

# Decaf: A Deconfounding Causal Generative Model

Anonymous Authors<sup>1</sup>

## Abstract

Causal generative models (CGMs) have recently emerged as capable approaches to simulate the causal mechanisms generating our observations, enabling causal inference. Unfortunately, existing approaches either are *overly restrictive*, assuming the absence of hidden confounders, or *lack generality*, being tailored to a particular query and graph. In this work, we introduce Decaf, a CGM that accounts for hidden confounders in a single amortized training process using only observational data and the causal graph. Importantly, Decaf can provably identify all causal queries with a valid adjustment set or sufficiently informative proxy variables. Remarkably, for the first time to our knowledge, we show that a confounded counterfactual query is identifiable, and thus solvable by Decaf, as long as its interventional counterpart is as well. Our empirical results on diverse settings—including the Ecoli70 dataset, with 3 independent hidden confounders, tens of observed variables and hundreds of causal queries—show that Decaf outperforms existing approaches, while demonstrating its out-of-the-box flexibility.

## 1 Introduction

Causal queries, or *what if* questions, seek to determine how changes in one variable affect another, which is crucial to evaluate the effects of interventions in fields such as healthcare (Feuerriegel et al., 2024), marketing policies (Varian, 2016) or education (Zhao & Heffernan, 2017). Importantly, when empirical trials are infeasible due to ethical, financial, or practical constraints, answering causal queries from observational data becomes essential.

To address this challenge, causal generative models (CGMs) (Javaloy et al., 2023; Chao et al., 2023; Khemakhem et al., 2021) have recently emerged as powerful and flexible tools for modelling structural causal models (SCMs), allowing

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

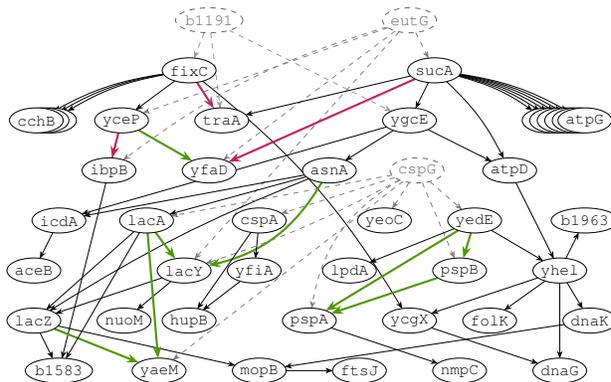


Figure 1: **Decaf can be effortlessly applied to highly complex causal graphs**, such as that of the Ecoli70 dataset (Schäfer & Strimmer, 2005), with multiple independent hidden confounders and dozens of variables. We dash hidden confounders, and highlight direct *confounded* effects that are now **identifiable**, or still **unidentifiable**, with Decaf.

for efficiently sampling interventional and counterfactuals distributions, and enabling the estimation of any causal query of interest. However, all existing CGMs also assume *causal sufficiency*, i.e., that all confounders are observed.

However, *causal sufficiency* is rarely satisfied in practice, making *hidden confounding* a major challenge in causality, as it generally renders causal queries *unidentifiable*, i.e., that they cannot be uniquely expressed as a function of the observations. While recent advances have shown that some confounded causal queries are identifiable if there exist sufficiently informative proxies of the hidden confounders (Miao et al., 2018; 2023; Wang & Blei, 2021), these approaches are still limited to specific intervention-outcome pairs and do not allow for counterfactual estimation.

Our objective is to bridge the gap between these two lines of work. To this end, we introduce the *deconfounding causal normalizing flow* (Decaf ) , to the best of our knowledge, the first CGM that allows the estimation of any *identifiable* causal query—including *counterfactuals*—in the presence of *hidden confounders*, with only observational data, the causal graph, and a single amortized training process. More in detail, Decaf resembles variational autoencoders (Kingma & Welling, 2014) as it is trained with an ELBO and comprises: **i**) a causal normalizing flow (CNF) (Javaloy et al.,

2023) as “decoder”, adapted to be conditioned on the (potentially many) hidden confounders; and **ii**) a conditional normalizing flow (Winkler et al., 2019) as “encoder”, computing the posterior distribution of the hidden confounders.

Furthermore, we theoretically demonstrate that Decaf accurately estimates all identifiable causal queries (interventional and *counterfactual*) for which we can find a valid adjustment set or sufficiently informative proxy variables, significantly extending existing results from prior works (Miao et al., 2018; Wang & Blei, 2021; Javaloy et al., 2023).

All of the above is well illustrated in the Ecoli70 dataset (Schäfer & Strimmer, 2005), whose causal graph is depicted in Fig. 1. Specifically, by training Decaf *once* on this dataset, we can efficiently model all 43 observed variables and 3 independent hidden confounders and, most importantly, *compute any causal query on demand* during deployment. Out of all the direct causal effects (i.e., edges) in Fig. 1, Decaf can accurately estimate all unconfounded effects, as well as 8 out of the 11 confounded ones. In stark contrast with previous works, Decaf also estimates counterfactual queries, increasing the previous count to 16 identifiable queries.

In order to assist practitioners, we provide algorithms to easily check whether a particular query of interest is identifiable in our framework, and we will make our code publicly available upon acceptance. Moreover, we empirically validate our claims on semi-synthetic and real-world experiments, demonstrating that Decaf outperforms existing alternatives while being widely applicable. Therefore, *Decaf offers a practical and efficient solution for causal inference in the presence of hidden confounding*, bridging the gap between general CGMs and specialized solutions.

## 2 Related Works

We discuss the most relevant works to put Decaf into context, and provide a more detailed literature review in App. D.

**Generative causal models.** In order to faithfully learn a SCM, one common approach consists modeling each variable as a function of its causal parents with an independent model, starting from the root nodes. As of the choice for modeling these functions, prior works range from simple yet well-established additive noise models (ANMs) (Hoyer et al., 2008), to more complex but powerful diffusion-based causal models (DCMs) (Chao et al., 2023), among others (Kocaoglu et al., 2018; Yang et al., 2020; Pawlowski et al., 2020; Parafita & Vitrià, 2022). Due to its nature, this approach typically is parameter-intensive, and can easily overfit and propagate errors to descendant variables.

Alternatively, recent works have explored using a single (structurally-constrained) network to model the SCM at once, e.g., using autoregressive flows (Khemakhem et al., 2021; Javaloy et al., 2023), or graph neural networks (GNNs) (Zečević et al., 2021; Sánchez-Martín et al., 2022).

Among these, the causal normalizing flow (CNF) deserves special attention, given its flexibility and theoretical guarantees, which we discuss later in §3.2. Most importantly, all the approaches above assume *causal sufficiency*, i.e. the absence of hidden confounders, limiting their applicability in settings with hidden confounding.

**Causal inference with latent confounders.** Another line of work relies on structural assumptions for correctly answering causal queries. However, these approaches typically deal only with interventional queries (i.e., not counterfactual ones) and are tailored to a specific causal graph and a single treatment-outcome pair, requiring us to train one model per query. In particular, existing works exploit instrumental variables (IVs) (Angrist & Pischke, 2009) or mediators (Pearl, 2009) to achieve this goal and, more recently, a body of works exploit proxy variables to account for latent confounding (Allman et al., 2009; Kuroki & Pearl, 2014; Kallus et al., 2018; Louizos et al., 2017; Miao et al., 2023; 2018). Of particular interest is the Deconfounder by Wang & Blei (2021), a probabilistic model that interprets multiple treatments as null proxies to find a substitute of the hidden confounders and estimate causal queries.

## 3 Background

### 3.1 Confounded Structural Causal Models

Next, we introduce some ideas from the causality literature used throughout this work to model the causal structure of the data and answer causal queries of interest.

**Definition 1.** A (*confounded*) *Structural Causal Model (SCM)* is a triplet  $\mathcal{M} := (\mathbf{f}, P_{\mathbf{u}}, P_{\mathbf{z}})$  describing a data-generating process over a set of  $D$  observed (endogenous) variables  $\mathbf{x} := (x_1, x_2, \dots, x_D)$  as

$$x_i := f_i(\text{pa}(i), \mathbf{u}_i, \mathbf{z}) \quad \text{for } i = 1, 2, \dots, D, \quad (1)$$

with  $\mathbf{u} := (u_1, u_2, \dots, u_D) \sim P_{\mathbf{u}}, \mathbf{z} \sim P_{\mathbf{z}},$

and where  $f_i$  represents the structural equation to compute the  $i$ -th endogenous variable,  $x_i$ , from its observed *causal parents*,  $\text{pa}(i)$ , the  $i$ -th exogenous variable,  $u_i$ , and the vector of *hidden confounders*,  $\mathbf{z}$ .<sup>1</sup>

Note that, while we make the dependence on the hidden confounders explicit for all observed variables in Eq. 1, we assume w.l.o.g. that a subset of them may not be directly affected by the hidden confounders. Furthermore, given a SCM  $\mathcal{M}$ , we denote by  $\mathcal{G}$  the *faithful* causal graph that it induces, representing *only* the direct causal relationships between pairs of endogenous and hidden variables and, when necessary, also exogenous variables.

One key element in causality is the do operator (Pearl, 2012), denoted by  $\text{do}(\cdot)$ , which conceptualizes the action of ex-

<sup>1</sup>Bold denotes random vectors.

ternally intervening on a treatment variable  $t$ , i.e., to set  $t$  to a fixed value independently of its parents. In turn, the do operator enables the computation of interventional and counterfactual queries in SCMs (Peters et al., 2017), i.e., of population and instance-wise *what if* questions.

**Definition 2.** A causal query  $Q(\mathcal{M}) := p(y \mid \text{do}(t), \mathbf{c})$  is a distribution over  $y \in \mathbf{x}$  (the *outcome* variable), as a result of intervening upon the variable  $t \in \mathbf{x}$  (the *treatment* variable). Additionally,  $Q(\mathcal{M})$  denotes an *interventional* or *counterfactual* query if the variable  $\mathbf{c}$  is, respectively, the empty set or the vector of factual observed values,  $\mathbf{x}^f$ .

However, in the presence of *hidden confounders*, one cannot simply apply the do-operator to evaluate causal queries, as the computations involve the causal parents and the unaccounted confounders would bias the results. Instead, one needs to find alternative ways to compute these quantities if possible, as we discuss in §2.

### 3.2 Causal Normalizing Flows

Causal normalizing flows (CNFs) (Javaloy et al., 2023) play an important role in this work, as they form the basic building blocks of Decaf, given their identifiability guarantees despite a mild set of assumptions.

Similar to Eq. 1, a CNF is defined as a pair  $(T_\theta, P_{\mathbf{u}})$  forming a data-generating process that yields a set of  $D$  endogenous variables as  $\mathbf{x} := T_\theta^{-1}(\mathbf{u})$ , where  $\mathbf{u} \sim P_{\mathbf{u}}$  and  $T_\theta: \mathbb{R}^D \rightarrow \mathbb{R}^D$  is a normalizing flow (Papamakarios et al., 2021). In particular,  $T_\theta$  is a normalizing flow with additional structural constraints, ensuring that it induces the same causal graph as the underlying SCM.

Javaloy et al. (2023) demonstrated that CNFs form a general class of identifiable SCMs, and that they can approximate the underlying SCM as closely as required simply by maximizing the observed joint evidence, i.e.,  $\max_\theta \log p_\theta(\mathbf{x})$ . Moreover, CNFs also allow for efficient sampling of any interventional and counterfactual distribution, enabling their use for complex causal-inference task.

Unfortunately, as discussed in §1, CNFs need to assume causal sufficiency to provide the above guarantees, thus limiting their applicability. In this work, we attempt to address this limitation and account for the presence of hidden confounders without losing theoretical guarantees.

## 4 Problem Statement

In this work, we assume the existence of an unobserved confounded SCM,  $\mathcal{M}$ , as in Def. 1, of which we have access to  $N$  i.i.d. observations and its induced causal graph,  $\mathcal{G}$ .

Our objective is therefore to design a CGM that can *faithfully answer as many causal queries from the original SCM as possible*, despite the presence of unobserved hidden confounders. In other words, to find a substitute model of  $\mathcal{M}$  that we can use to accurately perform causal inference.

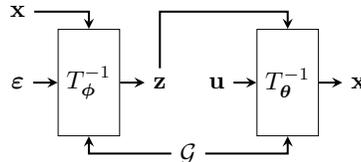


Figure 2: **Sketch of Decaf architecture.**  $T_\phi$  and  $T_\theta$  are conditional normalizing flows, with the top input as condition;  $\mathcal{G}$  is the causal graph, and  $\varepsilon$  is a non-causal random variable needed by the normalizing flow to sample  $\mathbf{z}$ .

**Assumptions.** Regarding the underlying SCM  $\mathcal{M}$ , we simply assume that it **i)** has  $C^1$ -diffeomorphic structural equations,<sup>2</sup> and **ii)** induces an acyclic graph. We denote the family of SCMs meeting these assumptions by  $\mathbb{M}$ .

## 5 Deconfounding Causal Normalizing Flows

To help bridge the gap between CGMs and tailored hidden-confounding solutions, we now present the *deconfounding causal normalizing flow* (or Decaf ) .

Intuitively, Decaf takes a well-grounded CGM such as the causal normalizing flow (Javaloy et al., 2023), which can provably approximate unconfounded SCMs and perform causal inference, and expand it such that it accounts for hidden confounding by building data-driven substitutes of these confounders, an idea that has been successfully explored in the past (Wang & Blei, 2019; 2021; Bica et al., 2020).

Decaf achieves the above by following a similar structure as that of a variational autoencoder (Kingma & Welling, 2014). That is, Decaf comprises two main components. First, an inference network which approximates the *intractable* posterior distribution of the hidden confounders, given their observed children. Second, a generative network that exploits structural constraints to accurately model the underlying SCM, given a substitute for hidden confounders. Each of these parts comes with their own challenges, however, which we now explain in detail:

**Generative network.** As mentioned in §3.2, we use CNFs (Javaloy et al., 2023) as our starting point. However, since our generative model needs to take in hidden confounders as conditional inputs, we adapt CNFs to use conditional normalizing flows (Winkler et al., 2019), instead of unconditional ones. The resulting model,  $T_\theta$ , is thus an invertible transformation describing a data-generating process, conditioned on  $\mathbf{z}$ , which can map a set of exogenous variables  $\mathbf{u}$  to our observations and vice versa, i.e.,

$$T_\theta(\mathbf{x}, \mathbf{z}) = \mathbf{u} \sim P_{\mathbf{u}} \quad \text{and} \quad \mathbf{x} = T_\theta^{-1}(\mathbf{u}, \mathbf{z}), \quad (2)$$

<sup>2</sup>That is, that  $\mathbf{f}$  has inverse and both  $\mathbf{f}$  and  $\mathbf{f}^{-1}$  are continuously differentiable w.r.t. the exogenous variables.

where we further exploit the given causal graph to ensure that the generative process is faithful, i.e., such that

$$p_{\theta}(\mathbf{x} | \mathbf{z}) = \prod_{i=1}^D p_{\theta}(x_i | \text{pa}(i), \mathbf{z}), \quad (3)$$

defining now a process similar to that given in Def. 1. Just as in Def. 1, only the children of  $\mathbf{z}$  will actually condition on  $\mathbf{z}$  in Eq. 3. Furthermore,  $T_{\theta}$  allows us to write down the exact likelihood of the data given  $\mathbf{z}$ ,

$$\log p_{\theta}(\mathbf{x} | \mathbf{z}) = p_{\mathbf{u}}(T_{\theta}(\mathbf{x}, \mathbf{z})) |\det(\nabla_{\mathbf{x}} T_{\theta}(\mathbf{x}, \mathbf{z}))|. \quad (4)$$

**Deconfounding network.** To model the posterior distribution of the hidden confounders given our observations, i.e., the abduction step needed to compute counterfactuals (Pearl, 2009), we use another conditional normalizing flow (Winkler et al., 2019), as it can approximate the true posterior distribution arbitrarily well. Once again, we exploit prior knowledge about the causal graph and mask the resulting network,  $T_{\phi}$ , such that it models each *independent* hidden confounder  $\mathbf{z}_k$  using only its observed children, i.e.,

$$q_{\phi}(\mathbf{z} | \mathbf{x}) = \prod_{k=1}^{D_{\mathbf{z}}} q_{\phi}(\mathbf{z}_k | \text{ch}(\mathbf{z}_k)), \quad (5)$$

where  $D_{\mathbf{z}}$  is the number of independent hidden confounders.

**Training process.** We jointly train both networks defined above as it would be typically done in deep latent-variable models, i.e., during training we *maximize* the evidence lower bound (ELBO) (Kingma & Welling, 2014):

$$\begin{aligned} \mathcal{L}(\theta, \phi) &= \mathbb{E}_{q_{\phi}}[\log p_{\theta}(\mathbf{x} | \mathbf{z})] - \text{KL}[q_{\phi}(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z})] \quad (6) \\ &= \mathbb{E}_{q_{\phi}}[\log p_{\theta}(\mathbf{x}, \mathbf{z})] + \text{H}(q_{\phi}(\mathbf{z} | \mathbf{x})), \quad (7) \end{aligned}$$

where  $p(\mathbf{z})$  is the prior distribution of  $\mathbf{z}$ , KL the Kullback-Leibler divergence (Kullback & Leibler, 1951), and H the differential entropy (Kolmogorov, 1956).

The motivation for this choice is three-fold. First, we want the generative network to explain the observations given samples from  $q_{\phi}$  (first term of Eq. 6). Second, as we do not know the optimal size for  $\mathbf{z}$ , we need to prevent the deconfounding network from allocating information exclusive of  $\mathbf{x}$  in  $\mathbf{z}$  (entropy term in Eq. 7). Finally, all the theory in §6 relies on Decaf matching the data evidence,  $p_{\text{data}}(\mathbf{x})$ , which we encourage Decaf to do since

$$\begin{aligned} \max_{\phi, \theta} \mathcal{L}(\phi, \theta) &= \min_{\phi, \theta} \text{KL}[p_{\text{data}}(\mathbf{x}) \| p_{\theta}(\mathbf{x})] \\ &\quad + \text{KL}[q_{\phi}(\mathbf{z} | \mathbf{x}) \| p_{\theta}(\mathbf{z} | \mathbf{x})]. \quad (8) \end{aligned}$$

**Causal inference.** Since the tuple  $(T_{\theta}, P_{\mathbf{u}}, P_{\mathbf{z}})$  defines a confounded SCM as defined in Def. 1, we can use Decaf

to efficiently sample from observational and interventional distributions by: **i**) sampling  $\mathbf{z}$  from  $p(\mathbf{z})$ ; and **ii**) sampling  $\mathbf{x}$  from either  $p_{\theta}(\mathbf{x} | \mathbf{z})$  or  $p_{\theta}(\mathbf{x} | \mathbf{z}, \text{do}(\mathbf{t}))$ , as proposed by Javaloy et al. (2023). For counterfactual inference, we can use the deconfounding network to perform the induction step, as the second KL term in Eq. 8 shows that it approximates the posterior induced by  $T_{\theta}$  (i.e., its  $\mathbf{z}$ -inverse given  $\mathbf{x}$ ). Therefore, to generate counterfactual samples we simply need to: **i**) sample from  $q_{\phi}(\mathbf{z} | \mathbf{x}^f)$ ; and **ii**) sample again from  $p_{\theta}(\mathbf{x} | \mathbf{z}, \text{do}(\mathbf{t}))$ . We provide more details about these steps and the do-operator in App. C.

## 6 Theoretical Results

We take advantage that our work is at intersection of CGMs and hidden-confounding solutions to leverage and expand the theory of both research fields. While we present here an intuitive summary of our main theoretical results, formal statements and derivations can be found in App. A.

Note that, throughout this section, we assume that Decaf matches the true data evidence, i.e.,  $p_{\text{data}}(\mathbf{x}) = p_{\theta}(\mathbf{x})$ . Given that CNFs (and hence Decaf) are universal density approximators (Papamakarios et al., 2021), we should be able to always meet this assumption, provided enough resources.

### 6.1 Causal Query Identifiability

First, we study which queries are identifiable with Decaf. We call a query identifiable if we are guaranteed to produce the same query distribution as the original SCM by matching the data evidence. More formally, we adopt the following definition (Pearl, 2009, Def. 3.2.4):

**Definition 3.** Let  $Q(\mathcal{M})$  be a causal query of a model  $\mathcal{M}$ . We call  $Q$  *identifiable* if, for any two models  $\mathcal{M}_1, \mathcal{M}_2 \in \mathbb{M}$ ,  $Q(\mathcal{M}_1) = Q(\mathcal{M}_2)$  whenever  $p_{\mathcal{M}_1}(\mathbf{x}) = p_{\mathcal{M}_2}(\mathbf{x}) > 0$ .

Another relevant concept for this section is that of a *valid adjustment set* (Peters et al., 2017, Def. 6.38). In plain terms, if we were to compute a causal query, say  $p(y | \text{do}(\mathbf{t}))$ , a valid adjustment set  $\mathbf{b}$  is a subset of variables such that: **i**) it blocks all backdoor paths between  $y$  and  $\mathbf{t}$ , and **ii**) it is independent of the variable  $\mathbf{t}$  after severing all incoming edges in  $\mathbf{t}$  in the associated causal graph. As a consequence, we can use  $\mathbf{b}$  to apply the adjustment formula,

$$p(y | \text{do}(\mathbf{t})) = \int p(y | \mathbf{t}, \mathbf{b}) p(\mathbf{b}) d\mathbf{b}. \quad (9)$$

Additionally, we refer to  $\mathbf{b}$  as *invalid* if only **i**) holds.

#### 6.1.1 INTERVENTIONAL QUERIES

We first look at the identifiability of interventional queries, i.e., queries of the form  $Q(\mathcal{M}) = p_{\mathcal{M}}(y | \text{do}(\mathbf{t}))$ , where  $y, \mathbf{t} \in \mathbf{x}$  are any two endogenous variables. We summarize our findings in the following proposition, which we properly formalize in App. A.2:

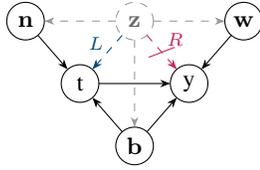


Figure 3: **Generic causal graph** where we are interested in the interventional query  $Q(\mathcal{M}) = p(y \mid \text{do}(t))$ . Blue and red edges play a crucial role, as their **presence** or **absence** induce different types of identifiability conditions.

**Proposition 6.1** (Informal). *Decaf is able to identify a given interventional causal query if one of the following exists:*

- i) a valid adjustment set  $\mathbf{b}$  not containing  $\mathbf{z}$ ,
- ii) an invalid one where  $p(\mathbf{b} \mid \text{do}(t))$  is identifiable, or
- iii) sufficiently informative proxy and null proxy variables.

To help us go through the requirements in Prop. 6.1, let us break them down with the example depicted in Fig. 3 where, depending on the **presence** or **absence** of edges L and R, we face qualitatively different identifiability scenarios:

**1. Unconfounded case, LR.** If neither treatment nor outcome are directly influenced by  $\mathbf{z}$ , then we can always find a valid adjustment set that does not include  $\mathbf{z}$ . We extend the results of Javaloy et al. (2023) to show that Decaf can identify any interventional causal query of this type.

**2. Confounded-treatment case, LR.** If only the treatment is directly affected by  $\mathbf{z}$ , we run into two possible scenarios. First, if we are able to find a valid adjustment set e.g.,  $\mathbf{b}$  and  $\mathbf{w}$  in Fig. 3, then Decaf can always identify the interventional query. Otherwise, Decaf could still identify the query if we find an *invalid* adjustment set where  $p(\mathbf{b} \mid \text{do}(t))$  is still identifiable by Decaf.

**3. Confounded-outcome case, LR.** When only the outcome variable directly depends on  $\mathbf{z}$ , Decaf can identify any interventional query, as it necessarily exists a valid adjustment set not containing  $\mathbf{z}$ . In our running example, variables  $\mathbf{n}$  and  $\mathbf{b}$  would block all backdoor paths in Fig. 3, and Decaf would properly estimate the interventional query.

**4. Fully-confounded case, LR.** When both variables directly depend on  $\mathbf{z}$ , identifiability is more challenging, as any adjustment set necessarily involves the hidden confounder. In this case, we extend in Prop. A.2 the results from Miao et al. (2018) and Wang & Blei (2021) to allow for *general causal graphs and additional covariates*. In short, we find that an interventional query is identifiable if we can find: **i)** a proxy  $\mathbf{w}$ , independent of  $\mathbf{t}$ , to distinguish  $\mathbf{z}$  from the exogenous variables  $\mathbf{u}$ ; and **ii)** a null proxy  $\mathbf{n}$ , independent of  $\mathbf{y}$  given  $\mathbf{t}$  and  $\mathbf{z}$ , to discern the correct structural equation. Additionally, as in prior works (Miao et al., 2018; Wang & Blei, 2021),  $\mathbf{z}$  should be *complete* given the proxies (refer to

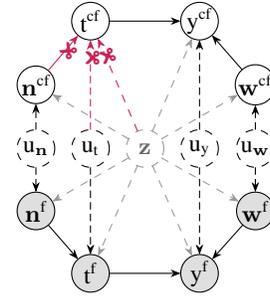


Figure 4: **Twin counterfactual network**, with observed nodes in gray. By duplicating the structural equations, we prove query identifiability in the counterfactual world while conditioning on the factual one.

Def. 5 for a formal definition). That is, both proxies should be sufficiently informative to accurately approximate  $\mathbf{z}$ .

### 6.1.2 COUNTERFACTUAL QUERIES

We focus next on the identifiability of counterfactual queries, i.e., queries of the form  $Q(\mathcal{M}) = p_{\mathcal{M}}(y^{cf} \mid \text{do}(t^{cf}), \mathbf{x}^f)$ , where  $\mathbf{x}^f$  is the observed factual, and where we are interested in *the distribution the outcome would have had, had we intervened on the treatment variable*. We demonstrate, for the first time to our knowledge, that counterfactual query identifiability holds for as many queries as for the interventional case. More specifically, we show that:

**Proposition 6.2** (Informal). *When an interventional query  $p(y \mid \text{do}(t))$  is identifiable by Decaf, then it equally identifies the counterfactual query  $p(y^{cf} \mid \text{do}(t^{cf}), \mathbf{x}^f)$ .*

The formal result can be found in Prop. A.7. In short, our result means that, if we can identify an interventional query, then we can identify its counterfactual counterpart as well.

Our result exploits the notion of twin SCM (Balke & Pearl, 1994), which duplicates the structural equations for the factual and counterfactual worlds while sharing the exogenous variables, and the fact that Prop. A.2 allows for queries with additional covariates as long as they do not form colliders, which is always the case with  $\mathbf{x}^f$  in  $p_{\mathcal{M}}(y^{cf} \mid \text{do}(t^{cf}), \mathbf{x}^f)$ , as we show in the example twin network from Fig. 4.

### 6.2 Identifying Exogenous Distributions

Besides causal query identifiability, another question of interest is whether Decaf recovers the true exogenous variables, up to component-wise transformations, disentangling the sources of variability of each endogenous variable. In App. A.1, we expand the results of Javaloy et al. (2023) to prove that Decaf identifies<sup>3</sup> the underlying SCM for those variables not directly affected by  $\mathbf{z}$ , i.e.:

**Corollary 6.3** (Informal). *Decaf identifies the underlying SCM, restricted to every variable other than  $\text{ch}(\mathbf{z})$ , up to an*

<sup>3</sup>In the sense of Xi & Bloem-Reddy (2023).

element-wise transformation of the exogenous distribution.

Moreover, we conjecture that Decaf should in most cases properly disentangle the rest of exogenous variables and  $\mathbf{z}$ . Although we do not formally prove it, we refer to the use case of §7.3 to illustrate that the exogenous variables and the latent variables extracted by Decaf. Our intuition is that, if some children of  $\mathbf{z}$  are conditionally independent, the information common to them can only be explained via  $\mathbf{z}$ . In addition, the entropy term in Eq. 7 discourages Decaf from using the components of  $\mathbf{z}$  that are not necessary for explaining the observations. Recent works proved similar results under slightly stronger assumptions (von Kügelgen et al., 2021; Zheng et al., 2022; Brady et al., 2023).

### 6.3 Practical Guidelines & Implications

In this section, we outline the different aspects to consider for the successful application of Decaf to solve causal queries in real-world scenarios.

**Training.** One key advantage of Decaf is that it needs to train only once per dataset. However, maximizing the ELBO makes it also susceptible to posterior collapse (Wang et al., 2021), i.e., to the KL term in Eq. 6 vanishing, and hence the posterior equating the prior distribution. Fortunately, we can leverage existing solutions, e.g., implement regularization terms as the one proposed by Vahdat & Kautz (2020). Recall also that, following §6, model selection should use an observational goodness-of-fit metric as selection criterion.

**Solving causal queries.** Whilst Decaf can compute *any* causal query, unidentifiable causal queries may still lead to incorrect estimates. To ensure reliability, we must verify the identifiability of each specific query of interest, for which we provide algorithms that check identifiability in the causal graph in App. E. Namely, Alg. 5 checks if a query that involves a specific treatment-outcome pair, which includes average treatment effects and counterfactuals, is identifiable. If we were interested in a query on all variables, e.g., as samples from an interventional distribution, we should evaluate the identifiability of the causal effects between the treatment and all its descendants, as proposed in Alg. 6.

**Limitations.** Decaf relaxes the assumption of *causal sufficiency*, but it still relies on completeness for the proxies, a common condition for nonparametric identification in causal inference (D’Haultfoeuille, 2011; Chen et al., 2014). This condition is untestable with observational data alone, though collecting additional proxies can help satisfy completeness (Andrews, 2011). Moreover, we assume that the true SCM is  $C^1$ -diffeomorphic with respect to the exogenous variables, which precludes theoretical guarantees for modeling discrete variables, although Javaloy et al. (2023); de Vassimon Manela et al. (2024) show that, in practice, CNFs effectively approximate discrete distributions.

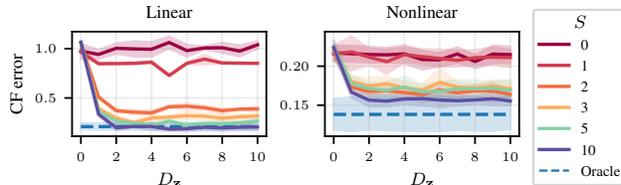


Figure 5: **Ablation.** Counterfactual error as we change the number of proxy variables and the latent dimensionality. We show means and 95 % confidence intervals over 5 realizations, intervening on the 25th, 50th, and 75th percentile of  $t$ .

## 7 Empirical Evaluation

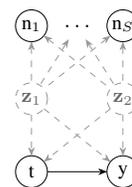
In this section, we assess the performance of Decaf comparatively to existing methods. Namely, we show that Decaf accurately estimate interventional and counterfactual queries when the requirements of Prop. 6.1 are met, and that it effectively estimates the exogenous information. We provide all experimental details in App. B.

**Common evaluation.** For all experiments, we estimate the performance on the interventional and counterfactual regimes via the mean absolute error (MAE) of, respectively, the average treatment effect (ATE) and the counterfactual samples, with respect to the ground-truth values. Moreover, we use as reference a CNF that *does observe* the hidden confounders, which we refer to as *oracle*. We also account for differences across observed variables by computing all errors over the standardized variables.

### 7.1 Ablation study

First, we conduct a simple ablation to understand how misspecifying the dimensionality of  $\mathbf{z}$  may affect Decaf, as well as its sensitivity to the number of available proxies. For additional details and results, refer to App. B.1.

**Experimental setup.** We consider two synthetic SCMs, linear and non-linear, that follow the causal graph depicted in the inset figure, comprising two independent hidden confounders affecting every variable, and  $S$  null proxies. Then, we evaluate how well Decaf estimates the direct effect of  $t$  on  $y$  while changing the number of proxy variables,  $S$ , and the specified latent dimensionality,  $D_z$ .



**Results.** Fig. 5 shows the counterfactual error for all cases, where we clearly observe that increasing the number of proxies reduces them, and with a drastic change as we add the second proxy, corroborating Prop. 6.1.

Similarly, we observe that underestimating  $D_z$  increases the error (especially if we assume causal sufficiency,  $D_z = 0$ ) while overestimating it does not. This indicates that, indeed, the entropy term in Eq. 6 prevents non-shared information from being modeled through  $\mathbf{z}$ , as discussed in §5.

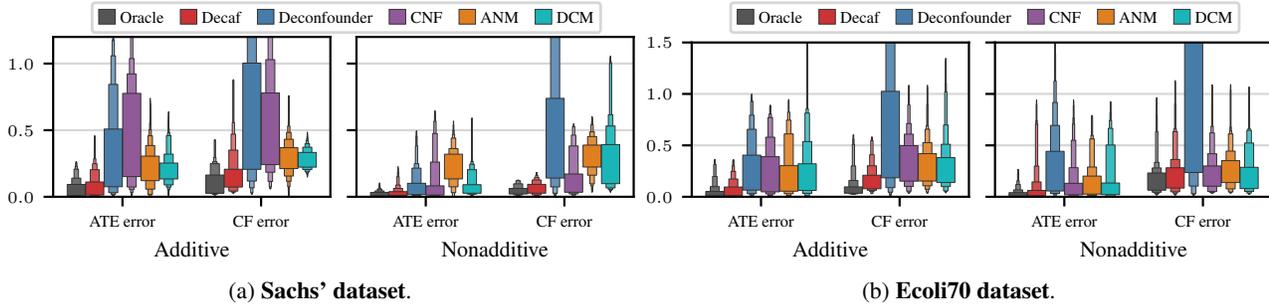


Figure 6: Error boxenplots for different CGMs, averaged over all **identifiable** direct effects of the Sachs (Fig. 7) (a) and Ecoli70 (Fig. 1) (b) datasets, after intervening in their 25th, 50th, and 75th percentiles in 5 random initializations.

### 7.2 Semi-synthetic Experiments

Next, we evaluate how Decaf performs relatively to existing approaches and, to this end, we consider semi-synthetic datasets for which we have access to the ground-truth SCMs. Additional details can be found in Apps. B.2 and B.3.

**Baselines.** We compare Decaf with three CGMs which assume causal sufficiency and are thus *unaware* of the hidden confounders: CNFs (Javaloy et al., 2023); ANMs (Hoyer et al., 2008); and DCMs (Chao et al., 2023); and with the Deconfounder (Wang & Blei, 2019), which uses proxies to provide unbiased ATE estimates under hidden confounding, yet it requires a model per treatment-outcome pair. We use the oracle as reference model to lower bound the error.

#### 7.2.1 PROTEIN-SIGNALLING NETWORKS

We first conduct a similar semi-synthetic experiment as that of Chao et al. (2023), based on a protein-signalling network dataset (Sachs et al., 2005). Specifically, we randomly generate a non-linear SCM that induces the same causal graph as the original dataset, depicted in Fig. 7, except for the root nodes, for which we use the original data. As a result, we have a hidden confounder with two dimensions, PKC and PKA, and three treatment variables to intervene upon, Raf, Mek, and Erk. We consider additive and non-additive structural equations, measure the effect of interventions on the downstream nodes and, more importantly, ensure that the randomized effect of the hidden confounder is perceptible.

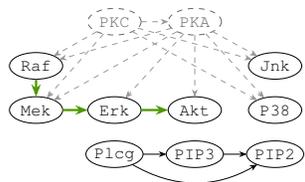


Figure 7: **Sachs' causal graph.** Green denotes **identifiable** confounded effects.

**Results.** We present a summary of the results in Fig. 6a, where we can observe that Decaf outperforms every approach in all cases, for both ATE and counterfactual errors, remaining fairly close to the oracle model. Moreover, we appreciate a great difference in performance between Decaf and CNFs, which corroborates the importance of the pro-

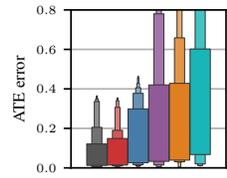
posed encoder and variational training employed by Decaf, since a CNF is equivalent to Decaf with  $D_z = 0$ .

#### 7.2.2 GENE NETWORKS

Next, we repeat a similar experiment as in the previous section, considering this time the causal graph of the Ecoli70 dataset (Schäfer & Strimmer, 2005) as reference, shown in Fig. 1, representing a gene network from E. coli data. This time, we replace root nodes with Gaussian variables.

**Results.** Similar to the previous case, the results presented in Fig. 6b demonstrate that Decaf is indeed able to closely match the performance of the oracle model, outperforming existing approaches. However, the non-additive case also shows significant long-tailed error distributions for all models, showing that Decaf can suffer the same problems as any data-centric approach, and that it is still needed to put attention on its effective training.

It is also worth-pointing out that the striking performance of the Deconfounder is a result of evaluating causal queries that cannot be identified by the model. As we discuss in App. B, the Deconfounder offers guarantees regarding ATE estimation and with more restrictive assumptions. If we plot instead the ATE error evaluated on only those paths that meet the assumptions placed by the Deconfounder, as shown in the inset figure, we see that it now achieves significantly lower errors that the *unaware* approaches.



Remarkably, this experiment highlights every strength of the proposed approach, since Decaf: **i)** models several hidden confounders affecting different sets of variables; **ii)** identifies all causal queries for which we have some proxy information; and **iii)** achieves the above in an agnostic manner, i.e., training out-of-the-box and *one single time*, despite the graph having 43 observed variables.

### 7.3 Fairness Real-world Use Case

Taking inspiration from the experiments by Kusner et al. (2017) and Javaloy et al. (2023), we aim to show how model-

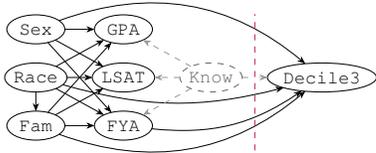


Figure 8: Causal graph assumed for the law school dataset. Decile3 is only used by the classifiers.

Table 1: Test RMSE in Decile3 prediction and MMD comparing different-group prediction distributions.

	Unfair	Unaware	Decaf	Fair $K$	Fair Add	Mean
RMSE	1.477	1.479	1.652	2.818	2.817	2.83
MMD	0.110	0.102	0.0018	$10^{-6}$	$10^{-8}$	0

ling confounded SCMs with Decaf can be leveraged beyond causal query estimation and, in particular, for counterfactual fairness prediction. See App. B.4 for further details.

**Dataset and objective.** Our aim is to train a predictor, using the law school dataset (Wightman, 1998) which comprises information of 21 790 law students, that remains accurate while being fair—using *demographic parity* as fairness criterion (Feldman et al., 2015)—toward the sensitive attributes of the students. In particular, we are interested in predicting the decile of a student in its 3rd year of university, given their undergraduate and 1st year grades, family income, race, and sex.

**Experimental setup.** First, we train Decaf assuming a causal graph such as the one in Fig. 8, excluding Decile3, where all grades are affected by a common “knowledge” hidden confounder. Then, we train a simple predictor using as input the hidden confounder and non-sensitive exogenous variables estimated by Decaf. If, as discussed in §6.2, Decaf successfully recovers the exogenous variables, we expect the predictor to be fair yet slightly less accurate, since Decile3 is directly affected by the sensitive attributes.

**Results.** Tab 1 shows the prediction error (RMSE) and the difference between groups (MMD) for the proposed predictor using Decaf, comparing with an *unfair* predictor that uses sensitive attributes; an *unaware* predictor that excludes sensitive attributes, and two fair predictors—*Fair K* and *Fair Add*—proposed by Kusner et al. (2017).

As shown in Fig. 9, Decaf provides a much fairer predictor than the *unfair* and the *unaware* predictors at the cost of slightly higher RMSE. We can also appreciate that the other two fair approaches are so by predicting a constant value for every individual, which can be also observed comparing the RMSE obtained by these predictors with a naive predictor that predicts the mean of the distribution in Tab 1.

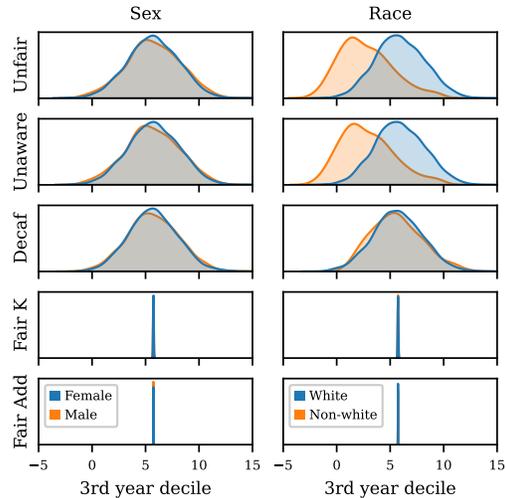


Figure 9: **Distribution of predicted Decile3.** Fairer predictors yield similar distributions across the two considered groups on each attribute (Sex and Race).

## 8 Concluding Remarks

In this work, we have bridged the current gap between CGMs, which fail to account for hidden confounders, and hidden-confounding solutions, which are tailored to a specific causal query and thus need to train once per query. To this end, we have introduced Decaf, and theoretically shown that it can accurately estimate causal queries in the presence of hidden confounders, if there exists a valid adjustment set or sufficiently informative proxies, extending prior results (Miao et al., 2018) to also consider counterfactuals. We have empirically shown that Decaf outperforms all considered baselines, better estimating confounded causal queries shown to be identifiable, and properly identifying exogenous distributions to train fair classifiers. Finally, we have provided algorithms to check the identifiability of causal queries which, along Decaf, provides practitioners with a powerful pipeline to perform causal inference in the presence of hidden confounders.

**Future work.** Our work opens many intriguing venues, e.g., integrating alternative identification strategies, such as instrumental variables (Hartford et al., 2017), to expand the range of identifiable queries that Decaf can estimate. We also find it interesting to apply Decaf to settings with time-varying treatments, where multiple interventions have to be performed. In real-world scenarios, it would be exciting to include interventional data during training, and seeing Decaf applied to real-world problems such as decision support systems (Sanchez et al., 2022), educational analysis (Murnane, 2010), or policy making (Fougère & Jacquemet, 2021), yet always validating them with interventional data.

## Impact statement

This research contributes to advance causal inference in machine learning, particularly enhancing the ability to estimate causal effects despite unobserved variables. Thus, this work supports more informed decision-making in scenarios where controlled experimentation is impractical or unethical, such as healthcare or education. As with all advances in causal inference, practitioners should be aware of the limitations and assumptions of causal models. Particularly, in sensitive applications, where decisions are based on accurate causal conclusions, validation with interventional data should be prioritized whenever possible to ensure reliability. Overall, this work aligns with the broader goal of improving machine learning and does not introduce significant ethical risk beyond those traditionally associated with the field.

## BIBLIOGRAPHY

- Allman, E. S., Matias, C., and Rhodes, J. A. Identifiability of parameters in latent structure models with many observed variables. *The Annals of Statistics*, 37(6A):3099–3132, 2009. ISSN 00905364, 21688966. URL <http://www.jstor.org/stable/25662188>. (Cited in pages 2, 29, and 30.)
- Andrews, D. W. Examples of L2-complete and boundedly-complete distributions. 2011. (Cited in pages 6 and 15.)
- Angrist, J. D. and Pischke, J.-S. *Mostly harmless econometrics: An empiricist’s companion*. Princeton university press, 2009. (Cited in pages 2 and 29.)
- Balke, A. and Pearl, J. Probabilistic Evaluation of Counterfactual Queries. *Probabilistic and Causal Inference*, 1994. URL <https://api.semanticscholar.org/CorpusID:18845266>. (Cited in page 5.)
- Bica, I., Alaa, A. M., and van der Schaar, M. Time Series Deconfounder: Estimating Treatment Effects over Time in the presence of hidden confounders. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 884–895. PMLR, 2020. URL <http://proceedings.mlr.press/v119/bica20a.html>. (Cited in page 3.)
- Bingham, E., Chen, J. P., Jankowiak, M., Obermeyer, F., Pradhan, N., Karaletsos, T., Singh, R., Szerlip, P., Horsfall, P., and Goodman, N. D. Pyro: Deep universal probabilistic programming. *Journal of machine learning research*, 20(28):1–6, 2019. (Cited in page 27.)
- Brady, J., Zimmermann, R. S., Sharma, Y., Schölkopf, B., von Kügelgen, J., and Brendel, W. Provably Learning Object-Centric representations. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 3038–3062. PMLR, 2023. URL <https://proceedings.mlr.press/v202/brady23a.html>. (Cited in page 6.)
- Carrasco, M., Florens, J.-P., and Renault, E. Linear inverse problems in structural econometrics estimation based on spectral decomposition and regularization. *Handbook of econometrics*, 6:5633–5751, 2007. (Cited in page 16.)
- Chao, P., Blöbaum, P., and Kasiviswanathan, S. P. Interventional and counterfactual inference with diffusion models. *ArXiv preprint*, abs/2302.00860, 2023. URL <https://arxiv.org/abs/2302.00860>. (Cited in pages 1, 2, 7, and 22.)
- Chen, X., Chernozhukov, V., Lee, S., and Newey, W. K. Local identification of nonparametric and semiparametric models. *Econometrica*, 82(2):785–809, 2014. (Cited in page 6.)
- de Vassimon Manela, D., Battaglia, L., and Evans, R. J. Marginal Causal Flows for Validation and Inference. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. (Cited in page 6.)
- D’Haultfoeuille, X. On the completeness condition in non-parametric instrumental problems. *Econometric Theory*, 27(3):460–471, 2011. (Cited in page 6.)
- Feldman, M., Friedler, S. A., Moeller, J., Scheidegger, C., and Venkatasubramanian, S. Certifying and Removing Disparate Impact. In Cao, L., Zhang, C., Joachims, T., Webb, G. I., Margineantu, D. D., and Williams, G. (eds.), *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, pp. 259–268. ACM, 2015. doi: 10.1145/2783258.2783311. URL <https://doi.org/10.1145/2783258.2783311>. (Cited in page 8.)
- Feuerriegel, S., Frauen, D., Melnychuk, V., Schweisthal, J., Hess, K., Curth, A., Bauer, S., Kilbertus, N., Kohane, I. S., and van der Schaar, M. Causal machine learning for predicting treatment outcomes. *Nature Medicine*, 30(4): 958–968, 2024. (Cited in page 1.)
- Fougère, D. and Jacquemet, N. Policy evaluation using causal inference methods. In *Handbook of Research Methods and Applications in Empirical Microeconomics*, pp. 294–324. Edward Elgar Publishing, 2021. (Cited in page 8.)
- Hartford, J. S., Lewis, G., Leyton-Brown, K., and Taddy, M. Deep IV: A Flexible Approach for Counterfactual prediction. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of*

- 495 *the 34th International Conference on Machine Learning,*  
 496 *ICML 2017, Sydney, NSW, Australia, 6-11 August 2017,*  
 497 *volume 70 of Proceedings of Machine Learning Research,*  
 498 *pp. 1414–1423. PMLR, 2017. URL <http://proceedings.mlr.press/v70/hartford17a.html>.*  
 499 *(Cited in pages 8 and 29.)*  
 500
- 501 Hoyer, P. O., Janzing, D., Mooij, J. M., Peters, J., and  
 502 Schölkopf, B. Nonlinear causal discovery with additive  
 503 noise models. In Koller, D., Schuurmans, D., Bengio,  
 504 Y., and Bottou, L. (eds.), *Advances in Neural Informa-*  
 505 *tion Processing Systems 21, Proceedings of the Twenty-*  
 506 *Second Annual Conference on Neural Information Pro-*  
 507 *cessing Systems, Vancouver, British Columbia, Canada,*  
 508 *December 8-11, 2008*, pp. 689–696. Curran Associates,  
 509 Inc., 2008. URL <https://proceedings.neurips.cc/paper/2008/hash/f7664060cc52bc6f3d620bcedc94a4b6-Abstract.html>. (Cited  
 510 in pages 2 and 7.)  
 511
- 512 Javaloy, A., Sánchez-Martín, P., and Valera, I. Causal nor-  
 513 malizing flows: from theory to practice. In Oh, A., Nau-  
 514 mann, T., Globerson, A., Saenko, K., Hardt, M., and  
 515 Levine, S. (eds.), *Advances in Neural Information Pro-*  
 516 *cessing Systems 36: Annual Conference on Neural In-*  
 517 *formation Processing Systems 2023, NeurIPS 2023, New*  
 518 *Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL  
 519 [http://papers.nips.cc/paper\\_files/paper/2023/hash/b8402301e7f06bdc97a31bfaa653dc32-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/b8402301e7f06bdc97a31bfaa653dc32-Abstract-Conference.html).  
 520 (Cited in pages 1, 2, 3, 4, 5, 6, 7, 15, 18, 19, 25, 27,  
 521 and 28.)  
 522
- 523 Kallus, N., Mao, X., and Udell, M. Causal Inference with  
 524 Noisy and Missing Covariates via Matrix factorization.  
 525 In Bengio, S., Wallach, H. M., Larochelle, H., Grauman,  
 526 K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances*  
 527 *in Neural Information Processing Systems 31: Annual*  
 528 *Conference on Neural Information Processing Systems*  
 529 *2018, NeurIPS 2018, December 3-8, 2018, Montréal,*  
 530 *Canada*, pp. 6921–6932, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/86a1793f65aeef4aeef4b479fc9b2bca-Abstract.html>. (Cited in pages 2, 29, and 30.)  
 531
- 532 Kallus, N., Mao, X., and Zhou, A. Interval Estimation of  
 533 Individual-Level Causal Effects Under Unobserved con-  
 534 founding. In Chaudhuri, K. and Sugiyama, M. (eds.), *The*  
 535 *22nd International Conference on Artificial Intelligence*  
 536 *and Statistics, AISTATS 2019, 16-18 April 2019, Naha,*  
 537 *Okinawa, Japan*, volume 89 of *Proceedings of Machine*  
 538 *Learning Research*, pp. 2281–2290. PMLR, 2019. URL  
 539 <http://proceedings.mlr.press/v89/kallus19a.html>. (Cited in pages 29 and 30.)  
 540
- 541 Kaltenpoth, D. and Vreeken, J. Nonlinear Causal Discovery  
 542 with Latent confounders. In Krause, A., Brunskill, E.,  
 543 Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.),  
 544 *International Conference on Machine Learning, ICML*  
 545 *2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume  
 546 202 of *Proceedings of Machine Learning Research*, pp.  
 547 15639–15654. PMLR, 2023. URL <https://proceedings.mlr.press/v202/kaltenpoth23a.html>. (Cited in page 22.)  
 548
- 549 Khemakhem, I., Monti, R. P., Leech, R., and Hyvärinen, A.  
 Causal Autoregressive flows. In Banerjee, A. and Fuku-  
 mizu, K. (eds.), *The 24th International Conference on*  
*Artificial Intelligence and Statistics, AISTATS 2021, April*  
*13-15, 2021, Virtual Event*, volume 130 of *Proceedings*  
*of Machine Learning Research*, pp. 3520–3528. PMLR,  
 2021. URL <http://proceedings.mlr.press/v130/khemakhem21a.html>. (Cited in pages 1 and 2.)
- Kingma, D. P. and Welling, M. Auto-Encoding Vari-  
 ational Bayes. In Bengio, Y. and LeCun, Y. (eds.),  
*2nd International Conference on Learning Representa-*  
*tions, ICLR 2014, Banff, AB, Canada, April 14-16,*  
*2014, Conference Track Proceedings*, 2014. URL <http://arxiv.org/abs/1312.6114>. (Cited in pages 1,  
 3, and 4.)
- Kocaoglu, M., Snyder, C., Dimakis, A. G., and Vishwanath,  
 S. Causalgan: Learning Causal Implicit Generative Mod-  
 els with Adversarial training. In *6th International Con-*  
*ference on Learning Representations, ICLR 2018, Van-*  
*couver, BC, Canada, April 30 - May 3, 2018, Conference*  
*Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=BJE-4xW0W>.  
 (Cited in page 2.)
- Kolmogorov, A. On the Shannon theory of information  
 transmission in the case of continuous signals. *IRE*  
*Transactions on Information Theory*, 2(4):102–108, 1956.  
 (Cited in page 4.)
- Kullback, S. and Leibler, R. A. On information and suf-  
 ficiency. *The annals of mathematical statistics*, 22(1):  
 79–86, 1951. (Cited in page 4.)
- Kuroki, M. and Pearl, J. Measurement bias and effect res-  
 toration in causal inference. *Biometrika*, 101(2):423–437,  
 2014. (Cited in pages 2, 29, and 30.)
- Kusner, M. J., Loftus, J. R., Russell, C., and Silva, R. Coun-  
 terfactual fairness. In Guyon, I., von Luxburg, U., Bengio,  
 S., Wallach, H. M., Fergus, R., Vishwanathan, S. V. N.,  
 and Garnett, R. (eds.), *Advances in Neural Information*  
*Processing Systems 30: Annual Conference on Neural*  
*Information Processing Systems 2017, December 4-9,*  
*2017, Long Beach, CA, USA*, pp. 4066–4076, 2017. URL  
<https://proceedings.neurips.cc/paper>

- 550 /2017/hash/a486cd07e4ac3d270571622f4  
551 f316ec5-Abstract.html. (Cited in pages 7, 8, 25,  
552 26, and 27.)
- 553 Long, C. P. and Antoniewicz, M. R. Metabolic flux analysis  
554 of *Escherichia coli* knockouts: lessons from the  
555 Keio collection and future outlook. *Current opinion in  
556 biotechnology*, 28:127–133, 2014. (Cited in page 23.)
- 558 Louizos, C., Shalit, U., Mooij, J. M., Sontag, D. A., Zemel,  
559 R. S., and Welling, M. Causal Effect Inference with Deep  
560 Latent-Variable models. In Guyon, I., von Luxburg, U.,  
561 Bengio, S., Wallach, H. M., Fergus, R., Vishwanathan,  
562 S. V. N., and Garnett, R. (eds.), *Advances in Neural In-  
563 formation Processing Systems 30: Annual Conference on  
564 Neural Information Processing Systems 2017, December  
565 4-9, 2017, Long Beach, CA, USA*, pp. 6446–6456, 2017.  
566 URL [https://proceedings.neurips.cc/p  
567 aper/2017/hash/94b5bde6de888ddf9cde6  
568 748ad2523d1-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/94b5bde6de888ddf9cde6748ad2523d1-Abstract.html). (Cited in pages 2,  
569 29, and 30.)
- 570 Luo, R. and Zhao, H. Bayesian hierarchical modeling for  
571 signaling pathway inference from single cell interven-  
572 tional data. *The annals of applied statistics*, 5:725–745,  
573 2011. doi: 10.1214/10-AOAS425. (Cited in page 22.)
- 575 Miao, W., Geng, Z., and Tchetgen Tchetgen, E. J. Identifying  
576 causal effects with proxy variables of an unmeasured  
577 confounder. *Biometrika*, 105(4):987–993, 2018. (Cited  
578 in pages 1, 2, 5, 8, 16, 17, 18, 29, and 30.)
- 580 Miao, W., Hu, W., Ogburn, E. L., and Zhou, X.-H. Identifying  
581 effects of multiple treatments in the presence of  
582 unmeasured confounding. *Journal of the American Stat-  
583 istical Association*, 118(543):1953–1967, 2023. (Cited in  
584 pages 1, 2, 15, 29, and 30.)
- 585 Murnane, R. *Methods matter: Improving causal inference  
586 in educational and social science research*. Oxford Uni-  
587 versity Press, 2010. (Cited in page 8.)
- 589 Nasr-Esfahany, A., Alizadeh, M., and Shah, D. Counterfac-  
590 tual Identifiability of Bijective Causal models. In Krause,  
591 A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S.,  
592 and Scarlett, J. (eds.), *International Conference on Ma-  
593 chine Learning, ICML 2023, 23-29 July 2023, Honolulu,  
594 Hawaii, USA*, volume 202 of *Proceedings of Machine  
595 Learning Research*, pp. 25733–25754. PMLR, 2023. URL  
596 [https://proceedings.mlr.press/v202/n  
597 asr-esfahany23a.html](https://proceedings.mlr.press/v202/nasr-esfahany23a.html). (Cited in page 15.)
- 598 Papamakarios, G., Nalisnick, E. T., Rezende, D. J., Mo-  
599 hamed, S., and Lakshminarayanan, B. Normalizing Flows  
600 for Probabilistic Modeling and inference. *J. Mach. Learn.  
601 Res.*, 22:57:1–57:64, 2021. URL [http://jmlr.org  
602 /papers/v22/19-1028.html](http://jmlr.org/papers/v22/19-1028.html). (Cited in pages 3  
603 and 4.)
- 604 Parafita, Á. and Vitrià, J. Estimand-agnostic causal query  
estimation with deep causal graphs. *IEEE Access*, 10:  
71370–71386, 2022. (Cited in page 2.)
- Pawlowski, N., de Castro, D. C., and Glocker, B. Deep  
Structural Causal Models for Tractable Counterfactual  
inference. In Larochelle, H., Ranzato, M., Hadsell, R.,  
Balcan, M., and Lin, H. (eds.), *Advances in Neural In-  
formation Processing Systems 33: Annual Conference  
on Neural Information Processing Systems 2020, Neur-  
IPS 2020, December 6-12, 2020, virtual*, 2020. URL  
[https://proceedings.neurips.cc/paper  
/2020/hash/0987b8b338d6c90bbedd8631b  
c499221-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/0987b8b338d6c90bbedd8631bc499221-Abstract.html). (Cited in page 2.)
- Pearl, J. *Causality*. Cambridge university press, 2009.  
(Cited in pages 2, 4, 20, and 29.)
- Pearl, J. The Do-Calculus Revisited. In de Freitas, N.  
and Murphy, K. P. (eds.), *Proceedings of the Twenty-  
Eighth Conference on Uncertainty in Artificial Intelli-  
gence, Catalina Island, CA, USA, August 14-18, 2012*, pp.  
3–11. AUAI Press, 2012. URL [https://dslpitt.  
org/uai/displayArticleDetails.jsp?mm  
nu=1&smnu=2&article\\_id=2330&proceedi  
ng\\_id=28](https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=2330&proceeding_id=28). (Cited in page 2.)
- Pearl, J., Glymour, M., and Jewell, N. P. *Causal inference  
in statistics: A primer*. John Wiley & Sons, 2016. (Cited  
in page 28.)
- Peters, J., Janzing, D., and Schölkopf, B. *Elements of causal  
inference: foundations and learning algorithms*. The MIT  
Press, 2017. (Cited in pages 3, 4, 16, 17, 19, and 20.)
- Ranganath, R. and Perotte, A. Multiple causal inference  
with latent confounding. *ArXiv preprint*, abs/1805.08273,  
2018. URL [https://arxiv.org/abs/1805.0  
8273](https://arxiv.org/abs/1805.08273). (Cited in page 30.)
- Sachs, K., Perez, O., Pe’er, D., Lauffenburger, D. A., and  
Nolan, G. P. Causal Protein-Signaling Networks Derived  
from Multiparameter Single-Cell Data. *Science*, 308  
(5721):523–529, 2005. doi: 10.1126/science.1105809.  
URL [https://www.science.org/doi/abs/  
10.1126/science.1105809](https://www.science.org/doi/abs/10.1126/science.1105809). (Cited in pages 7  
and 22.)
- Sanchez, P., Voisey, J. P., Xia, T., Watson, H. I., O’Neil,  
A. Q., and Tsafaris, S. A. Causal machine learning for  
healthcare and precision medicine. *Royal Society Open  
Science*, 9(8):220638, 2022. (Cited in page 8.)
- Sánchez-Martín, P., Rateike, M., and Valera, I. VACA:  
Designing Variational Graph Autoencoders for Causal  
queries. In *Thirty-Sixth AAAI Conference on Artificial  
Intelligence, AAAI 2022, Thirty-Fourth Conference on In-  
novative Applications of Artificial Intelligence, IAAI 2022,*

- 605 *The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February*  
606 *22 - March 1, 2022*, pp. 8159–8168. AAAI Press, 2022.  
607 URL <https://ojs.aaai.org/index.php/AAI/article/view/20789>. (Cited in page 2.)
- 610 Schäfer, J. and Strimmer, K. A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Statistical applications in genetics and molecular biology*, 4(1), 2005. (Cited in pages 1, 2, 7, and 23.)
- 611 Scutari, M. Learning Bayesian Networks with the bnlearn R Package. *Journal of Statistical Software*, 35(3):1–22, 2010. doi: 10.18637/jss.v035.i03. (Cited in pages 22 and 23.)
- 612 Spirtes, P., Glymour, C., and Scheines, R. *Causation, prediction, and search*. MIT press, 2001. (Cited in page 17.)
- 613 Tchetgen, E. J. T., Ying, A., Cui, Y., Shi, X., and Miao, W. An introduction to proximal causal learning. *ArXiv preprint*, abs/2009.10982, 2020. URL <https://arxiv.org/abs/2009.10982>. (Cited in page 29.)
- 614 Vahdat, A. and Kautz, J. NVAE: A Deep Hierarchical Variational autoencoder. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/e3b21256183cf7c2c7a66be163579d37-Abstract.html>. (Cited in page 6.)
- 615 Varian, H. R. Causal inference in economics and marketing. *Proceedings of the National Academy of Sciences*, 113(27):7310–7315, 2016. (Cited in page 1.)
- 616 von Kügelgen, J., Sharma, Y., Gresele, L., Brendel, W., Schölkopf, B., Besserve, M., and Locatello, F. Self-Supervised Learning with Data Augmentations Provably Isolates content from style. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 16451–16467, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/8929c70f8d710e412d38da624b21c3c8-Abstract.html>. (Cited in page 6.)
- 617 Wang, Y. and Blei, D. M. The blessings of multiple causes. *Journal of the American Statistical Association*, 114(528): 1574–1596, 2019. (Cited in pages 3, 7, 23, 24, and 30.)
- 618 Wang, Y. and Blei, D. M. A Proxy Variable View of Shared confounding. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10697–10707. PMLR, 2021. URL <http://proceedings.mlr.press/v139/wang21c.html>. (Cited in pages 1, 2, 3, 5, 15, 16, 18, 23, and 30.)
- 619 Wang, Y., Blei, D. M., and Cunningham, J. P. Posterior Collapse and Latent Variable Non-identifiability. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 5443–5455, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/2b6921f2c64dee16ba21ebf17f3c2c92-Abstract.html>. (Cited in page 6.)
- 620 Wightman, L. F. Lsac National Longitudinal Bar Passage Study. LSAC Research Report Series. 1998. (Cited in page 8.)
- 621 Winkler, C., Worrall, D., Hoogeboom, E., and Welling, M. Learning Likelihoods with Conditional Normalizing flows. *ArXiv preprint*, abs/1912.00042, 2019. URL <https://arxiv.org/abs/1912.00042>. (Cited in pages 2, 3, and 4.)
- 622 Xi, Q. and Bloem-Reddy, B. Indeterminacy in Generative Models: Characterization and Strong identifiability. In Ruiz, F. J. R., Dy, J. G., and van de Meent, J. (eds.), *International Conference on Artificial Intelligence and Statistics, 25-27 April 2023, Palau de Congressos, Valencia, Spain*, volume 206 of *Proceedings of Machine Learning Research*, pp. 6912–6939. PMLR, 2023. URL <https://proceedings.mlr.press/v206/xi23a.html>. (Cited in page 5.)
- 623 Yang, M., Liu, F., Chen, Z., Shen, X., Hao, J., and Wang, J. Causalvae: Structured causal disentanglement in variational autoencoder. *ArXiv preprint*, abs/2004.08697, 2020. URL <https://arxiv.org/abs/2004.08697>. (Cited in page 2.)
- 624 Zečević, M., Dhimi, D. S., Velivcković, P., and Kersting, K. Relating graph neural networks to structural causal models. *ArXiv preprint*, abs/2109.04173, 2021. URL <https://arxiv.org/abs/2109.04173>. (Cited in page 2.)
- 625 Zhao, S. and Heffernan, N. Estimating Individual Treatment Effect from Educational Studies with Residual Counterfactual Networks. *International Educational Data Mining Society*, 2017. (Cited in page 1.)

660 Zheng, Y., Ng, I., and Zhang, K. On the Identifiability  
661 of Nonlinear ICA: Sparsity and beyond. In Koyejo, S.,  
662 Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and  
663 Oh, A. (eds.), *Advances in Neural Information Processing*  
664 *Systems 35: Annual Conference on Neural Information*  
665 *Processing Systems 2022, NeurIPS 2022, New Orleans,*  
666 *LA, USA, November 28 - December 9, 2022, 2022.* URL  
667 [http://papers.nips.cc/paper\\_files/p](http://papers.nips.cc/paper_files/paper/2022/hash/6801fa3fd290229efc490ee0cf1c5687-Abstract-Conference.html)  
668 [aper/2022/hash/6801fa3fd290229efc490](http://papers.nips.cc/paper_files/paper/2022/hash/6801fa3fd290229efc490ee0cf1c5687-Abstract-Conference.html)  
669 [ee0cf1c5687-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/6801fa3fd290229efc490ee0cf1c5687-Abstract-Conference.html).  
670 (Cited in page 6.)  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
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

# Appendix

## Table of Contents

---

<b>A Causal identifiability</b>	<b>15</b>
A.1 Model identifiability . . . . .	15
A.2 Query identifiability . . . . .	15
A.3 Counterfactual query identifiability . . . . .	20
<b>B Experimental details and additional results</b>	<b>21</b>
B.1 Ablation study . . . . .	21
B.2 Semi-synthetic Sachs' dataset . . . . .	22
B.3 Semi-synthetic Ecoli70 dataset . . . . .	23
B.4 Law school fairness use-case . . . . .	25
<b>C Do-operator</b>	<b>27</b>
C.1 Do-operator in causal normalizing flows . . . . .	27
C.2 Do-operator in interventional distributions with Decaf . . . . .	28
C.3 Do-operator in counterfactuals with Decaf . . . . .	28
<b>D Additional details on related work of causal inference with hidden confounders</b>	<b>29</b>
<b>E Algorithms for causal query identification</b>	<b>30</b>
E.1 Pipeline for using Decaf . . . . .	31

---

## A Causal identifiability

### A.1 Model identifiability

In this section, we briefly discuss the identifiability of those variables that are indirectly confounded by  $\mathbf{z}$  or not confounded at all, i.e., of those variables that are not children of any hidden confounder. As we discuss now, we can reduce our SCM to a conditional SCM that only models these variables, recovering the identifiability guarantees from Javaloy et al. (2023).

To prove model identifiability, we resort to what we call the induced conditional SCM, which intuitively represents the original SCM where we restrict our view to these variables, and assume the rest of the variables are given.

**Definition 4** (Induced conditional SCM). Given a SCM  $\mathcal{M} = (\mathbf{f}, P_{\mathbf{u}}, P_{\mathbf{z}})$ , and a subset of observed variables  $\mathbf{x}' \subset \mathbf{x}$ , we define the *induced conditional SCM of  $\mathcal{M}$  given  $\mathbf{x}'$* , denoted by  $\mathcal{M}_{|\mathbf{x}'}$ , to the SCM result of having observed  $\mathbf{x}'$ , and where causal generators and exogenous variables are restricted to only those components associated with the rest of variables.

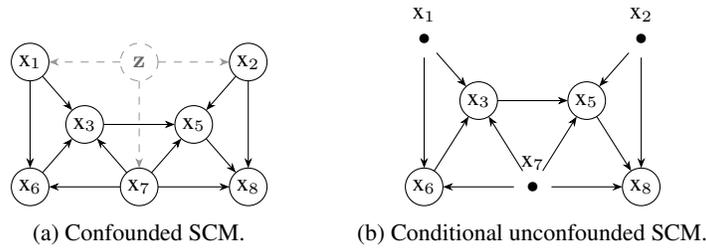


Figure 10: Example of: (a) a confounded SCM  $\mathcal{M}$ ; and (b) its induced conditional counterpart,  $\mathcal{M}_{|\mathbf{x}'}$ , when the children of the hidden confounder are observed and fixed. Note that  $\mathcal{M}_{|\mathbf{x}'}$  has no hidden confounding.

We provide a visual depiction of this idea in Fig. 10. Using this definition, we can observe that, if we were to condition of the children of the hidden confounder, we would be left with a (conditional) unconfounded SCM, as the influence of the hidden confounder has been completely blocked by conditioning on its children. Now, if we have two models that perfectly match their marginal distributions, this means that they perfectly match their induced conditional SCM, no matter which value we observed for  $\text{ch}(\mathbf{z})$ , and we can thus leverage existing results from Javaloy et al. (2023) for unconfounded SCMs.

**Corollary A.1.** Assume that we have two SCMs  $\mathcal{M} := (\mathbf{f}, P_{\mathbf{u}}, P_{\mathbf{z}})$  and  $\tilde{\mathcal{M}} := (\tilde{\mathbf{f}}, P_{\tilde{\mathbf{u}}}, P_{\tilde{\mathbf{z}}})$  that are compatible, i.e., they induce the same causal graph, and which coincide in their marginal distributions,  $p(\mathbf{x}) = \tilde{p}(\mathbf{x})$ . Then, both SCMs, restricted to every variable other than  $\text{ch}(\mathbf{z})$ , are equal up to an element-wise transformation of the exogenous distributions.

*Proof.* The proof follows almost directly from (Javaloy et al., 2023, Theorem 1). First, note that the two induced conditional SCMs are no longer influenced by  $\mathbf{z}$  once that we have observed a specific realization of  $\text{ch}(\mathbf{z})$ , so that we can drop  $\mathbf{z}$  from their structure, i.e., we can denote them by  $\mathcal{M}_{|\text{ch}(\mathbf{z})} = (\mathbf{f}_{|\text{ch}(\mathbf{z})}, P_{\mathbf{u}_{|\text{ch}(\mathbf{z})}})$  and  $\tilde{\mathcal{M}}_{|\text{ch}(\mathbf{z})} = (\tilde{\mathbf{f}}_{|\text{ch}(\mathbf{z})}, P_{\tilde{\mathbf{u}}_{|\text{ch}(\mathbf{z})}})$ . To ease notation, let us call  $\mathbf{x}^{\text{c}} := \mathbf{x} \setminus \text{ch}(\mathbf{z})$  the variables that are not children of  $\mathbf{z}$ .

Next, note that for almost every realization of  $\text{ch}(\mathbf{z})$ , we have that  $p(\mathbf{x}^{\text{c}} | \text{ch}(\mathbf{z})) = \tilde{p}(\mathbf{x}^{\text{c}} | \text{ch}(\mathbf{z}))$  since  $p(\mathbf{x}) = \tilde{p}(\mathbf{x})$  by assumption and  $p(\mathbf{x}) = p(\mathbf{x}^{\text{c}} | \text{ch}(\mathbf{z}))p(\text{ch}(\mathbf{z}))$ . As a result, for each realization of  $\text{ch}(\mathbf{z})$  we can apply Theorem 1 of Javaloy et al. (2023), which yields that the two induced conditional SCMs are equal up to an element-wise transformation of the exogenous distribution.

Finally, since the causal generators and exogenous distributions of the induced SCMs are, for almost every  $\text{ch}(\mathbf{z})$ , identical to their counterparts in the original SCMs (as we have just discarded those components associated with  $\text{ch}(\mathbf{z})$ ), we get that the elements in the two SCMs associated with every variable except those in  $\text{ch}(\mathbf{z})$  are identical up to said (possibly  $\text{ch}(\mathbf{z})$ -dependent) transformation.  $\square$

### A.2 Query identifiability

We now prove the identifiability of the causal queries considered in the main text.

To this end, one key property that we will use in the following is that of completeness (see, e.g., the work of Wang & Blei (2021)). Intuitively, we say that a random variable  $\mathbf{z}$  is complete given another random variable  $\mathbf{n}$  if “any infinitesimal change in  $\mathbf{z}$  is accompanied by variability in  $\mathbf{n}$ ” (Miao et al., 2023), yielding enough information to recover the posterior distribution of  $\mathbf{z}$ . This concept is similar in spirit to that of variability in the case of discrete random variables (Nasr-Esfahany et al., 2023). In practice, completeness is more likely to be achieved the more proxies we measure (Andrews, 2011).

**Definition 5** (Completeness). We say that a random variable  $\mathbf{z}$  is complete given  $\mathbf{n}$  for all  $\mathbf{c}$  if, for any square-integrable function  $g(\cdot)$  and almost all  $\mathbf{c}$ ,  $\int g(\mathbf{z}, \mathbf{c})p(\mathbf{z} | \mathbf{c}, \mathbf{n}) d\mathbf{z} = 0$  for almost all  $\mathbf{n}$ , if and only if  $g(\mathbf{z}, \mathbf{c}) = 0$  for almost all  $\mathbf{z}$ .

The following proposition is a generalization of the results previously presented by Miao et al. (2018) and Wang & Blei (2021), where we include an additional covariate  $\mathbf{c}$  to the causal query, and make no implicit assumptions on the causal graph allowing, e.g., for the treatment and outcome variables to share some observed parents. However, note that  $\mathbf{c}$  cannot be a collider (e.g., forming a subgraph of the form  $\mathbf{n} \rightarrow \mathbf{c} \leftarrow \mathbf{y}$ ) as, otherwise, conditioning on it would make independent variables dependent (in the example,  $\mathbf{y}$  and  $\mathbf{n}$ ), and the causal effect of  $\mathbf{t}$  on  $\mathbf{y}$  would not be identifiable (Peters et al., 2017).

**Proposition A.2** (Query identifiability). Assume that we have two SCMs  $\mathcal{M} := (\mathbf{f}, P_{\mathbf{u}}, P_{\mathbf{z}})$  and  $\tilde{\mathcal{M}} := (\tilde{\mathbf{f}}, P_{\tilde{\mathbf{u}}}, P_{\tilde{\mathbf{z}}})$  that are compatible, i.e., they induce the same causal graph, and which coincide in their marginal distributions,  $p(\mathbf{x}) = \tilde{p}(\mathbf{x})$ . Then, they compute the same causal query,  $p(\mathbf{y} | do(\mathbf{t}), \mathbf{c}) = \tilde{p}(\mathbf{y} | do(\mathbf{t}), \mathbf{c})$ , where  $\mathbf{y}, \mathbf{t}, \mathbf{c} \subset \mathbf{x}$ , if there exists two proxies  $\mathbf{w}, \mathbf{n} \subset \mathbf{x}$  and a variable  $\mathbf{b} \subset \mathbf{x}$ , none of them overlapping nor containing variables from the previous subsets, such that:

- i)  $\mathbf{w}$  is conditionally independent of  $(\mathbf{t}, \mathbf{n})$  given  $\mathbf{b}, \mathbf{z}$  and  $\mathbf{c}$ . That is,  $\mathbf{w} \perp\!\!\!\perp (\mathbf{t}, \mathbf{n}) | \mathbf{b}, \mathbf{z}, \mathbf{c}$ .
- ii)  $\mathbf{n}$  is conditionally independent of  $\mathbf{y}$  given  $\mathbf{t}, \mathbf{b}, \mathbf{z}$  and  $\mathbf{c}$ . That is,  $\mathbf{y} \perp\!\!\!\perp \mathbf{n} | \mathbf{t}, \mathbf{b}, \mathbf{z}, \mathbf{c}$ .
- iii)  $(\mathbf{b}, \mathbf{z})$  forms a valid adjustment set for the query  $p(\mathbf{y} | do(\mathbf{t}), \mathbf{c})$ . That is, given  $\mathbf{c}$ , they are independent of  $\mathbf{t}$  after severing any incoming edges to it,  $do(\mathbf{t}) \perp\!\!\!\perp (\mathbf{b}, \mathbf{z}) | \mathbf{c}$ , and they block every backdoor path from  $\mathbf{t}$  to  $\mathbf{y}$ .

iv)  $\mathbf{z}$  is complete given  $\mathbf{n}$  for all  $\mathbf{t}, \mathbf{b}$ , and  $\mathbf{c}$ ,

v)  $\tilde{\mathbf{z}}$  is complete given  $\mathbf{w}$  for all  $\mathbf{b}$  and  $\mathbf{c}$ ,

and the following regularity conditions also hold:

vi)  $\iint \tilde{p}(\tilde{\mathbf{z}} | \mathbf{w}, \mathbf{b}, \mathbf{c})\tilde{p}(\mathbf{w} | \tilde{\mathbf{z}}, \mathbf{b}, \mathbf{c}) d\tilde{\mathbf{z}} d\mathbf{w} < \infty$  for all  $\mathbf{b}, \mathbf{c}$ , and

vii)  $\int \tilde{p}(\mathbf{y} | \mathbf{t}, \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c})^2 \tilde{p}(\tilde{\mathbf{z}} | \mathbf{b}, \mathbf{c}) d\tilde{\mathbf{z}} < \infty$  for all  $\mathbf{t}, \mathbf{b}$ , and  $\mathbf{c}$ .

*Proof.* First, note that the first three independence assumptions hold for both models,  $\mathcal{M}$  and  $\tilde{\mathcal{M}}$ , as they induce the same causal graph. Following the same arguments as Miao et al. (2018, Proposition 1), we have that assumptions **v)**, **vi)**, and **vii)** guarantee the existence of a function  $\tilde{h}$  such that it solves the integral equation over  $\tilde{\mathcal{M}}$ ,

$$\tilde{p}(\mathbf{y} | \mathbf{t}, \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c}) = \int \tilde{h}(\mathbf{y}, \mathbf{t}, \mathbf{b}, \mathbf{w}, \mathbf{c})\tilde{p}(\mathbf{w} | \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c}) d\mathbf{w}, \quad (10)$$

since assumption **vi)** ensures that the conditional expectation operator is compact (Carrasco et al., 2007), assumption **v)** that all square-integrable functions are in the image of the operator (i.e., the operator is surjective), and assumption **vii)** that  $\tilde{p}(\mathbf{y} | \mathbf{t}, \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c})$  is indeed part of the image.

We can show that  $\tilde{h}$  also solves a similar integral equation, this time over the other SCM,  $\mathcal{M}$ , as follows:

$$p(\mathbf{y} | \mathbf{t}, \mathbf{b}, \mathbf{n}, \mathbf{c}) = \tilde{p}(\mathbf{y} | \mathbf{t}, \mathbf{b}, \mathbf{n}, \mathbf{c}) \quad [\text{equal marginals}] \quad (11)$$

$$= \int \tilde{p}(\mathbf{y} | \mathbf{t}, \mathbf{b}, \mathbf{n}, \tilde{\mathbf{z}}, \mathbf{c})\tilde{p}(\tilde{\mathbf{z}} | \mathbf{t}, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\tilde{\mathbf{z}} \quad [\text{augment with } \tilde{\mathbf{z}}] \quad (12)$$

$$= \int \tilde{p}(\mathbf{y} | \mathbf{t}, \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c})\tilde{p}(\tilde{\mathbf{z}} | \mathbf{t}, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\tilde{\mathbf{z}} \quad [\text{assumption ii)}] \quad (13)$$

$$= \iint \tilde{h}(\mathbf{y}, \mathbf{t}, \mathbf{b}, \mathbf{w}, \mathbf{c})\tilde{p}(\mathbf{w} | \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c})\tilde{p}(\tilde{\mathbf{z}} | \mathbf{t}, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\tilde{\mathbf{z}} d\mathbf{w} \quad [\text{plug Eq. 10}] \quad (14)$$

$$= \iint \tilde{h}(\mathbf{y}, \mathbf{t}, \mathbf{b}, \mathbf{w}, \mathbf{c})\tilde{p}(\mathbf{w} | \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{t}, \mathbf{n}, \mathbf{c})\tilde{p}(\tilde{\mathbf{z}} | \mathbf{t}, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\tilde{\mathbf{z}} d\mathbf{w} \quad [\text{assumption i)}] \quad (15)$$

$$= \int \tilde{h}(\mathbf{y}, \mathbf{t}, \mathbf{b}, \mathbf{w}, \mathbf{c})p(\mathbf{w} | \mathbf{t}, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\mathbf{w}. \quad [\text{equal marginals}] \quad (16)$$

Similarly, we can relate the expression for the interventional distribution of both models:

$$\tilde{p}(\mathbf{y} | do(\mathbf{t}), \mathbf{c}) = \int \tilde{p}(\mathbf{y} | do(\mathbf{t}), \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c})\tilde{p}(\mathbf{b}, \tilde{\mathbf{z}} | \mathbf{c}) d\mathbf{b} d\tilde{\mathbf{z}} \quad [\text{augment and assumption iii)}] \quad (17)$$

$$= \int \tilde{p}(\mathbf{y} | \mathbf{t}, \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c})\tilde{p}(\mathbf{b}, \tilde{\mathbf{z}} | \mathbf{c}) d\mathbf{b} d\tilde{\mathbf{z}} \quad [\text{backdoor criterion}] \quad (18)$$

$$= \iint \tilde{h}(y, t, \mathbf{b}, \mathbf{w}, \mathbf{c}) \tilde{p}(\mathbf{w} | \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c}) \tilde{p}(\mathbf{b}, \tilde{\mathbf{z}} | \mathbf{c}) d\mathbf{b} d\mathbf{w} d\tilde{\mathbf{z}} \quad [\text{plug Eq. 10}] \quad (19)$$

$$= \int \tilde{h}(y, t, \mathbf{b}, \mathbf{w}, \mathbf{c}) p(\mathbf{b}, \mathbf{w} | \mathbf{c}) d\mathbf{b} d\mathbf{w} \quad [\text{equal marginals}] \quad (20)$$

$$= p(y | \text{do}(t), \mathbf{c}), \quad (21)$$

where the last equality is a consequence of Eq. 16 as we will show now. More specifically, we have that

$$p(y | t, \mathbf{b}, \mathbf{n}, \mathbf{c}) = \int \tilde{h}(y, t, \mathbf{b}, \mathbf{w}, \mathbf{c}) p(\mathbf{w} | t, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\mathbf{w} \quad [\text{Eq. 16}] \quad (22)$$

$$= \iint \tilde{h}(y, t, \mathbf{b}, \mathbf{w}, \mathbf{c}) p(\mathbf{w} | \mathbf{b}, \mathbf{z}, t, \mathbf{n}, \mathbf{c}) p(\mathbf{z} | t, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\mathbf{w} d\mathbf{z}, \quad [\text{augment with } \mathbf{z}] \quad (23)$$

$$= \iint \tilde{h}(y, t, \mathbf{b}, \mathbf{w}, \mathbf{c}) p(\mathbf{w} | \mathbf{b}, \mathbf{z}, \mathbf{c}) p(\mathbf{z} | t, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\mathbf{w} d\mathbf{z}. \quad [\text{assumption i}] \quad (24)$$

Similarly, we have that

$$p(y | t, \mathbf{b}, \mathbf{n}, \mathbf{c}) = \int p(y | t, \mathbf{b}, \mathbf{n}, \mathbf{z}, \mathbf{c}) p(\mathbf{z} | t, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\mathbf{z} \quad [\text{augment with } \mathbf{z}] \quad (25)$$

$$= \int p(y | t, \mathbf{b}, \mathbf{z}, \mathbf{c}) p(\mathbf{z} | t, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\mathbf{z}. \quad [\text{assumption ii}] \quad (26)$$

Now, equating both expressions we have that

$$0 = \iint \left\{ p(y | t, \mathbf{b}, \mathbf{z}, \mathbf{c}) - \int \tilde{h}(y, t, \mathbf{b}, \mathbf{w}, \mathbf{c}) p(\mathbf{w} | \mathbf{b}, \mathbf{z}, \mathbf{c}) d\mathbf{w} \right\} p(\mathbf{z} | t, \mathbf{b}, \mathbf{n}, \mathbf{c}) d\mathbf{z}, \quad (27)$$

which, due to assumption iv), implies that

$$p(y | t, \mathbf{b}, \mathbf{z}, \mathbf{c}) \stackrel{\text{a.e.}}{=} \int \tilde{h}(y, t, \mathbf{b}, \mathbf{w}, \mathbf{c}) p(\mathbf{w} | \mathbf{b}, \mathbf{z}, \mathbf{c}) d\mathbf{w}. \quad (28)$$

Finally, putting all together we see that we can write the interventional distribution of the original model using  $\tilde{h}$ ,

$$p(y | \text{do}(t), \mathbf{c}) = \iint p(y | \text{do}(t), \mathbf{b}, \mathbf{z}, \mathbf{c}) p(\mathbf{b}, \mathbf{z} | \mathbf{c}) d\mathbf{b} d\mathbf{z} \quad [\text{augment and assumption iii}] \quad (29)$$

$$= \iint p(y | t, \mathbf{b}, \mathbf{z}, \mathbf{c}) p(\mathbf{b}, \mathbf{z} | \mathbf{c}) d\mathbf{b} d\mathbf{z} \quad [\text{backdoor criterion}] \quad (30)$$

$$= \iint \tilde{h}(y, t, \mathbf{b}, \mathbf{w}, \mathbf{c}) p(\mathbf{w} | \mathbf{b}, \mathbf{z}, \mathbf{c}) p(\mathbf{b}, \mathbf{z} | \mathbf{c}) d\mathbf{b} d\mathbf{z} d\mathbf{w} \quad [\text{Eq. 28}] \quad (31)$$

$$= \int \tilde{h}(y, t, \mathbf{b}, \mathbf{w}, \mathbf{c}) p(\mathbf{b}, \mathbf{w} | \mathbf{c}) d\mathbf{b} d\mathbf{w}, \quad [\text{equal marginals}] \quad (32)$$

which justifies the last equality in Eq. 21. □

Using a causal graph similar to the one presented by Miao et al. (2018), we now provide some intuition on the semantics of each random variable in Prop. A.2. More specifically, consider the causal graph that we depict in Fig. 11, and say that we want to identify the causal query  $p(y | \text{do}(t))$  (that is, the same query as in Prop. A.2 but with  $\mathbf{c} = \emptyset$ ). As it is common in the causal inference literature (Peters et al., 2017; Spirtes et al., 2001),  $t$  and  $y$  represent the treatment and outcome random variables. More specific to Prop. A.2 are  $\mathbf{n}$  and  $\mathbf{w}$ . The variable  $\mathbf{w}$  is a proxy variable whose role is that of distinguishing the information from  $\mathbf{z}$  and other variables, to reconstruct the information of  $\mathbf{z}$  and block the backdoor path that  $\mathbf{z}$  would usually block. Similarly, the variable  $\mathbf{n}$  is another proxy variable which, in this case, serves the purpose of verifying that the substitute formed

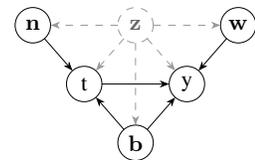


Figure 11: Example for which Prop. A.2 applies, and where  $\mathbf{b}$  is not the empty set.

with  $\mathbf{w}$  is indeed a good substitute. Finally, the variable  $\mathbf{b}$  serves the purpose of blocking all the remaining backdoor paths that  $\mathbf{z}$  may not block, so that we can apply the backdoor criterion.

Moreover, note that for all interventional queries we will let  $\mathbf{c}$  be the empty set, similar to the results proved by Miao et al. (2018) and Wang & Blei (2021). We will consider cases when  $\mathbf{c}$  is not empty later in App. A.3 to prove counterfactual identifiability. Note also that Prop. A.2 reduces to the existing results when we have that  $\mathbf{c} = \mathbf{b} = \emptyset$ .

Using this general proposition, we can now reason about causal identifiability in a wide range of scenarios, where  $t$  and  $y$  may or may not be directly caused by the hidden confounder, as we show in the following subsections.

### A.2.1 UNCONFOUNDED CASE

First, we consider the case where neither  $t$  nor  $y$  are directly affected by the hidden confounder, i.e.,  $\mathbf{z} \notin \text{ch}(\mathbf{z})$ . In this case, the proof can be simplified and drop the requirement of finding valid proxy variables.

**Corollary A.3** (Unconfounded case). *Assume that we have two SCMs  $\mathcal{M} := (\mathbf{f}, P_{\mathbf{u}}, P_{\mathbf{z}})$  and  $\tilde{\mathcal{M}} := (\tilde{\mathbf{f}}, P_{\tilde{\mathbf{u}}}, P_{\tilde{\mathbf{z}}})$  that are compatible, i.e., they induce the same causal graph, and which coincide in their marginal distributions,  $p(\mathbf{x}) = \tilde{p}(\mathbf{x})$ . Assume that  $y, t \notin \text{ch}(\mathbf{z})$ . Then,  $p(y \mid \text{do}(t), \mathbf{c}) = \tilde{p}(y \mid \text{do}(t), \mathbf{c})$ , where  $y, t, \mathbf{c} \subset \mathbf{x}$ .*

*Proof.* The proof follows directly by applying Prop. A.2 with the minimal subset  $\mathbf{b} \subset \text{pa}(t) \setminus \{\mathbf{c}\}$  that blocks all the backdoor paths, and by noticing that in this case there is no need to use the variables  $\mathbf{z}$  and  $\tilde{\mathbf{z}}$ . That is, we can go from Eq. 17 to Eq. 21 directly by using only  $\mathbf{b}$  and the equal-marginals assumption:

$$\tilde{p}(y \mid \text{do}(t), \mathbf{c}) = \int \tilde{p}(y \mid \text{do}(t), \mathbf{b}, \mathbf{c}) \tilde{p}(\mathbf{b} \mid \mathbf{c}) d\mathbf{b} \tag{33}$$

$$= \int \tilde{p}(y \mid t, \mathbf{b}, \mathbf{c}) \tilde{p}(\mathbf{b} \mid \mathbf{c}) d\mathbf{b} \tag{34}$$

$$= \int p(y \mid t, \mathbf{b}, \mathbf{c}) p(\mathbf{b} \mid \mathbf{c}) d\mathbf{b} \tag{35}$$

$$= p(y \mid \text{do}(t), \mathbf{c}). \tag{36}$$

□

Even though we can leverage and simplify Prop. A.2 as shown above, it is worth remarking that, for this particular case, the model identifiability results described in App. A.1 are stronger, as it provides results on the identifiability of the causal generators and exogenous distributions, and therefore of any causal query derived from them.

### A.2.2 FULLY CONFOUNDED CASE

In the case where both variables are directly confounded by  $\mathbf{z}$ , we cannot do much but to see whether we can apply Prop. A.2 with  $\mathbf{c} = \emptyset$  and a valid  $\mathbf{b}$ . If we manage to find two proxies  $\mathbf{w}$  and  $\mathbf{n}$  that hold the independence conditions from Prop. A.2 and that change the posterior of  $\mathbf{z}$  enough, then we can use the proposition to ensure the identifiability of the query. Otherwise, the query is not identifiable and the model might or might not estimate the query correctly.

### A.2.3 CONFOUNDED OUTCOME CASE

For the case where only the outcome random variable is directly affected by the hidden variable, we can apply a similar reasoning as we did in the case with no direct confounding, although this time we cannot leverage the model identifiability results from Javaloy et al. (2023). More specifically:

**Corollary A.4** (Confounded-outcome case). *Assume that we have two SCMs  $\mathcal{M} := (\mathbf{f}, P_{\mathbf{u}}, P_{\mathbf{z}})$  and  $\tilde{\mathcal{M}} := (\tilde{\mathbf{f}}, P_{\tilde{\mathbf{u}}}, P_{\tilde{\mathbf{z}}})$  that are compatible, i.e., they induce the same causal graph, and which coincide in their marginal distributions,  $p(\mathbf{x}) = \tilde{p}(\mathbf{x})$ . Assume that  $t \notin \text{ch}(\mathbf{z})$ . Then,  $p(y \mid \text{do}(t), \mathbf{c}) = \tilde{p}(y \mid \text{do}(t), \mathbf{c})$ , where  $y, t, \mathbf{c} \subset \mathbf{x}$ .*

*Proof.* The proof is identical to that of Cor. A.3. □

**Front-door example.** While the proof above is trivial given the previous results, it is worth stressing that for them to hold it is necessary to model the hidden confounder as we do in this work with the proposed Decaf, and that other approaches may not work for all cases. As an example, consider the SCM depicted in Fig. 12, where we have that the outcome is directly confounded by  $\mathbf{z}$ , while  $t$  is not. In this case, a Decaf should be able to identify the true causal query  $p(y \mid \text{do}(t))$ , using  $\tilde{\mathbf{z}}$  to model the influence of  $\mathbf{b}$

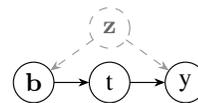


Figure 12: Textbook example of a front-door in a SCM.

onto  $y$  that is not explained through  $t$ . Other models that do not model  $z$  (e.g., an unaware causal normalizing flow (Javaloy et al., 2023)), would not be able to match the observed marginal likelihood as they assume that  $y \perp\!\!\!\perp \mathbf{b} \mid t$  yet we know that  $y \not\perp\!\!\!\perp \mathbf{b} \mid t$  in the true model. Even more, with those models we would have that  $p(y \mid \text{do}(t)) = p(y \mid t)$  which is clearly false by just looking at Fig. 12.

To be even more explicit, in this case we would have a factorization of the form

$$\tilde{p}(\mathbf{b}, t, y, \tilde{\mathbf{z}}) = \tilde{p}(\tilde{\mathbf{z}})\tilde{p}(\mathbf{b} \mid \tilde{\mathbf{z}})\tilde{p}(t \mid \mathbf{b})\tilde{p}(y \mid t, \tilde{\mathbf{z}}). \quad (37)$$

Then, the estimated interventional distribution that a Decaf estimates as  $\int \tilde{p}(y \mid t, \tilde{\mathbf{z}}) d\tilde{\mathbf{z}}$  equals the true one:

$$p(y \mid \text{do}(t)) = \int p(y \mid t, \mathbf{b})p(\mathbf{b}) d\mathbf{b} \quad [\mathbf{b} \text{ forms a valid adjustment set}] \quad (38)$$

$$= \int \left\{ \int \tilde{p}(y \mid t, \mathbf{b}, \tilde{\mathbf{z}})\tilde{p}(\tilde{\mathbf{z}} \mid t, \mathbf{b}) d\tilde{\mathbf{z}} \right\} \tilde{p}(\mathbf{b}) d\mathbf{b} \quad [\text{latent factorization and equal marginals}] \quad (39)$$

$$= \iint \tilde{p}(y \mid t, \tilde{\mathbf{z}})\tilde{p}(\tilde{\mathbf{z}} \mid \mathbf{b})\tilde{p}(\mathbf{b}) d\mathbf{b} d\tilde{\mathbf{z}} \quad [\text{causal graph factorization in Eq. 37}] \quad (40)$$

$$= \int \tilde{p}(y \mid t, \tilde{\mathbf{z}})\tilde{p}(\tilde{\mathbf{z}}) d\tilde{\mathbf{z}} \quad [\text{marginalize } \mathbf{b}] \quad (41)$$

$$= \tilde{p}(y \mid \text{do}(t)). \quad (42)$$

Remarkably, the identification of  $p(y \mid \text{do}(t))$  allows us to solve also the query  $p(y \mid \text{do}(\mathbf{b}))$  leveraging the frontdoor criterion (Peters et al., 2017).

$$p(y \mid \text{do}(\mathbf{b})) = \int p(t \mid \mathbf{b})p(y \mid \text{do}(t)) dt \quad [\text{frontdoor criterion}] \quad (43)$$

$$= \int \tilde{p}(t \mid \mathbf{b}) \int \tilde{p}(y \mid t, \tilde{\mathbf{z}})\tilde{p}(\tilde{\mathbf{z}}) d\tilde{\mathbf{z}} dt \quad [\text{plug in Eq. 41 and equal marginals}] \quad (44)$$

$$= \tilde{p}(y \mid \text{do}(\mathbf{b})) \quad (45)$$

#### A.2.4 CONFOUNDED TREATMENT CASE

When only the treatment variable  $t$  is directly confounded, we can find two different scenarios: if we are able to find a valid adjustment set  $\mathbf{b}$  blocking all confounded paths, in which case we can reason just as in the other partially confounded case, and otherwise, where we rely on the identifiability with respect to this invalid adjustment set. For example, if it happens to be a parent of  $y$  which is directly caused by the treatment variable  $t$  and the hidden confounder  $z$  as in Fig. 13, we cannot find a valid adjustment set for the causal query, but an invalid one may still serve us if we can identify the same query with the adjustment set as outcome.

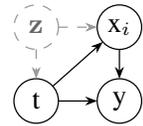


Figure 13: Case where no valid adjustment set can be found.

**Corollary A.5** (Confounded-treatment case). *Assume that we have two compatible SCMs  $\mathcal{M} := (\mathbf{f}, P_{\mathbf{u}}, P_{\mathbf{z}})$  and  $\tilde{\mathcal{M}} := (\tilde{\mathbf{f}}, P_{\tilde{\mathbf{u}}}, P_{\tilde{\mathbf{z}}})$ , i.e., they induce the same causal graph, and which coincide in their marginal distributions,  $p(\mathbf{x}) = \tilde{p}(\mathbf{x})$ . Assume also that  $y \notin \text{ch}(\mathbf{z})$ . Then,  $p(y \mid \text{do}(t), \mathbf{c}) = \tilde{p}(y \mid \text{do}(t), \mathbf{c})$ , where  $y, t, \mathbf{c} \subset \mathbf{x}$  if there exists a subset  $\mathbf{b} \subset \mathbf{x}$  not containing variables from the previous subsets, such that one of the following two conditions are true:*

- i)  $\mathbf{b}$  forms a valid adjustment set for the query  $p(y \mid \text{do}(t), \mathbf{c})$ .
- ii)  $\mathbf{b}$  forms an invalid adjustment set for the query  $p(y \mid \text{do}(t), \mathbf{c})$  but the query  $p(\mathbf{b} \mid \text{do}(t), \mathbf{c})$  is identifiable. That is,  $\mathbf{b}$  blocks all the backdoor paths, and  $p(\mathbf{b} \mid \text{do}(t), \mathbf{c}) = \tilde{p}(\mathbf{b} \mid \text{do}(t), \mathbf{c})$ .

*Proof.* If condition **i**) holds, then we have a valid adjustment set, and the proof is identical to that of Cor. A.3.

Otherwise, if condition **ii**) holds, we have that the interventional query on  $y$  equals the observational query when conditioned on  $\mathbf{b}$ , but that now  $\mathbf{b}$  is not independent of  $\text{do}(t)$ , i.e.,

$$\tilde{p}(y \mid \text{do}(t), \mathbf{c}) = \int \tilde{p}(y \mid \text{do}(t), \mathbf{b}, \mathbf{c})\tilde{p}(\mathbf{b} \mid \text{do}(t), \mathbf{c}) d\mathbf{b} \quad (46)$$

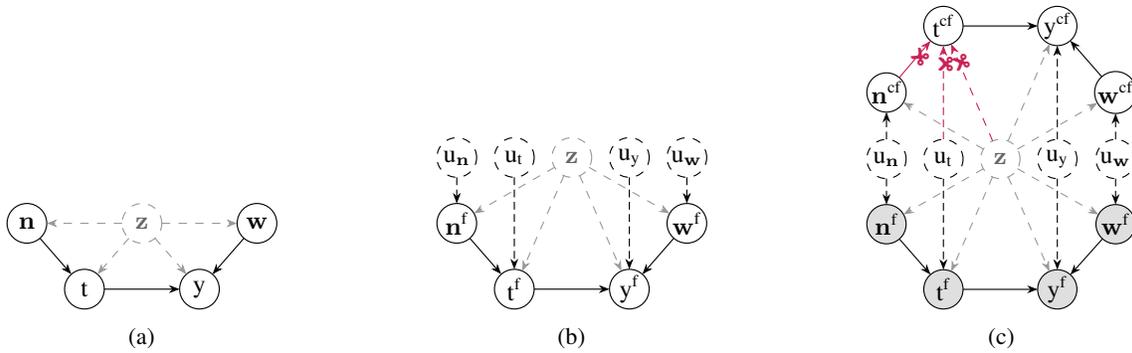


Figure 14: Example of the transition from (a) the regular depiction of a (confounded) SCM, to (b) an explicit SCM where the exogenous variables are drawn, and (c) a counterfactual twin SCM where the data-generating process is replicated in the “factual and counterfactual worlds”. Besides, figure (c) also depicts which nodes are observed and which edges are severed, in order to compute a counterfactual query of the type  $p(y^{cf} \mid \text{do}(t^{cf}), \mathbf{x}^f)$ .

$$= \int \tilde{p}(y \mid t, \mathbf{b}, \mathbf{c}) \tilde{p}(\mathbf{b} \mid \text{do}(t), \mathbf{c}) d\mathbf{b} \quad (47)$$

$$= \int p(y \mid t, \mathbf{b}, \mathbf{c}) p(\mathbf{b} \mid \text{do}(t), \mathbf{c}) d\mathbf{b} \quad (48)$$

$$= p(y \mid \text{do}(t), \mathbf{c}), \quad (49)$$

where we needed to use that the query  $p(\mathbf{b} \mid \text{do}(t), \mathbf{c})$  is identifiable in the third equality.  $\square$

### A.3 Counterfactual query identifiability

In this section, we show that counterfactual query identifiability is a direct result of the interventional query identifiability from the previous section.

In order to formally define counterfactuals, in this section we introduce the concept of counterfactual SCMs in a somewhat novel way. Namely, we combine the concepts of twin networks from Pearl (2009) (which replicates the data-generating process) and that of counterfactual SCMs from Peters et al. (2017) (which defines a counterfactual prior to the intervention).

**Definition 6** (Counterfactual twin SCM). Given a SCM  $\mathcal{M} = (\mathbf{f}, P_{\mathbf{u}}, P_{\mathbf{z}})$ , we define its counterfactual twin SCM as a SCM  $\mathcal{M}^{cf}$  where all structural equations are duplicated, and the exogenous noise is shared across replications, and where additionally one of the halves is observed (“the factual world”), and the other half is unobserved (“the counterfactual world”).

We provide in Fig. 14 a more intuitive depiction on the construction of these counterfactual twin networks. From this definition, one can recover the counterfactual SCM defined by Peters et al. (2017) by just focusing on the replicated part of the counterfactual twin network, and conditioning the exogenous noise and hidden confounder on the observed half, i.e.,  $(\mathbf{f}, P_{\mathbf{u} \mid \mathbf{x}^f}, P_{\mathbf{z} \mid \mathbf{x}^f})$ . Similarly, one can compute the usual counterfactual query by performing an intervention on the counterfactual twin network, i.e., by replacing the intervened equations by the constant intervened value, and computing the query conditioned on the factual variables,  $p(y^{cf} \mid \text{do}(t^{cf}), \mathbf{x}^f)$ . This is visually represented in Fig. 14c.

In order to prove query identifiability in the counterfactual setting, we need to use the following technical result regarding the completeness of a random variable:

**Lemma A.6.** *If a random variable  $\mathbf{z}$  is complete given  $\mathbf{n}$  for all  $\mathbf{b}$ , as given by Def. 5, then it is complete given  $\mathbf{n}$  for all  $\mathbf{b}$  and  $\mathbf{c}$ , where  $\mathbf{c}$  is another continuous random variable.*

*Proof.* We prove this result by contradiction. Assume that the result does not hold, then there must exist a non-zero measure subset of the space of  $\mathbf{b} \times \mathbf{c}$  for which there exists a square-integrable function  $g(\cdot)$  such that  $\int g(\mathbf{z}, \mathbf{b}, \mathbf{c}) p(\mathbf{z} \mid \mathbf{b}, \mathbf{c}, \mathbf{n}) d\mathbf{z} = 0$  for almost all  $\mathbf{n}$ , but  $g(\mathbf{z}, \mathbf{b}, \mathbf{c}) \neq 0$  for almost all  $\mathbf{z}$ .

Since this subset has positive measure, there must contain an  $\varepsilon$ -ball within. If we now focus on the  $\mathbf{b}$ -projection of this ball where we fix  $\mathbf{c}$  to its value on the centre, we have that it is a subset of non-zero measure in the space of  $\mathbf{b}$  (as otherwise it would be zero-measure in the Cartesian-product measure), where the function  $g(\cdot, \mathbf{c})$  breaks our initial assumption of the completeness of  $\mathbf{z}$ . Thus, we reach a contradiction.  $\square$

Given Def. 6, it is rather intuitive that, if a causal query is identifiable in a SCM  $\mathcal{M}$ , then it has to be identifiable in both halves of its induced counterfactual twin SCM  $\mathcal{M}^{cf}$ , as they are identical. More importantly, we can now leverage again Prop. A.2, this time with  $\mathbf{c} = \mathbf{x}^f$ , to prove counterfactual query identifiability whenever we have interventional query identifiability.

**Proposition A.7** (Counterfactual identifiability). *Assume that we have two SCMs  $\mathcal{M} := (\mathbf{f}, P_{\mathbf{u}}, P_{\mathbf{z}})$  and  $\tilde{\mathcal{M}} := (\tilde{\mathbf{f}}, P_{\tilde{\mathbf{u}}}, P_{\tilde{\mathbf{z}}})$  that are compatible, i.e., they induce the same causal graph, and which coincide in their marginal distributions,  $p(\mathbf{x}) = \tilde{p}(\mathbf{x})$ . Then, if a query  $p(y | do(t))$  is identifiable in the sense of Prop. A.2, where  $y, t \subset \mathbf{x}$ , then the query  $p(y^{cf} | do(t^{cf}), \mathbf{x}^f)$  is also identifiable in their induced counterfactual twin SCMs as long as the regularity conditions still hold, i.e., if:*

- i)  $\iint \tilde{p}(\tilde{\mathbf{z}} | \mathbf{w}, \mathbf{b}, \mathbf{c}) \tilde{p}(\mathbf{w} | \tilde{\mathbf{z}}, \mathbf{b}, \mathbf{c}) d\tilde{\mathbf{z}} d\mathbf{w} < \infty$  for all  $\mathbf{b}, \mathbf{c}$ , and
- ii)  $\int \tilde{p}(y | t, \mathbf{b}, \tilde{\mathbf{z}}, \mathbf{c})^2 \tilde{p}(\tilde{\mathbf{z}} | \mathbf{b}, \mathbf{c}) d\tilde{\mathbf{z}} < \infty$  for all  $t, \mathbf{b}$ , and  $\mathbf{c}$ .

*Proof.* We essentially need to prove that the independence and completeness assumptions keep holding when we add the factual covariate,  $\mathbf{c} = \mathbf{x}^f$ .

For the independence, we need to show that, if we have a set of variables that fulfil the independence conditions from Prop. A.2, then this set of variables keeps holding them if we include  $\mathbf{c} = \mathbf{x}^f$ . This is, however, easy to show since factual and counterfactual variables only have “tail-to-tail” dependencies, i.e., they are connected only through the shared exogenous variables. As a result, if two variables from the same half are conditionally independent given a third set of variables, conditioning on the other half cannot change this independence.

For the completeness, we need to show that introducing the factual variable retain the completeness assumed in Prop. A.2. However, this is direct to show using the technical result in Lemma A.6. Specifically, it holds that

- i)  $\mathbf{z}$  is complete given  $\mathbf{n}$  for all  $t, \mathbf{b}$ , and  $\mathbf{c}$ , and
- ii)  $\tilde{\mathbf{z}}$  is complete given  $\mathbf{w}$  for all  $\mathbf{b}$  and  $\mathbf{c}$ .

Therefore, we have shown that the requirements of Prop. A.2 hold when we append a factual variable in the twin network, and thus we can reapply all the results from the previous sections to the counterfactual cases.  $\square$

It is important to note that, while the results above provide counterfactual identifiability whenever we have interventional identifiability, we still rely on how much of a good approximation the encoder is to the inverse of the decoder in the proposed Decaf model. That is, the quality of the encoder determines how well we can perform the abduction step to compute counterfactuals. This consideration is unique to counterfactuals, as we just had to sample the latent variable as usual in the case of interventional queries.

## B Experimental details and additional results

### B.1 Ablation study

First, in Fig. 15 we present the ATE error committed for each combination of proxies and latent dimension, complementary to the figure of the presented in §7.1. If we observe the ATE error, we extract the same conclusion as observing counterfactual error, the causal effect is not recoverable with less than two proxies, and more proxies result in better estimates. On the other hand, the selection of the dimension of the latent space bigger than the true dimension of the latent confounders does not affect the performance negatively.

In addition, we show the equations that we have used for the ablation study. There exist two unobserved confounders,  $\mathbf{z}_1$  and  $\mathbf{z}_2$ . The set of all observed proxies  $\{x_3, x_4, \dots, x_{12}\}$  is represented in the graph of §7.1 as  $\mathbf{x}_3$ . Note that the proxies available in the nonlinear experiment are bounded or periodic, specially sigmoids and hyperbolic tangents saturate and  $\max(0, x)$  loses all the information about the confounder for negative values and sines and cosines are periodic functions. In other words, the distributions  $p(\mathbf{z} | x_i)$  are not complete, we lose information about  $\mathbf{z}$  when in the transformations to each  $x$ . However, if we add more

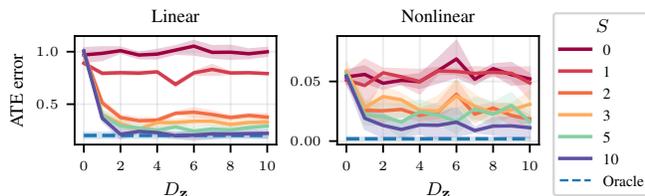


Figure 15: ATE absolute error varying the number of available proxies ( $S$ ) and the dimensionality of the latent space ( $D_z$ ). Mean and 95% confidence interval over 5 realizations and all interventions, made in percentiles 25, 50 and 75 of  $\mathbf{x}_1$ . Oracle represents a causal normalizing flow that observes  $\mathbf{z}$ .

proxies of the confounders, the information that the proxies contain about the confounder is higher, and the causal effect of  $x_1$  on  $x_2$  becomes recoverable.

1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182

Linear	Nonlinear
$\begin{cases} \mathbf{z}_1 \sim P_{\mathbf{z}_1} \\ \mathbf{z}_2 \sim P_{\mathbf{z}_2} \\ x_1 = 1.5 \cdot \mathbf{z}_1 + 0.5 \cdot \mathbf{z}_2 + 0.4 \cdot \mathbf{u}_1 \\ x_2 = -0.75 \cdot \mathbf{z}_1 + 0.6 \cdot \mathbf{z}_2 + 0.9 \cdot x_1 + 0.3 \cdot \mathbf{u}_2 \\ x_3 = -0.5 \cdot \mathbf{z}_1 + 0.3 \cdot \mathbf{z}_2 + 0.5 \cdot \mathbf{u}_3 \\ x_4 = 0.75 \cdot \mathbf{z}_1 - 0.4 \cdot \mathbf{z}_2 + 0.4 \cdot \mathbf{u}_4 \\ x_5 = -0.85 \cdot \mathbf{z}_1 + 0.6 \cdot \mathbf{z}_2 + 0.6 \cdot \mathbf{u}_5 \\ x_6 = 0.6 \cdot \mathbf{z}_1 + 0.6 \cdot \mathbf{z}_2 + 0.55 \cdot \mathbf{u}_6 \\ x_7 = -0.8 \cdot \mathbf{z}_1 + 0.4 \cdot \mathbf{z}_2 + 0.4 \cdot \mathbf{u}_7 \\ x_8 = 0.9 \cdot \mathbf{z}_1 - 0.7 \cdot \mathbf{z}_2 + 0.6 \cdot \mathbf{u}_8 \\ x_9 = -0.72 \cdot \mathbf{z}_1 + 0.5 \cdot \mathbf{z}_2 + 0.56 \cdot \mathbf{u}_9 \\ x_{10} = 0.78 \cdot \mathbf{z}_1 + 0.4 \cdot \mathbf{z}_2 + 0.58 \cdot \mathbf{u}_{10} \\ x_{11} = -0.55 \cdot \mathbf{z}_1 + 0.7 \cdot \mathbf{z}_2 + 0.6 \cdot \mathbf{u}_{11} \\ x_{12} = 0.88 \cdot \mathbf{z}_1 + 0.3 \cdot \mathbf{z}_2 + 0.4 \cdot \mathbf{u}_{12} \end{cases}$	$\begin{cases} \mathbf{z}_1 \sim P_{\mathbf{z}_1} \\ \mathbf{z}_2 \sim P_{\mathbf{z}_2} \\ x_1 = \frac{\mathbf{z}_1^2}{4} \cdot \sin\left(\frac{\mathbf{z}_2}{2}\right) + \mathbf{z}_1 + 0.6 \cdot \mathbf{u}_1 \\ x_2 = \frac{\mathbf{z}_1 \cdot x_1}{4} + 0.8 \cdot \mathbf{z}_2 + 0.5 \cdot x_1 + x_1 \cdot \mathbf{u}_2 \cdot 0.3 + 0.2 \cdot \mathbf{u}_2 \\ x_3 = 0.6 \cdot \mathbf{z}_1^2 + \left(\frac{\mathbf{z}_2}{4}\right)^3 + 0.3 \cdot \sin\left(\frac{\mathbf{z}_2}{2}\right) + 0.5 \cdot \mathbf{u}_3 \\ x_4 = \sin\left(\frac{\mathbf{z}_1}{2}\right) + \cos\left(\frac{\mathbf{z}_2}{3}\right) + 0.4 \cdot \mathbf{u}_4 \\ x_5 = \cos\left(\frac{\mathbf{z}_1}{2}\right) - \tanh\left(\frac{\mathbf{z}_2}{3}\right) + 0.6 \cdot \mathbf{u}_5 \\ x_6 = \tanh\left(\frac{\mathbf{z}_1}{2}\right) + \sigma\left(\frac{\mathbf{z}_2}{2}\right) + 0.55 \cdot \mathbf{u}_6 \\ x_7 = \sigma\left(\frac{\mathbf{z}_1}{2}\right) + \max(0, -\mathbf{z}_2) + 0.4 \cdot \mathbf{u}_7 \\ x_8 = \max(0, \mathbf{z}_1) - 0.5 \cdot \max(0, \mathbf{z}_2) + 0.6 \cdot \mathbf{u}_8 \\ x_9 = \max(0, -\mathbf{z}_1) + 0.3 \cdot \max(0, -\mathbf{z}_2) + 0.5 \cdot \mathbf{z}_1 \cdot \mathbf{u}_9 \\ x_{10} = 0.8 \cdot \max(0, \mathbf{z}_1) + 0.3 \cdot \max(0, \mathbf{z}_2) + 0.58 \cdot \mathbf{u}_{10} \\ x_{11} = 0.75 \cdot \max(0, -\mathbf{z}_1) + 0.5 \cdot \max(0, \mathbf{z}_2) + 0.6 \cdot \mathbf{u}_{11} \\ x_{12} = 0.3 \cdot \mathbf{z}_1^3 + 0.5 \cdot  \mathbf{z}_2  + 0.4 \cdot \mathbf{u}_{12} \end{cases}$

## B.2 Semi-synthetic Sachs’ dataset

This dataset represents a network of protein-signaling in human T lymphocytes. Every variable, except PKA and PLC $\beta$  can be intervened upon; therefore, there is not only one causal query of interest, but tens of possible causal queries can arise in this setting. This highlights one of the strengths of Decaf, because we only need a single trained model to answer all identifiable causal queries.

The original data contains a total of 853 observational samples; however, we have decided to evaluate our model on semi-synthetic data because of the following reasons:

- The original network of [Sachs et al. \(2005\)](#) contains cycles, which is violation of one of our assumptions. However, we have found different versions of the causal graph ([Kaltenpoth & Vreeken, 2023](#); [Luo & Zhao, 2011](#)) that do not contain cycles. Therefore, we have decided to employ the causal graph that appears in the library *bnlearn* ([Scutari, 2010](#))—a recognized library for Bayesian Network learning—as ground truth causal graph. The best way to ensure that the causal graph used is the ground truth is by generating samples according to the causal graph. In addition, that causal graph is the one used by [Chao et al. \(2023\)](#).
- We can compare our model with one of the baseline models, DCM, with the same dataset as [Chao et al. \(2023\)](#) used.
- Semi-synthetic data allows us to compute all metrics to evaluate causal queries, having the ground truth.

For generating the data in this experiment, we have followed the procedure proposed by [Chao et al. \(2023\)](#), where they take the causal graph of [Sachs et al. \(2005\)](#) and the empirical distribution of the root nodes, and generate the rest of the variables with random non-linear mechanisms. In addition, exogenous variables have been included in an additive and non-additive manner, respectively.

In the following, we complement the figures presented in §7 with a table that summarizes all the interesting metrics, evaluated on the **confounded identifiable** causal queries shown in [Fig. 7](#). Interventional distributions and counterfactuals have been computed intervening in percentiles 25, 50 and 75 of the intervened variable.

Since observational MMD is computed only once, the statistics given in [Tab 2](#) are calculated *only* over 5 runs. On the other hand, we have as many interventional MMDs per run as interventions have been made. However, the statistics of

1207  
1208  
1209

interventional MMD are computed over all the interventions of all intervened variables and 5 runs (5 runs  $\times$  3 intervened variables = 15 samples). Finally, statistics over counterfactual error and ate error aggregate all the intervention-outcome pairs over the five runs. For example, in this case we intervene on 3 variables, performing 3 different interventions and we evaluate on 3, 2, and 1 variable respectively for each intervened variable, and we have a total of  $(3+2+1)\times 3\times 5 = 90$  different measurements to compute the statistics.

Table 2: Performance metrics on Sachs datasets. Mean<sub>std</sub> over five runs and all causal queries of interest. Interventions on Raf, Mek and Akt and evaluating on **confounded identifiable** effects. Bold indicates significantly better results (95% CI from a Mann-Whitney U test). Lower error values indicate better performance.

Model		Additive				Nonadditive			
		MMD obs $\times 10^4$	MMD int $\times 10^4$	ATE err  $\times 10^2$	CF err  $\times 10^2$	MMD obs $\times 10^4$	MMD int $\times 10^4$	ATE err  $\times 10^2$	CF err  $\times 10^2$
Oracle	CNF	4.84 <sub>1.84</sub>	7.50 <sub>6.17</sub>	6.05 <sub>6.83</sub>	10.03 <sub>10.29</sub>	5.96 <sub>2.37</sub>	6.71 <sub>2.97</sub>	2.34 <sub>2.02</sub>	4.84 <sub>3.43</sub>
Aware	Decaf 	<b>2.15</b> <sub>0.54</sub>	<b>7.82</b> <sub>9.59</sub>	<b>17.89</b> <sub>17.72</sub>	<b>17.92</b> <sub>4.01</sub>	5.12 <sub>2.42</sub>	<b>5.39</b> <sub>3.33</sub>	<b>3.26</b> <sub>4.09</sub>	<b>6.82</b> <sub>4.65</sub>
	Deconfounder	–	–	34.34 <sub>33.45</sub>	71.13 <sub>86.98</sub>	–	–	8.14 <sub>10.69</sub>	63.15 <sub>79.12</sub>
Unaware	CNF	5.80 <sub>1.58</sub>	73.94 <sub>88.78</sub>	44.49 <sub>39.12</sub>	56.09 <sub>38.89</sub>	5.11 <sub>1.90</sub>	12.79 <sub>20.73</sub>	9.74 <sub>15.71</sub>	15.15 <sub>15.37</sub>
	ANM	83.86 <sub>13.41</sub>	110.28 <sub>112.43</sub>	22.42 <sub>14.06</sub>	29.40 <sub>12.22</sub>	81.90 <sub>7.21</sub>	60.40 <sub>144.08</sub>	23.88 <sub>13.94</sub>	28.97 <sub>12.44</sub>
	DCM	87.80 <sub>2.95</sub>	125.59 <sub>118.20</sub>	21.21 <sub>11.34</sub>	28.25 <sub>6.96</sub>	14.23 <sub>4.57</sub>	69.74 <sub>390.81</sub>	8.44 <sub>7.96</sub>	27.50 <sub>23.71</sub>

The metrics in Tab 2 indicate that Decaf outperforms all baselines across all interventional and counterfactual causal queries in both settings of the semi-synthetic datasets. However, as discussed in §8, a limitation of our empirical approach is that the differences in observational MMD, the selection criterion for CGMs, are marginal between the *oracle*, Decaf, and CNF. Notably, Decaf even achieves a lower MMD than the *oracle*. This discrepancy arises because the number of variables is large, and the MMD differences are on the order of  $10^{-4}$ .

### B.3 Semi-synthetic Ecoli70 dataset

The Ecoli 70 dataset represent the gene expression of 46 genes of the RNA-seq of *Escherichia coli* bacteria. The assumed causal graph comes from the study of (Schäfer & Strimmer, 2005), which provides insight into the regulatory mechanisms governing *E. coli* gene expression. Examples of interventions in these networks are gene knockout and gene overexpression (Long & Antoniewicz, 2014). A priori, there could be several variables in which intervening can be interesting in evaluating the effects in the cell.

For this experiment, we have generated the data in the same way as done with Sachs’ dataset with random mechanisms, but in this case, since we do not have enough samples, root nodes follow standard Gaussian distributions. We have included an additive and a nonadditive ways of including exogenous variables. In this case, we have used a semi-synthetic dataset because the real dataset available in *bnlearn* (Scutari, 2010) contains only 9 samples.

In Fig. 1 is presented the causal graph of this setting.

In addition, note that Fig. 1 has been extracted from our Alg. 5 of causal effect identifiability. That is, we have specified the causal graph and the variables that are unmeasured, and our Algorithm returns (in green) all the paths that are identifiable by Decaf. Consider that black arrows are also identifiable, not only by Decaf, but also for any CGM that approximates the observed data. In red, arrows that are not identifiable by Decaf because there are not enough proxies to infer an unbiased causal effect.

A table summarizing the results obtained in the estimation **confounded identifiable** causal queries are presented in Tab 3. The statistics have been computed in the same way as in Sachs’ dataset. In the case of ATE and CF error, they have been computed only on the *direct* confounded identifiable paths, i.e., the green paths in Fig. 1.

Decaf significantly outperforms the baselines in ATE and counterfactual estimation in the additive setting and in ATE estimation in the nonadditive setting. The MMD differences, both observational and interventional, are negligible between the *oracle*, Decaf, and CNF, likely due to the high number of variables diluting estimation bias. Counterfactual differences in the nonadditive setting are also insignificant. However, compared to the *oracle*, the gap between the *oracle* and *unaware* CGMs is smaller than in the additive case. While Decaf reaches an intermediate point, the difference remains insignificant.

#### B.3.1 COMMENT ON DECONFOUNDER RESULTS

One may realize that the errors committed by the Deconfounder of (Wang & Blei, 2019; 2021) are greater than the errors committed by the unaware models. First of all, we want to underline that, although the Deconfounder allows us to predict

Table 3: Performance metrics on Ecoli70 dataset. ATE and CF error statistics computed aggregating all causal queries and 5 runs. Intervened and evaluated on the direct **confounded identifiable** causal effects of Fig. 1. Bold indicates significantly better results (95% CI from a Mann-Whitney U test). Lower error values indicate better performance.

Model	Additive				Nonadditive			
	MMD obs $\times 10^4$	MMD int $\times 10^4$	ATE err  $\times 10^2$	CF err  $\times 10^2$	MMD obs $\times 10^4$	MMD int $\times 10^4$	ATE err  $\times 10^2$	CF err  $\times 10^2$
Oracle CNF	2.34 <sub>0.62</sub>	6.05 <sub>5.28</sub>	5.04 <sub>7.42</sub>	9.91 <sub>12.46</sub>	1.49 <sub>0.57</sub>	4.05 <sub>8.22</sub>	3.51 <sub>4.84</sub>	1.67 <sub>1.64</sub>
Aware Decaf 	2.97 <sub>0.34</sub>	10.00 <sub>9.11</sub>	<b>7.29</b> <sub><b>8.61</b></sub>	<b>17.48</b> <sub><b>12.98</b></sub>	1.63 <sub>0.46</sub>	9.19 <sub>21.42</sub>	<b>8.72</b> <sub><b>17.61</b></sub>	2.15 <sub>2.10</sub>
Aware Deconfounder	–	–	27.35 <sub>26.17</sub>	82.15 <sub>116.90</sub>	–	–	30.00 <sub>33.24</sub>	9.90 <sub>9.47</sub>
Unaware CNF	2.98 <sub>1.15</sub>	10.25 <sub>12.13</sub>	23.91 <sub>25.16</sub>	34.02 <sub>23.90</sub>	1.95 <sub>0.77</sub>	10.20 <sub>20.87</sub>	12.72 <sub>19.21</sub>	2.45 <sub>2.06</sub>
Unaware ANM	32.80 <sub>2.81</sub>	44.33 <sub>17.62</sub>	21.88 <sub>23.89</sub>	31.33 <sub>20.64</sub>	13.17 <sub>3.95</sub>	27.56 <sub>31.57</sub>	15.04 <sub>18.18</sub>	2.71 <sub>1.88</sub>
Unaware DCM	31.65 <sub>0.27</sub>	49.50 <sub>36.83</sub>	24.45 <sub>33.31</sub>	30.22 <sub>24.83</sub>	18.78 <sub>6.01</sub>	33.37 <sub>36.14</sub>	15.07 <sub>22.37</sub>	2.36 <sub>2.08</sub>

counterfactuals, the algorithm does not present any guarantees of a correct counterfactual estimation because it does not model the exogenous variables of the SCM. That is the reason of the bad performance in counterfactual estimation.

On the other hand, let us justify some of the other paths where the errors of the Deconfounder are greater than unaware models. In Sachs’ dataset to model the causal effect  $E_{kt} \rightarrow A_{kt}$ , the factorization model of the deconfounder uses  $R_{af}$ ,  $M_{ek}$ ,  $J_{nk}$  and  $P_{38}$  to extract the substitute confounder; the factorization model assumes that all those variables are independent conditioned to  $\tilde{z}$ , while that is not the case in the true SCM and, therefore, this SCM violates the independence assumption of (Wang & Blei, 2019). The same argument is valid for the paths  $y_{ceP} \rightarrow y_{faD}$ ,  $lacA \rightarrow y_{aeM}$ ,  $y_{ceP} \rightarrow y_{faD}$ ,  $y_{deE} \rightarrow p_{spA}$  and  $p_{spB} \rightarrow p_{spA}$ .

On the other hand, the paths  $lacZ \rightarrow y_{aeM}$ ,  $asnA \rightarrow lacY$  are frontdoor paths that Decaf can identify because it models the hidden confounder following the true causal graph. However, the Deconfounder is not designed to model this paths. To evaluate its performance for frontdoor paths, Deconfounder uses the same variables as Decaf to extract the substitute of the confounder. However, the Deconfounder assumes independence conditioned to the substitute confounder and that is not the case; therefore, we are violating the independence assumption again.

The only two paths that meet the Deconfounder assumptions in Fig. 1 are  $lacA \rightarrow lacY$  and  $y_{deE} \rightarrow p_{spB}$ . And we can observe that in those paths, the Deconfounder performs at least as well as unaware methods. On the other hand, all the factor models used for the Deconfounder implementation (PPCA, Deep exponential families and Variational autoencoder) assume additive noise. Therefore, interventional distributions in nonadditive settings are not computable theoretically with these models.

Table 4: Performance metrics on Ecoli70 dataset. Statistics computed an all samples over 5 runs, intervening and evaluating only in the causal effects that Deconfounder should solve. Bold indicates significantly better results (95% CI from a Mann-Whitney U test). Lower error values indicate better performance.

Model	ATE err  $\times 10^2$	CF err  $\times 10^1$
Oracle CNF	8.31 <sub>10.95</sub>	1.49 <sub>1.86</sub>
Aware Decaf 	<b>9.18</b> <sub><b>10.42</b></sub>	<b>2.18</b> <sub><b>2.02</b></sub>
Aware Deconfounder	14.35 <sub>15.24</sub>	12.03 <sub>15.81</sub>
Unaware CNF	27.82 <sub>30.17</sub>	4.01 <sub>3.62</sub>
Unaware ANM	27.63 <sub>29.74</sub>	3.64 <sub>3.15</sub>
Unaware DCM	42.45 <sub>54.23</sub>	4.08 <sub>4.12</sub>

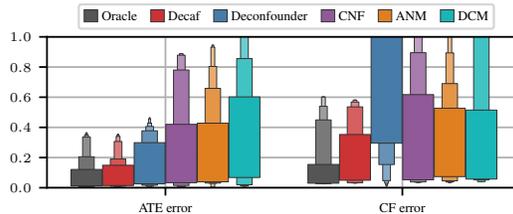


Figure 16: ATE and CF error evaluating only links where deconfounder should work in the additive case.

B.3.2 METRICS ON THE OTHER PATHS

In this subsection we include a comparison between all the models in the *unconfounded* and the *unidentifiable* effects. For *unconfounded effects*, our expectation is to observe that all the CGMs achieve a performance comparable with the *oracle*. On the other hand, we expect to have higher errors in *unidentifiable effects*, since we do not have theoretical guarantees.

**Unconfounded Effects.** The results for *unconfounded effects* are summarized in Fig. 17 and Tab 5, considering only direct effects for ATE and counterfactual error computations. As expected, Decaf and CNF achieve metrics comparable to the *oracle* in both ATE and counterfactual estimations, particularly evident in Fig. 17, where error distributions are nearly identical. 5 does not show statistically significant differences between Decaf and CNF. Notably, architectures based on causal normalizing flows outperform ANM and DCM, which model each causal mechanism,  $f_i$ , with separate networks. This difference is crucial in settings with many variables and complex relations, where scalability is essential. Unlike ANM and DCM, which suffer from error propagation and limited scalability, causal normalizing flows leverage a single amortized model, making them more efficient in high-dimensional scenarios.

Finally, note that the Deconfounder has not been included in these metrics because it is not designed for *unconfounded queries* and there are many queries, while one Deconfounder model is needed for each query.

Table 5: Performance metrics on Ecoli70 dataset. Statistics computed on all *unconfounded* direct effects and 5 runs. Bold indicates significantly better results (95% CI from a Mann-Whitney U test). Lower error values indicate better performance.

		Additive			Nonadditive		
	Model	MMD int $\times 10^4$	ATE err  $\times 10^2$	CF err  $\times 10^2$	MMD int $\times 10^4$	ATE err  $\times 10^2$	CF err  $\times 10^2$
	Oracle						
	CNF	3.72 <sub>3.73</sub>	2.00 <sub>2.27</sub>	1.27 <sub>3.49</sub>	1.94 <sub>2.96</sub>	1.92 <sub>1.99</sub>	1.76 <sub>4.10</sub>
Aware	Decaf	4.59 <sub>5.58</sub>	2.11 <sub>2.39</sub>	1.42 <sub>3.81</sub>	2.76 <sub>7.61</sub>	1.93 <sub>1.87</sub>	1.75 <sub>4.04</sub>
	CNF	4.77 <sub>6.09</sub>	2.02 <sub>2.21</sub>	1.22 <sub>3.18</sub>	2.97 <sub>7.64</sub>	1.95 <sub>1.92</sub>	1.71 <sub>3.93</sub>
Unaware	ANM	34.72 <sub>8.56</sub>	3.57 <sub>3.02</sub>	2.02 <sub>4.09</sub>	15.13 <sub>12.57</sub>	3.53 <sub>3.15</sub>	2.64 <sub>5.34</sub>
	DCM	36.23 <sub>14.29</sub>	3.48 <sub>2.75</sub>	2.69 <sub>2.30</sub>	21.22 <sub>13.68</sub>	3.42 <sub>2.63</sub>	3.00 <sub>3.42</sub>

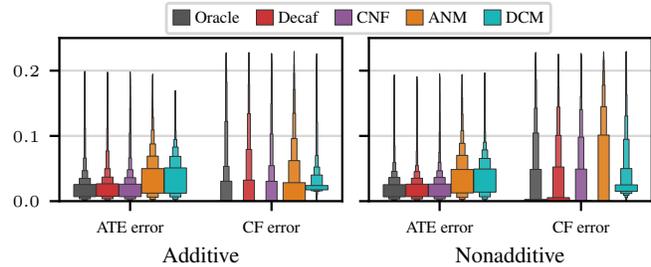


Figure 17: Error boxenplots on the Ecoli70 dataset for different CGMs, averaged over all *unconfounded* direct effects (see Fig. 1) after intervening in their 25th, 50th, and 75th percentiles and 5 random realizations of the experiment.

**Unidentifiable Effects.** The results for *unidentifiable effects*—causal queries that violate the assumptions in §6—are summarized in Fig. 18 and Tab 6. Notably, the *oracle* performs significantly better than the other CGMs. As seen in Fig. 18, error distributions are highly skewed, with ATE and counterfactual errors reaching extreme values—considering that metrics are computed on the standardized variables. Tab 6 shows no significant differences between the metrics achieved by Decaf and CNF.

B.4 Law school fairness use-case

The experiment with real-world data was inspired by Kusner et al. (2017) and Javaloy et al. (2023).

The purpose is to find a fair estimator of the decile that the grades of each student will occupy in their third year

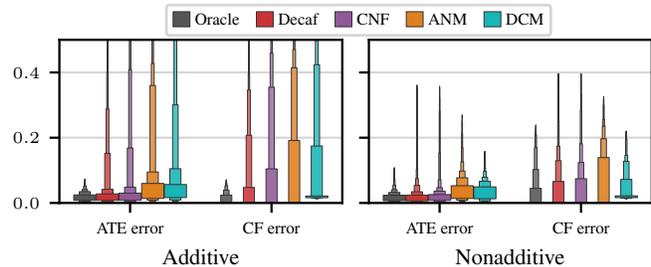


Figure 18: Error boxenplots on the Ecoli70 dataset for different CGMs, averaged over all *unidentifiable* direct effects (see Fig. 1) after intervening in their 25th, 50th, and 75th percentiles and 5 random realizations of the experiment.

Table 6: Performance metrics on Ecoli70 dataset. Statistics computed on all **unidentifiable** direct effects and 5 runs. Bold indicates significantly better results (95% CI from a Mann-Whitney U test). Lower error values indicate better performance

Model	Additive			Nonadditive		
	MMD int $\times 10^4$	ATE err  $\times 10^2$	CF err  $\times 10^3$	MMD int $\times 10^5$	ATE err  $\times 10^2$	CF err  $\times 10^2$
Oracle CNF	3.71 <sub>3.52</sub>	1.79 <sub>1.36</sub>	5.88 <sub>15.16</sub>	16.98 <sub>6.87</sub>	1.75 <sub>1.59</sub>	1.62 <sub>4.57</sub>
Aware Decaf 	3.83 <sub>3.93</sub>	3.76 <sub>7.90</sub>	31.21 <sub>92.38</sub>	18.67 <sub>5.64</sub>	2.18 <sub>3.73</sub>	2.11 <sub>6.13</sub>
Unaware CNF	4.54 <sub>4.81</sub>	4.75 <sub>10.65</sub>	44.76 <sub>126.36</sub>	20.22 <sub>6.68</sub>	2.32 <sub>3.80</sub>	2.13 <sub>6.25</sub>
ANM	34.38 <sub>5.17</sub>	7.43 <sub>12.64</sub>	52.70 <sub>137.99</sub>	130.71 <sub>41.64</sub>	4.01 <sub>3.82</sub>	2.93 <sub>7.21</sub>
DCM	35.49 <sub>4.95</sub>	7.67 <sub>13.93</sub>	67.46 <sub>132.21</sub>	198.23 <sub>58.62</sub>	3.43 <sub>2.76</sub>	3.29 <sub>3.92</sub>

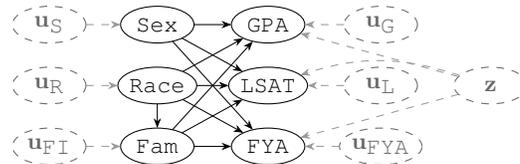


Figure 19: Confounded SCM modeled by Decaf.

of university.

The dataset contains information on 27,000 law students who were admitted by the Law School Admissions Council (LSAC) from 1991 through 1997. We have performed an experiment similar to that carried out by [Kusner et al. \(2017\)](#), where race and sex were treated as sensitive attributes. We have considered the following variables to include in our study:

- **Race**: binary indicator of the race that distinguish between white and non-white.
- **Sex**: binary indicator of the sex that distinguish between male and female.
- **Fam**: family income.
- **LSAT**: the grade achieved in the Law School Admission Test (LSAT).
- **UGPA**: the undergraduate grade point average (GPA) of the student previous to the admission.
- **FYA**: first-year average grade.
- **Decile3**: the decile of the grades in the third year of university. This is the variable to predict.

For our purpose, we consider that an estimator,  $\hat{y}$ , is fair if it meets *Demographic parity*, defined in ([Kusner et al., 2017](#), Def. 3) as follows. A predictor  $\hat{y}$  satisfies demographic parity if the predicted distributions for different values of a sensitive attribute are equal:  $p(\hat{y} | t = 0) = p(\hat{y} | t = 1)$ . We evaluate the difference between predicted distributions using MMD—a lower distance between the predictions for the two groups of a sensitive attributes denotes a fairer predictor.

The assumed causal graph is slightly different from that of [Kusner et al. \(2017\)](#), since their purpose is to make a fair prediction FYA accounting only for Race, Sex, LSAT and UGPA. However, we include Fam and FYA as predictors and the task is to predict Decile3 and the assumed causal graph is the one of [Fig. 8](#).

**Proposed fair predictor with Decaf.** We propose to model the confounded SCM presented in [Fig. 19](#), where are explicitly shown the exogenous variables, that are independent of the other variables of the graph except of their associated endogenous variable.

Afterwards, we predict the outcome, Decile3 from the extracted latent variable that acts as substitute of the knowledge and the exogenous variables of FYA and Fam, following the causal graph of [Fig. 8](#), using a generalized linear model:  $\tilde{p}(\text{Decile3} | \mathbf{u}_{FI}, \mathbf{u}_{FYA}, \mathbf{z})$ . Decaf models  $\mathbf{z}$  and the exogenous variables as independent from Race and Sex. Therefore, the prediction of Decile3 should be fair.

**Baselines.** The baselines used to compare our approach are the methods *Fair K* and *Fair add* proposed in [Kusner et al. \(2017\)](#).

1430 *Fair K* is a fair predictor categorized in Level 2 in [Kusner et al. \(2017\)](#), which postulates that the student’s knowledge, `know`  
 1431 affects GPA, LSAT, FYA and Decile 3, following the distributions described below.

$$\begin{aligned}
 1432 & \text{Fam} \sim \mathcal{N}(b_{\text{Fam}} + w_{\text{FamRace}}^R, 1), \\
 1433 & \text{GPA} \sim \mathcal{N}(b_G + w_G^K \text{know} + w_G^R \text{Race} + w_G^S \text{Sex} + w_G^{\text{Fam}} \text{Fam}, \sigma_G^2), \\
 1434 & \text{LSAT} \sim \text{Poisson}(\exp(b_L + w_L^K \text{know} + w_L^R \text{Race} + w_L^S \text{Sex} + w_L^{\text{Fam}} \text{Fam})), \\
 1435 & \text{FYA} \sim \mathcal{N}(w_F^K \text{know} + w_F^R \text{Race} + w_F^S \text{Sex} + w_F^{\text{Fam}} \text{Fam}, 1), \\
 1436 & \text{Decile3} \sim \text{Poisson}(\exp(w_D^K \text{know} + w_D^R \text{Race} + w_D^S \text{Sex} + w_D^{\text{Fam}} \text{Fam})), \\
 1437 & \text{know} \sim \mathcal{N}(0, 1).
 \end{aligned} \tag{50}$$

1441 Then, the posterior distribution `know` is inferred using MonteCarlo with the probabilistic programming language Pyro  
 1442 ([Bingham et al., 2019](#)). The outcome is predicted using the inferred `know` using a generalized linear model:  $\hat{p}(\text{Decile3} |$   
 1443 `know`).

1444 On the other hand, *Fair Add* predicts the outcome from the residuals of predicting each variable with each parent, which  
 1445 guarantees that these residuals are independents of Race and Sex. That is, the predictor estimates the distribution  
 1446  $p(\text{Decile3} | \mathbf{r}_{\text{Fam}}, \mathbf{r}_{\text{UGPA}}, \mathbf{r}_{\text{LSAT}}, \mathbf{r}_{\text{FYA}})$ , where these residuals are computed as:

$$\begin{aligned}
 1447 & \mathbf{r}_{\text{Fam}} = \text{Fam} - \mathbb{E}[\text{Fam} | \text{Sex}, \text{Race}] \\
 1448 & \mathbf{r}_{\text{UGPA}} = \text{UGPA} - \mathbb{E}[\text{GPA} | \text{Sex}, \text{Race}, \text{Fam}] \\
 1449 & \mathbf{r}_{\text{LSAT}} = \text{LSAT} - \mathbb{E}[\text{LSAT} | \text{Sex}, \text{Race}, \text{Fam}] \\
 1450 & \mathbf{r}_{\text{FYA}} = \text{FYA} - \mathbb{E}[\text{FYA} | \text{Sex}, \text{Race}, \text{Fam}]
 \end{aligned} \tag{51}$$

1451 All predictors used are generalized linear models.

1452 **Discussion of Results.** Although the *fair* methods proposed by [Kusner et al. \(2017\)](#) achieve significantly better *demo-*  
 1453 *graphic parity* than our approach using Decaf (as indicated by a much lower MMD), their predictive performance is  
 1454 substantially inferior. Specifically, their performance is comparable to predicting the outcome using only the mean of the  
 1455 distribution, which serves as a baseline in [Tab 1](#). In contrast, Decaf achieves a 98% reduction in MMD while incurring only  
 1456 an 11% increase in RMSE, as illustrated in [Fig. 9](#).

1457 These experiments demonstrate that leveraging Decaf to model confounded Structural Causal Models is beneficial beyond  
 1458 causal query estimation, leading to improved overall performance.

## 1459 C Do-operator

1460 We introduce here the algorithms that Decaf employ to generate interventional samples and counterfactuals. But first, we  
 1461 include those of [Javaloy et al. \(2023\)](#), since we leverage these CNFs as building blocks for Decaf. Note that the notation  
 1462 applied for Decaf is slightly different from the that used in the causal flows, naming the intervened variable as `t`, instead of  
 1463 `xi`, in order to be consistent with the notation used in [§3](#) and [§6](#). However, note that both variables play the same role, and  
 1464 that `t`  $\subset$  `x`.

### 1465 C.1 Do-operator in causal normalizing flows

1466 **Algorithm 1** Algorithm to sample from the interventional distribution,  $P(\mathbf{x} | \text{do}(x_i = \alpha))$ . From [Javaloy et al. \(2023\)](#).

```

1467 1: function SAMPLEINTERVENEDDIST(i,  $\alpha$ )
1468 2:    $\mathbf{u} \sim P_{\mathbf{u}}$  ▷ Sample a value from the observational distribution.
1469 3:    $\mathbf{x} \leftarrow T_{\theta}^{-1}(\mathbf{u})$ 
1470 4:    $x_i \leftarrow \alpha$  ▷ Set  $x_i$  to the intervened value  $\alpha$ .
1471 5:    $\mathbf{u}_i \leftarrow T_{\theta}(\mathbf{x})_i$  ▷ Change the  $i$ -th value of  $\mathbf{u}$ .
1472 6:    $\mathbf{x} \leftarrow T_{\theta}^{-1}(\mathbf{u})$ 
1473 7:   return  $\mathbf{x}$  ▷ Return the intervened sample.
1474 8: end function

```

1485 The computation of counterfactuals follows the steps of *abduction*, *action* and *prediction* postulated by (Pearl et al., 2016).  
 1486 The *abduction* step consists of using the observations to determine the value of the exogenous variables. Then, *action* is  
 1487 computing the intervention, modifying the causal mechanism of the intervened variable and *prediction* consist of using the  
 1488 exogenous variables and the modified SCM to compute the counterfactual.

---

1491 **Algorithm 2** Algorithm to sample from the counterfactual distribution,  $P(\mathbf{x}^{\text{cf}} \mid \text{do}(x_i = \alpha), \mathbf{x}^{\text{f}})$ . From Javaloy et al. (2023).

```

1492 1: function GETCOUNTERFACTUAL( $\mathbf{x}^{\text{f}}, i, \alpha$ )
1493 2:    $\mathbf{u} \leftarrow T_{\theta}(\mathbf{x}^{\text{f}})$                                 ▷Abduction: Get  $\mathbf{u}$  from the factual sample.
1494 3:    $x_i^{\text{f}} \leftarrow \alpha$                                 ▷Action: Set  $x_i$  to the intervened value  $\alpha$ .
1495 4:    $\mathbf{u}_i \leftarrow T_{\theta}(\mathbf{x}^{\text{f}})_i$                             ▷Action: Change the  $i$ -th value of  $\mathbf{u}$ .
1496 5:    $\mathbf{x}^{\text{cf}} \leftarrow T_{\theta}^{-1}(\mathbf{u})$                         ▷Prediction: Get counterfactual
1497 6:   return  $\mathbf{x}^{\text{cf}}$                                         ▷ Return the counterfactual value.
1498 7: end function
    
```

---

### 1502 C.2 Do-operator in interventional distributions with Decaf

1503 The sampling process consists of sampling first from the prior distribution of the latent variables and from the distribution of  
 1504 the exogenous variables. Then, one can use the generative network ( $T_{\theta}$ ) to take samples of the rest of variables, changing  
 1505 the components of  $\mathbf{u}$  associated with  $t$ . Note that  $\mathbf{z}$  is not the input of the normalizing flow, but a condition (or *context*).  
 1506 Therefore,  $\mathbf{z}$  is transformed neither in the forward nor the reverse pass of the flow.

---

1509 **Algorithm 3** Algorithm to sample from the interventional distribution,  $P(\mathbf{x} \mid \text{do}(t = \alpha))$  with Decaf.

```

1510 1: function SAMPLEINTERVENEDDIST( $t, \alpha$ )
1511 2:    $\mathbf{z} \sim P_{\mathbf{z}}$                                             ▷ Sample a value from the prior of  $\mathbf{z}$ .
1512 3:    $\mathbf{u} \sim P_{\mathbf{u}}$                                             ▷ Sample a value from the observational distribution.
1513 4:    $\mathbf{x} \leftarrow T_{\theta}^{-1}(\mathbf{u}, \mathbf{z})$ 
1514 5:    $t \leftarrow \alpha$                                         ▷ Set  $t$  to the intervened value  $\alpha$ .
1515 6:    $\mathbf{u}_t \leftarrow T_{\theta}(\mathbf{x}, \mathbf{z})_t$                             ▷ Change the component of  $\mathbf{u}$  associated with  $t$ .
1516 7:    $\mathbf{x} \leftarrow T_{\theta}^{-1}(\mathbf{u}, \mathbf{z})$ 
1517 8:   return  $\mathbf{x}$                                             ▷ Return the intervened sample.
1518 9: end function
    
```

---

1522 Additionally, the process to compute the average treatment effect (ATE) involves to generate interventional distributions.  
 1523 For example, to compute the ATE comparing two interventions ( $\alpha_1, \alpha_2$ ) in the variable  $t$ , we would generate samples of the  
 1524 interventional distributions,  $p(\mathbf{x} \mid \text{do}(t = \alpha_1)), p(\mathbf{x} \mid \text{do}(t = \alpha_2))$ , respectively, and approximate their expectations with  
 1525 MonteCarlo.

$$1528 \text{ATE} = \mathbb{E}[\mathbf{x} \mid \text{do}(t = \alpha_2)] - \mathbb{E}[\mathbf{x} \mid \text{do}(t = \alpha_1)] \quad (52)$$

1530 Unfortunately, if we were interested in evaluating the ATE on only one variable,  $y$ , the process would involve to generate  
 1531 samples of the whole interventional distribution and select only the samples of the interested variable.

### 1533 C.3 Do-operator in counterfactuals with Decaf

1535 As part of the abduction step, our model estimates the posterior distribution of hidden confounders given a factual datapoint,  
 1536  $q_{\phi}(\mathbf{z} \mid \mathbf{x}^{\text{f}})$ . Therefore, the counterfactual given by the model is no longer a single point but comes from a distribution. To  
 1537 obtain a single sample that allows us to compare it with the true counterfactual, we estimate the mean of the posterior  
 1538 distribution from samples using MonteCarlo, and we use that value to generate a counterfactual.

1540 **Algorithm 4** Algorithm to sample from the interventional distribution,  $P(\mathbf{x} \mid \text{do}(t = \alpha))$  with Decaf.

```

1541 1: function GETCOUNTERFACTUAL( $\mathbf{x}^f, t, \alpha$ )
1542 2:    $q_\phi(\mathbf{z} \mid \mathbf{x}^f) \leftarrow$  Deconfounding network( $\mathbf{x}^f$ )
1543 3:    $\mathbf{z} \sim \mathbb{E}_{q_\phi}[\mathbf{z} \mid \mathbf{x}^f]$ 
1544 4:    $\mathbf{u} \leftarrow T_\theta(\mathbf{x}^f, \mathbf{z})$ 
1545 5:    $t^f \leftarrow \alpha$ 
1546 6:    $\mathbf{u}_t \leftarrow T_\theta(\mathbf{x}^f, \mathbf{z})_t$ 
1547 7:    $\mathbf{x}^{cf} \leftarrow T_\theta^{-1}(\mathbf{u}, \mathbf{z})$ 
1548 8:   return  $\mathbf{x}^{cf}$ 
1549 9: end function

```

$\triangleright$  **Abduction:** Get  $\mathbf{z}$  from the factual sample.  
 $\triangleright$  **Abduction:** Estimate the mean of the distribution.  
 $\triangleright$  **Abduction:** Get  $\mathbf{u}$  from the factual sample.  
 $\triangleright$  **Action:** Set  $t$  to the intervened value  $\alpha$ .  
 $\triangleright$  **Action:** Change the component of  $\mathbf{u}$  associated with  $t$ .  
 $\triangleright$  **Prediction:** compute the counterfactual  
 $\triangleright$  Return the counterfactual value.

## 1552 D Additional details on related work of causal inference with hidden confounders

1553 In this section we go deeper into the methods of causal inference in scenarios where there are unobserved confounders.  
1554 First of all, we want to remark that all the following methods have been designed to address causal inferences in specific  
1555 causal graphs (or subgraphs), therefore they can be used when there exists the causal relationships presented in Fig. 20.

1556 In the following text, we assume the notation introduced in §3, where  $\mathbf{z}$  is the hidden confounder,  $t$  is the intervened variable  
1557 or treatment and  $y$  is the outcome, i.e. the variable where we want to evaluate the causal effects.

1558 We have classified the different approaches depending on the graph that they are designed to address. However, there are  
1559 two considerations that are common for all these approaches.

1560 First, the methods follow a two-stage process: **i)** extracting a substitute for the unobserved confounder,  $\tilde{\mathbf{z}}$ , using the variables  
1561 affected by the confounder or instrumental variables, and **ii)** estimating the outcome given this substitute,  $\tilde{y} \sim p(y \mid \tilde{\mathbf{z}}, t)$ . In  
1562 larger causal graphs, one predictor must be trained for each outcome, and one extractor must be trained per independent  
1563 confounder.

1564 Second, none of these methods shows the ability of identify *counterfactuals*, since they do not model exogenous variables.

1565 **Presence of null proxies independent of  $t$  (Fig. 20a).** We say  $\mathbf{n}$  to be a null proxy of  $\mathbf{z}$  if it is a child of  $\mathbf{z}$  independent  
1566 of the outcome,  $y$ , given  $\mathbf{z}$ :  $\mathbf{n} \perp\!\!\!\perp y \mid \mathbf{z}$ . Methods for estimating causal effects were developed when null proxies of the  
1567 confounder were available and those proxies are independent of the intervened variable:  $\mathbf{n} \perp\!\!\!\perp t \mid \mathbf{z}$ . We can use these proxies  
1568 to infer a substitute. Among these, Allman et al. (2009); Kuroki & Pearl (2014) studies the case in which the confounder is  
1569 categorical and uses matrix factorization to extract a substitute when there are at least three Gaussian proxies (Allman et al.,  
1570 2009), when the conditional distribution of the confounder given the proxy is known or when other proxies are available  
1571 (Kuroki & Pearl, 2014). Kallus et al. (2018) also employ matrix factorization to cases where the confounder is continuous  
1572 and the relation with the covariates and the treatment (but not with the outcome) is linear. In addition, Kallus et al. (2019)  
1573 uses kernel functions to extract the substitute confounder when the generators are nonlinear. The most relevant method based  
1574 on deep generative methods is proposed by Louizos et al. (2017), consisting of a VAE to extract the substitute confounder  
1575 when several null proxies are available, although there is no theoretical guarantee of its operation. Finally, Miao et al. (2023)  
1576 offers a regression-based approach to estimate the unobserved confounder under *equivalence*, which assumes that any model  
1577 of the joint achieves element-wise transformations of the latents, which is not feasible to check:  $\tilde{p}(t, \mathbf{z} \mid \mathbf{n}) = p(t, V(\mathbf{z}) \mid \mathbf{n})$ .  
1578 The graph in which all these methods operate can be found in Fig. 20a.

1579 **Presence of two proxies: null and not null (Fig. 20b).** When the null proxies affect treatment (see Fig. 20b: the proxy,  
1580  $\mathbf{n}$ , affects treatment  $t$ ), Miao et al. (2018) offers theoretic guarantees of causal identifiability in the presence of another  
1581 proxy,  $\mathbf{w}$ , and completeness conditions. The proxy  $\mathbf{w}$  can be active, that is, it can directly affect  $y$ . Practically, in Tchetgen  
1582 et al. (2020) the two-stage proximal least squares (P2SLS) we can find the method to infer the substitute confounder from  
1583  $p(\mathbf{w} \mid t, \mathbf{n})$ . P2SLS can be implemented using neural networks to achieve greater flexibility.

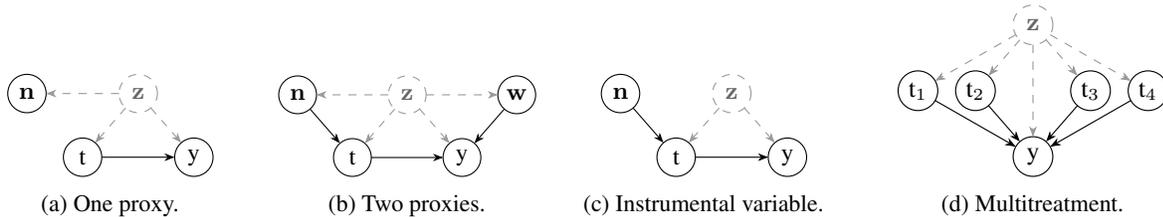
1584 **Instrumental variable (Fig. 20c).** Another condition that allows causal inference is the presence of instrumental variables  
1585 (IVs), i.e. variables that affect only the treatment and are independent of both the unobserved confounder and the outcome  
1586 given the treatment (in Fig. 20c,  $\mathbf{n}$  is an IV). In linear DGP, Pearl (2009); Angrist & Pischke (2009) demonstrates that a  
1587 two-stage regression process mitigates the confounding bias as the only effect that flows from the instrumental variable to  
1588 the outcome is through treatment. A substitute of the confounder is extracted by computing the conditional distribution of  
1589 the treatment given the instrumental variable:  $\tilde{\mathbf{z}} \sim p(t \mid \mathbf{n})$ . Furthermore, (Hartford et al., 2017) develops an extension of  
1590 this theory to include arbitrarily complex nonlinear DGP, designing a two-step deep approach, based on neural networks.

1595 **Multitreatment affected by a common confounder (Fig. 20d).** Finally, the multitreatment scenario (Fig. 20d) has been  
 1596 studied by Wang & Blei (2019); Ranganath & Perotte (2018). It is called multitreatment because all covariates can be seen as  
 1597 a treatment over the outcome,  $y$ . It is assumed that, in the true DGP, there exist several covariates that are independent given  
 1598 the unobserved confounder. Therefore Wang & Blei (2019) propose to use a factorization model, such as probabilistic PCA  
 1599 or Poisson Matrix Factorization, to infer the substitute confounder. A factorization model assumes that the distribution of all  
 1600 the treatments factorizes in the following way:  $p(\mathbf{t}, \mathbf{z}) = p(\mathbf{z}) \prod_{i=1}^d p(t_i | \mathbf{z})$ , which should allow to construct a substitute  
 1601 of the confounder from the posterior of  $\mathbf{z}$ :  $\tilde{\mathbf{z}} \sim \tilde{p}(\tilde{\mathbf{z}} | \mathbf{t})$ . However, these method only offers identifiability in the asymptotic  
 1602 setting where the number of treatments is infinite. On the other hand Ranganath & Perotte (2018) proposes a method that  
 1603 uses a VAE as a factorization model, adding a regularization term to reduce the additional mutual information between the  
 1604 estimated confounder and treatment  $t_j$  given the rest of treatments  $\mathbf{t}_{-j}$ . Again, the theoretical guarantees of this approach  
 1605 need an infinite number of treatments to achieve unbiased estimates of causal effects.

1606 Wang & Blei (2021) connect the ideas of Miao et al. (2018) and Wang & Blei (2019) ensuring causal identification in the  
 1607 multitreatment setting when it is known that some of the treatments *can act as null proxies*, that is, they do not affect the  
 1608 outcome. This assumption allows them to provide theoretical guarantees when the number of treatments does not tend to be  
 1609 infinite.

1611 **How is Deconfounder Wang & Blei (2019; 2021) related to our work.** As Decaf does, Deconfounder infers the posterior  
 1612 distribution of the substitute of the confounder from the observational data using a generative model. However, the  
 1613 application of a factorization model restricts the structural dependencies that we can model. For example, the Deconfounder  
 1614 cannot model the structural dependencies of Fig. 20b, since the factorization model assumes  $\mathbf{n} \perp\!\!\!\perp \mathbf{t} \perp\!\!\!\perp \mathbf{w} | \mathbf{z}$ . In contrast, the  
 1615 Decaf uses a causal flow, which does allow this dependencies because the causal graph is encoded in the flow.

1616 We also stress that Decaf models the whole confounded SCM, including the exogenous variables. This allows to compute  
 1617 *counterfactuals* and train in a query-agnostic manner. In contrast, Deconfounder cannot compute counterfactuals and needs  
 1618 of a separate model per query.



1628 Figure 20: *Ad-hoc* graphs. (a) Kuroki & Pearl (2014); Louizos et al. (2017); Miao et al. (2023); Kallus et al. (2018; 2019);  
 1629 Allman et al. (2009) address the case where  $\mathbf{n}$  is independent of  $\mathbf{t}$ . (b) Miao et al. (2018) is designed for the case where there  
 1630 exist two proxies. (c) Graph with an instrumental variable, but this graph is out of the scope of our framework. (d) Wang &  
 1631 Blei (2019; 2021); Ranganath & Perotte (2018) are designed for the multitreatment setting.

## 1634 E Algorithms for causal query identification

1635 As explained in §6.3, we can ask Decaf to solve any causal query, but we do not have the guarantee that the estimation that  
 1636 Decaf returns is correct unless the query is identifiable. Therefore, we provide the practitioner with algorithms to check the  
 1637 identifiability of causal queries.

1639 **Specific treatment-outcome pair.** We start presenting the Alg. 5 to identify a causal query specifying the pair treatment  
 1640 and outcome, which is valid for estimating the interventional distribution of the outcome— $p(y | \text{do}(t), \mathbf{c})$ —and the  
 1641 counterfactual— $p(y^{\text{cf}} | \text{do}(t), \mathbf{x}^{\text{f}})$ —, since we postulated in §6 that the latter is identifiable if the former is.

1643 We have employed this algorithm in all the paths of Sachs and Ecoli70 datasets to check the identifiability of all the direct  
 1644 causal effects—where  $y$  is a child of  $t$ —, in order to get a visual representation of the identifiable queries of a complex graph.  
 1645 However, due to the large number of possible causal queries resulting from all edge combinations in the 43-node Ecoli70  
 1646 dataset, we have not analyzed identifiability for all indirected queries.

1647 Trivially, if one is interested in evaluating a query which involves several outcomes,  $\{y_1, y_2, \dots, y_O\}$ , one causal query per  
 1648  $y_i$  should be evaluated.

1650 **Algorithm 5** Identification of causal queries that include intervention and outcome  $(t, y)$

---

1651 **Require:** Graph  $\mathcal{G}$ , intervention variable  $t$ , outcome variable  $y$ , covariates  $c$ , hidden variables  $z$

1652 **Ensure:** Boolean indicating if query is identifiable

1653 1:  $z \leftarrow$  hidden variables that are parents of both  $t$  and  $y$

1654 2: **return** True if  $z$  is  $\emptyset$  ▷ Unconfounded is identifiable

1655 3: **for all**  $z_k \in z$  **do**

1656 4:   **Comment:** Each  $z_k$  is an independent component of  $z$

1657 5:    $n$ -proxies  $\leftarrow$  children of  $z_k$   $d$ -separated from  $t$  given  $(z, c)$

1658 6:    $w$ -proxies  $\leftarrow$  children of  $z_k$   $d$ -separated from  $y$  given  $(z, c)$

1659 7:   **if** there exist  $n \in n$ -proxies and  $w \in w$ -proxies such that  $n$  is  $d$ -separated from  $w$  given  $(z, c)$  **then**

1660 8:      $z_k$  is deconfounded

1661 9:   **end if**

1662 10: **end for**

1663 11: **return** all  $z_k$  are deconfounded

---

1666 **Evaluation on all the variables.** Although the Alg. 6 consist of applying Alg. 5 iteratively, we also find it interesting to include the extension to identify causal queries evaluated on all variables in the dataset, which is useful for using Decaf as a generative model for the interventional distribution— $p(\mathbf{x} \mid \text{do}(t))$ —, or offering complete counterfactual samples— $p(\mathbf{x}^{\text{cf}} \mid \text{do}(t), \mathbf{x}^{\text{f}})$ —intervening in a specific variable,  $t \subset \mathbf{x}$ .

1671 **Algorithm 6** Identification of causal queries, intervening in  $t$  and evaluating in all variables

---

1672 **Require:** Graph  $\mathcal{G}$ , intervention variable  $t$ , hidden variables  $z$

1673 **Ensure:** Boolean indicating if the interventional distribution is identifiable

1674 1:  $z \leftarrow$  hidden variables that are parents of  $t$

1675 2: **for all**  $x_i \in$  descendants of  $t$  **do**

1676 3:   **Comment:** Evaluate only on descendants of the intervention

1677 4:   Check  $(t, x_i)$  identifiability with Alg. 5

1678 5: **end for**

1679 6: **return** all  $(t, x_i)$  are identifiable

---

### 1682 E.1 Pipeline for using Decaf

1683 Our framework provides a systematic approach to solving causal queries by integrating Decaf, a model trained on observational data, with algorithms designed for query identifiability analysis.

1684 As depicted in the pipeline, the framework takes as input a dataset  $\mathcal{D}$ , a causal graph  $\mathcal{G}$ , and a set of  $N$  interesting queries  $\{Q_i\}_{i=1}^N$ . The process begins by training Decaf on  $\mathcal{D}$  and  $\mathcal{G}$ , enabling it to learn the confounded SCM,  $\mathcal{M}$ .

1691 Simultaneously, the identifiability of each causal query  $Q_i$  is assessed using dedicated algorithms (Alg. 5 and Alg. 6). If  $Q_i$  is identifiable, the trained Decaf is used to estimate  $Q_i(\mathcal{M})$  (Alg. 3 and Alg. 4), yielding the estimated causal effect  $\hat{Q}_i(\mathcal{M})$ . If  $Q_i$  is not identifiable, the framework indicates that answering the query is not feasible given the available data and causal structure. Other causal queries can be answered by the model without retraining, provided that their identifiability is verified beforehand.

1699 This workflow ensures a principled approach to causal inference, leveraging both data-driven modelling and theoretical guarantees on identifiability. Both the Decaf model and the algorithms for query identifiability and estimation will be included in the code that we will provide upon acceptance.

1704

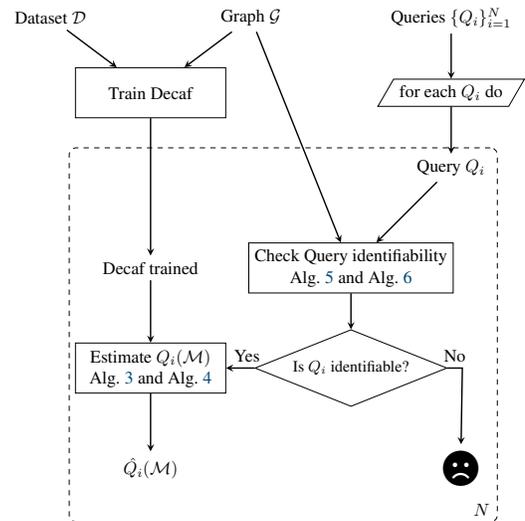


Figure 21: **Block diagram of the pipeline.**

1705 **Validation with interventional data.** As a final step in the pipeline in real-world scenarios, especially in sensitive  
1706 applications, we encourage practitioners to validate the framework with interventional data. Causal queries such as *average*  
1707 *treatment effects* (ATEs) can be validated if a randomized experiment is available in which interventions are carried out on  
1708 the treatment variable.

1709 However, in cases where experiments on the required variable are not available, our framework can still be partially validated  
1710 by assessing the completeness of the inferred hidden confounder given the observed proxies. This can be done by evaluating  
1711 causal effects in another causal query that shares the same hidden confounder. Specifically, if a causal query  $Q_1$  lacks  
1712 interventional data, but another query  $Q_2$  involving the same hidden confounder is identifiable, the inferred confounder  
1713 from  $Q_2$  can be postulated as a valid substitute for estimating  $Q_1$ . This indirect validation method provides a way to assess  
1714 the reliability of our framework without requiring direct interventions for every confounded query.  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759