x'' \rightarrow D)$ to
2462 observe $D_{x''} = d$; and (3) maps from $\{z, x', d_{x''}\} \mapsto x$, to perform the counterfactual intervention
2463 $\text{CTF-WRITE}(x \rightarrow Y)$ that maximizes $\mathbb{E}[Y_x \mid z, x', d_{x''}]$.

2464 For each mapping in (2), (3) yields a local optimum, in expectation over $X, Z, D_{x''}$. Optimizing over
2465 all choices of x'' in (2) yields a global optimum. Translating this to practice, we provide a general
2466 algorithm (Algorithm 6) that adapts any standard MAB solver to implement the optimal \mathcal{L}_3 -strategy
2467 π^{opt} . We provide examples using Thompson Sampling in the Appendix E.3.1 (Algorithms 4,5).

2468 The natural question is whether we can keep going higher up in Layer 3 of the PCH. Could we
2469 construct higher order strategies that map from $\{X, Z, D_{x''}, D_{x'''}\}$ by drawing samples from more
2470 refined counterfactuals? Sadly no, because a distribution like $P(Y_x, X, D_{x''}, D_{x'''})$ is not realizable
2471 (Def. 3.4), and the machinery we developed in Sec. 3 gives us the tools to reason about this.

2472 **Theorem F.2** (Optimality). *Given a decision problem following the MAB template (Fig. 21), π^{opt} is*
2473 *an optimal realizable strategy. I.e., $\mu_{\pi^{\text{opt}}} \geq \mu_\pi, \forall \pi \in \Pi$, the space of realizable strategies.* ■
2474

2475 The significance of this result is that it averts the need to apply Thm. 3.5 and Cor. 3.7 to search
2476 intractably over the space of all possible \mathcal{L}_3 -distributions for which ones are realizable. Of course,
2477 π^{opt} need not be *uniquely* optimal.

2478 **Corollary F.3** (\mathcal{L}_3 -dominance). *Given an MAB decision problem with causal diagram described*
2479 *by the MAB template (Fig. 21), the optimal \mathcal{L}_3 -strategy π^{opt} dominates the \mathcal{L}_1 -strategy π^{obs} and the*
2480 *optimal \mathcal{L}_2 -strategy π^{int} . I.e., $\mu_{\pi^{\text{opt}}} \geq \mu_{\pi^{\text{obs}}}$ and $\mu_{\pi^{\text{opt}}} \geq \mu_{\pi^{\text{int}}}$.* ■

2481 ⁵We use *strategy* interchangeably with *policy* when the context is clear. However, it should be noted that a
2482 *policy* is usually defined w.r.t. a certain policy space, mapping from a fixed domain to actions on X . Here, we
2483 consider different domains to map from.

⁶CTF-WRITE is simply the deterministic equivalent of CTF-RAND (Def. 2.3).

2484 For decades, the Fisherian RCT methodology used to enact π^{int} was deemed to be the "gold standard"
 2485 for decision-making. We show that the \mathcal{L}_3 -strategy π^{opt} is at least as good (and often better) than \mathcal{L}_2 -
 2486 strategies. This means that if MAB solvers UCB, EXP3 etc. were deployed to enact an \mathcal{L}_2 -strategy in
 2487 an environment where π^{opt} is better, the agent would incur linear cumulative regret, since the learning
 2488 approach comes no closer to discovering the optimal strategy as the number of trials increases.

2490 **Algorithm 6** MAB-OPT

2491 1: **Input:** MAB problem following Fig. 21; MAB solver (e.g. UCB, EXP3, TS); No. of rounds T ;
 2492 Obs. data $P(\mathbf{v})$
 2493 2: **for** each z, x' **do**
 2494 3: Initialize D -arms x''
 2495 4: **end for**
 2496 5: **for** each z, x', x'', d **do**
 2497 6: Initialize Y -arms x
 2498 7: If $x = x' = x''$, hot-start using $P(\mathbf{v})$
 2499 8: **end for**
 2500 9: **for** $t \in [T]$ **do**
 2501 10: Observe x', z
 2502 11: Draw D -arm x'' using MAB solver
 2503 12: Perform CTF-WRITE($x'' \rightarrow D$) and get $D_{x''} = d$
 2504 13: Draw Y -arm x using MAB solver
 2505 14: Perform CTF-WRITE($x \rightarrow Y$) and get $Y_x = y$
 2506 15: Update D -arms and Y -arms according to MAB solver rules using y
 2507 16: **end for**

2508
2509
2510
2511
2512 **F.1 PROOFS FOR SECTION THEOREM F.2 AND COROLLARY F.3**

2513 *Remark F.4.* In order for an agent to enact a non-trivial decision strategy $\pi : \{\mathbf{W}_* = \mathbf{w}\} \mapsto \mathcal{A}$, we
 2514 observe that (1) the distribution $P(Y_{\mathcal{A}}, \mathbf{W}_*)$ must be realizable (Def. 3.4); and (2) the agent must be
 2515 able to observe \mathbf{W}_* before performing actions \mathcal{A} . We call this a *realizable decision strategy*, and
 2516 notate the space of all realizable strategies in a MAB problem as Π . ■

2517
2518
2519 **Corollary F.3** : This result follows immediately from Theorem F.2, by simply recognizing that
 2520 $\pi^{\text{int}}, \pi^{\text{obs}} \in \Pi$, the space of realizable strategies (the agent is presumed capable of performing the
 2521 actions $\text{RAND}(X)$, $\text{WRITE}(X : x)$).
 2522

2523 Therefore, $\mu_{\pi^{\text{ctf}}}$ cannot be less than $\mu_{\pi^{\text{int}}}, \mu_{\pi^{\text{obs}}}$, by Theorem F.2. ■

2524
2525
2526 **Theorem F.2** :

2527
2528 From Lemma F.5, all strategies involve mappings, where each mapping maps to one of the following 5
 2529 possible action sets: (1) $\{\}$ (no action); (2) $\text{WRITE}(X : x)$, for some x ; (3) only $\text{CTF-WRITE}(x \rightarrow Y)$
 2530 for some x ; (4) only $\text{CTF-WRITE}(x'' \rightarrow D)$ for some x'' ; or (5) both $\text{CTF-WRITE}(x \rightarrow Y)$,
 2531 $\text{CTF-WRITE}(x'' \rightarrow D)$ for some x, x'' .

2532 Define Π_5 to be the space of strategies where every mapping of each strategy in Π_5 is mapping to a
 2533 pair of actions $\text{CTF-WRITE}(x \rightarrow Y)$, $\text{CTF-WRITE}(x'' \rightarrow D)$ for some x, x'' . I.e., all mappings only
 2534 involve possibility (5) under these strategies, for all any unit encountered in the decision problem.

2535 Let π_5 be an optimal strategy in this space. I.e., $\pi_5 \in \arg \max_{\pi \in \Pi_5} \mu_{\pi}$.

2536 By Lemma F.7, π_5 is also an optimal strategy in the space of all possible strategies. This means we
 2537 only need to consider strategies whose mappings are mappings to a pair of CTF-WRITE procedures.

2538 Let \mathbf{W}_* be the context used by π_5 . If \mathbf{W}_* does not already contain the natural variables X, Z , we
 2539 can always define π'_5 that use context $\mathbf{W}'_* = \mathbf{W}_* \cup \{X, Z\}$ s.t. $\mu_{\pi'_5} = \mu_{\pi_5}$, where π'_5 simply ignores
 2540 the extra context variable in the mapping.

2541 Such a move would not affect the realizability of π'_5 because CTF-WRITE does not override the
 2542 natural value of X , and both X, Z can be observed before decision-making.
 2543

2544 Combinatorially, there are only 3 possibilities for picking each mapping in π'_5 .

- 2545 1. Mapping from $\{x', z\}$ to a pair of actions $\{\text{CTF-WRITE}(x \rightarrow Y), \text{CTF-WRITE}(x'' \rightarrow D)\}$
- 2546 2. Mapping from $\{x', z\}$ to some $\text{CTF-WRITE}(x \rightarrow Y)$, observing $Y_x = y$, and mapping from
 2547 $\{x', z, y_x\}$ to $\text{CTF-WRITE}(x'' \rightarrow D)$; or
- 2548 3. Mapping from $\{x', z\}$ to some $\text{CTF-WRITE}(x'' \rightarrow D)$, observing $D_{x''} = d$, and mapping
 2549 from $\{x', z, d_{x''}\}$ to $\text{CTF-WRITE}(x \rightarrow Y)$
 2550

2551 We can use similar arguments to Lemma F.7, where we restricted our attention to the space of
 2552 strategies Π_5 which could mimic all other optimal strategies.
 2553

2554 Possibility 2 can be mimicked by some mapping following possibility 1 which maps to a joint pair of
 2555 actions. The two are equivalent in terms of outcome, because conditioning on y_x to choose x'' does
 2556 not affect the outcome Y . So we can restrict our attention to possibilities 1 and 2.

2557 Each mapping of possibility 1 can be mimicked by possibility 3, where the extra step of conditioning
 2558 on $d_{x''}$ just ignores the extra information about $d_{x''}$. Thus, we can replace all mappings in the optimal
 2559 strategy π'_5 with mappings of possibility 3, to get a strategy π''_5 that also performs optimally.
 2560

2561 Since there are two mappings in π''_5 , they must be the mappings which maximize the outcome.

2562 This is precisely the definition of the strategy π^{opt} given in Equation 42 and in the description
 2563 immediately following it.

2564 ■

2565
 2566 **Lemma F.5.** *Any decision strategy π for a decision problem having causal structure same as the*
 2567 *MAB template is s.t. each mapping of the strategy maps from domain of the context to one of the*
 2568 *five following possible sets of actions: (1) $\{\}$ (no action); (2) $\text{WRITE}(X : x)$, for some x ; (3)*
 2569 *only $\text{CTF-WRITE}(x \rightarrow Y)$ for some x ; (4) only $\text{CTF-WRITE}(x'' \rightarrow D)$ for some x'' ; or (5) both*
 2570 *$\text{CTF-WRITE}(x \rightarrow Y), \text{CTF-WRITE}(x'' \rightarrow D)$ for some x, x'' .*
 2571

2572 *Proof.* Since the physical action space only involves doing nothing, WRITE or CTF-WRITE.
 2573 Any other combination would be equivalent to one of the 5 above. E.g., $\text{WRITE}(X : x)$ and
 2574 $\text{CTF-WRITE}(x'' \rightarrow D)$ is the equivalent to the pair $\text{CTF-WRITE}(x \rightarrow Y)$ and $\text{CTF-WRITE}(x'' \rightarrow D)$
 2575 (see Remark D.5).

2576 We ignore randomized actions for simplicity. From standard results in learning theory, there is an
 2577 optimum to be found at a simplex corner so we need only search over the space of hard interventions.
 2578 ■

2580
 2581 **Lemma F.6.** *The context \mathbf{W}_* used in the strategy $\pi : \{\mathbf{W}_* = \mathbf{w}\} \mapsto \mathcal{A}$ can only possibly contain a*
 2582 *subset of $X, Z, D, D_{x''}$ for some x'' , and at most one potential response of D .*
 2583

2584 *Proof.* There are only 4 variables to consider: X, Y, Z, D .

2585 By the definition of a realizable strategy (Remark F.4), we need $P(Y_{\mathcal{A}}, \mathbf{W}_*)$ to be realizable. By Cor.
 2586 3.7 there cannot be two potential responses of the same variable in a realizable distribution. This
 2587 rules out any other potential response of Y , and ensures only one each of X, Z, D .

2588 Since the only possible actions are interventions involving X , which do not affect Z and X (natural
 2589 variable), these are the only potential responses that could appear involving these variables.
 2590

2591 Likewise, with D, D (natural value) and $D_{x''}$ are the only possible potential responses that could
 appear, and at most one of them can. ■

2592 **Lemma F.7.** *If π_5 is an optimal strategy in Π_5 , the set of all strategies which map to a pair of*
 2593 *CTF-WRITE procedures, then π_5 is also an optimal strategy in the set of all strategies possible in the*
 2594 *MAB decision problem.*
 2595

2596 *Proof.* Let Π_1 be the space of all strategies possible in the problem. Note that $\Pi_5 \subseteq \Pi_1$. Let $\pi_1 \notin \Pi_5$
 2597 be an optimal strategy. I.e. $\pi_1 \in \arg \max_{\pi \in \Pi_1} \mu_\pi$. If no such π_1 the Lemma stands proved.

2598 Let \mathbf{W}_* be the context used by π_1 . If \mathbf{W}_* does not already contain the natural variable X , we can
 2599 always define π'_1 that uses context $\mathbf{W}'_* = \mathbf{W}_* \cup \{X\}$ s.t. $\mu_{\pi'_1} = \mu_{\pi_1}$, where π'_1 simply ignores the
 2600 extra context variable in the mapping. For now, it doesn't matter whether such π'_1 is realizable or not.
 2601 Just that it is also an optimal strategy.
 2602

2603 Each mapping in the strategy π'_1 maps from the domain of \mathbf{W}'_* to one of the five possible action
 2604 sets mentioned in Lemma F.5. E.g., for some $\mathbf{W}'_* = \mathbf{w}$, the strategy π'_1 maps this to $\mathbf{w} \mapsto \{\}$ or
 2605 $\mathbf{w} \mapsto \text{WRITE}(X)$.

2606 Consider a mapping in π'_1 from the domain of $\mathbf{W}'_* = \{x', \dots\}$ to possibility (1), empty set of actions
 2607 (recall, the context includes natural X). Such a mapping can be mimicked by an equivalent mapping
 2608 $\mathbf{W}'_* = \{x', \dots\} \mapsto \{\text{CTF-WRITE}(x' \rightarrow Y), \text{CTF-WRITE}(x' \rightarrow D)\}$. By the consistency property if
 2609 $X(\mathbf{u}) = x'$, then $Y_{x'}(\mathbf{u}) = Y(\mathbf{u})$ and $D_{x'}(\mathbf{u}) = D(\mathbf{u})$.

2610 Thus, we can replace all the mappings in π'_1 that involve a mapping to the empty set of actions, with
 2611 an equivalent pair of CTF-WRITE using the natural value of X observed in the context. Call this new
 2612 strategy π_2 . π_2 is as good as π'_1 because the mappings are all equivalent. Thus, π_2 is also optimal in
 2613 Π_1 . Again, it doesn't matter that π_2 may not be realizable, just that it is optimal.

2614 Next, consider a mapping in π_2 from the domain of $\mathbf{W}'_* = \{x', \dots\}$ to possibility (2), some action
 2615 $\text{WRITE}(X : x)$. Such a mapping can be mimicked by an equivalent mapping $\mathbf{W}'_* = \{x', \dots\} \mapsto$
 2616 $\{\text{CTF-WRITE}(x \rightarrow Y), \text{CTF-WRITE}(x \rightarrow D)\}$. The evaluation of $f_Y(x, Z, \mathbf{u})$ in both scenarios is
 2617 the same, with the only difference being that f_X is overwritten, which doesn't affect the outcome Y
 2618 for each \mathbf{u} . I.e., the outcome Y would be the same for every unit under both strategies.

2619 Thus, we can replace all the mappings in π_2 that involve a mapping to some action $\text{WRITE}(X : x)$,
 2620 with an equivalent pair of CTF-WRITE. Call this new strategy π_3 . π_3 is as good as π_2 because the
 2621 mappings are all equivalent in terms of outcome. Thus, π_3 is also optimal in Π_1 .
 2622

2623 Next, consider a mapping in π_3 from the domain of $\mathbf{W}'_* = \{x', \dots\}$ to possibility (3), some
 2624 action $\text{CTF-WRITE}(x \rightarrow Y)$. Such a mapping can be mimicked by an equivalent mapping
 2625 $\mathbf{W}'_* = \{x', \dots\} \mapsto \{\text{CTF-WRITE}(x \rightarrow Y), \text{CTF-WRITE}(x' \rightarrow D)\}$ for natural value x' . By the
 2626 consistency property, if $X(\mathbf{u}) = x'$ then $D_{x'}(\mathbf{u}) = D(\mathbf{u})$.

2627 Thus, we can replace all the mappings in π_3 that involve a mapping to some action $\text{CTF-WRITE}(x \rightarrow$
 2628 $Y)$, with an equivalent pair of CTF-WRITE. Call this new strategy π_4 . π_4 is as good as π_3 because
 2629 the mappings are all equivalent in terms of outcome. Thus, π_4 is also optimal in Π_1 .

2630 Next, consider a mapping in π_4 from the domain of $\mathbf{W}'_* = \{x', \dots\}$ to possibility (4), some
 2631 action $\text{CTF-WRITE}(x \rightarrow D)$. Such a mapping can be mimicked by an equivalent mapping
 2632 $\mathbf{W}'_* = \{x', \dots\} \mapsto \{\text{CTF-WRITE}(x' \rightarrow Y), \text{CTF-WRITE}(x \rightarrow D)\}$ for natural value x' . By the
 2633 consistency property, if $X(\mathbf{u}) = x'$ then $Y_{x'}(\mathbf{u}) = Y(\mathbf{u})$.

2634 Thus, we can replace all the mappings in π_4 that involve a mapping to some action $\text{CTF-WRITE}(x \rightarrow$
 2635 $D)$, with an equivalent pair of CTF-WRITE. Call this new strategy π'_5 . π'_5 is as good as π_4 because
 2636 the mappings are all equivalent in terms of outcome. Thus, π'_5 is also optimal in Π_1 .
 2637

2638 However, note that the only possible mappings in π'_5 are possibility (5) involving a pair of CTF-WRITE
 2639 actions. Which means $\pi'_5 \in \Pi$.

2640 Thus, we show that all optimal strategies in Π_5 are also optimal in the overall space of strategies. ■
 2641
 2642
 2643
 2644
 2645