

Appendices

Table of Contents

A	Extended Discussion of Related Work	20
A.1	Comparison of Related Identifiability Results	20
B	Proof of Minimality of the CRL Equivalence Class \sim_{CRL}	22
C	Identifiability Proofs	23
C.1	Auxiliary Lemmata	23
C.2	Proof of Thm. 3.2	24
C.3	Proof of Thm. 3.4	29
C.4	Proof of Thm. 4.2	31
D	Experimental Details and Additional Results	33
D.1	Experimental Details for § 6	33
D.2	Additional Results: Learning Nonlinear Latent SCMs from Partial Causal Order	34
E	Discussion of the Role of Our Assumptions	36

A Extended Discussion of Related Work

Prior work on causal representation learning with general nonlinear relationships (both among latents and between latents and observations) and without an explicit task or label typically relies on some form of *weak supervision*. One example of weak supervision is “multi-view” data consisting of tuples of related observations. von Kügelgen et al. [123] consider counterfactual pairs of observations arising from imperfect interventions through data augmentation, and prove identifiability for the invariant non-descendants of intervened-upon variables. Brehmer et al. [18] also use counterfactual pre- and post-intervention views and show that the latent SCM can be identified given all single-node perfect stochastic interventions. Another type of weak-supervision is temporal structure [2], possibly combined with nonstationarity [132, 133], interventions on known targets [80, 81], or observed actions inducing sparse mechanism shifts [74, 75, 110]. Other works use more explicit supervision in the form of annotations of the ground truth causal variables or a known causal graph [111, 126, 130].

A different line of work instead approaches causal representation learning from the perspective of causal discovery in the presence of latent variables [115]. Doing so from *purely observational i.i.d. data* requires additional constraints on the generative process, such as restrictions on the graph structure or particular parametric and distributional assumptions, and typically leverages the identifiability of linear ICA [27, 38, 79]. For *linear, non-Gaussian* models, Silva et al. [113] show that the causal DAG can be recovered up to Markov equivalence if all observed variables are “pure” in that they have a unique latent causal parent. Cai et al. [23] and Xie et al. [128, 129] extend this result to identify the full graph given two pure observed children per latent, and Adams et al. [1] provide sufficient and necessary graphical conditions for full identification. For *discrete* causal variables, Kivva et al. [68] introduce a similar “no twins” condition to reduce the task of learning the number and cardinality of latents and the Markov-equivalence class of the graph to mixture identifiability.

Other lines of work have investigated the relationship between causal models at different levels of coarse-graining or abstraction [6, 11, 12, 24–26, 35, 105, 134, 135], or learning invariant representations in a supervised setting [3, 8, 32, 33, 73, 86, 92, 124, 125], often for domain generalization.

Other concurrent works address, e.g., learning from soft interventions with polynomial mixing [136], or inferring both causal graph and the number of latents subject to graphical constraints [63].

A.1 Comparison of Related Identifiability Results

To complement the presentation of related multienvironment CRL works in § 1, we provide a structured overview of and comparison with existing identifiability results for causal representation learning in Tab. 1. The table categorizes work along different dimensions. First, we make a distinction based on the type of data (observational, interventional, or counterfactual) results rely on (colour coded in green, yellow, and red, respectively). These are also referred to as different rungs in the “ladder of causation” [97] or layers in the Pearl Causal Hierarchy [10]. Second, we categorize work depending on the assumptions placed on the latent causal model and the mixing function. As can be seen, works relying solely on observational data often require restrictive graphical assumptions on the mixing function. On the other hand, access to much more informative counterfactual data has allowed identification even for nonparametric causal models and mixing functions. Our work can be viewed as a step towards addressing the lack of nonparametric identifiability results in the interventional regime.

We emphasize that Tab. 1 is not exhaustive: certain relevant works did not easily fit into our categorization or the causal representation learning framework adopted in the present work. For example, not listed are works that leverage temporal structure [2, 74, 75, 80, 81], rely on heterogenous data and distribution shifts which are not directly expressed in terms of or linked to interventions [82, 132, 133], allow for edges from observed to latent variables [1], or require more direct supervision [111, 130].

Table 1: Comparison of Existing Identifiability Results for Causal Rerepresentation Learning. All of the listed works assume invertibility (or injectivity) of the mixing function, as well as causal sufficiency (Markovianity) for the causal latent variables. Most or all of the listed results require additional technical assumptions, and may provide additional results, which we omit for sake of readability; see the references for more details.

Work	Layer	Causal Model	Mixing Function	Main Identifiability Result
Cai et al. [23], Xie et al. [128, 129]	observational	linear, non-Gaussian	linear with non-Gaussian noise s.t. each V_i has 2 pure (obs. or unobs.) children	number of latents + G
Kivva et al. [68]	observational	discrete, nonparametric	identifiable mixture model s.t. obs. children of $V_i \not\subseteq$ obs. children of V_j	number, cardinality & dist. of discrete latents + G up to Markov equivalence
Ahuja et al. [5, Thm. 4]	observational	nonlinear w. independent support [103, 125]	finite-degree polynomial	V up to permutation, shift, & linear scaling
Squires et al. [117, Thms. 1 & 2]	interventional	linear	linear	G and V up to partial-order preserving permutations from obs. dist. & all single-node <i>perfect</i> interventions
Squires et al. [117, Thm. 1]	interventional	linear	linear	G up to transitive closure from obs. dist. & all single-node <i>imperfect</i> interventions
Varici et al. [122, Thm. 16]	interventional	nonparametric	linear	G and V up to partial-order preserving permutations from obs. dist. & all single-node <i>perfect</i> interventions
Ahuja et al. [5, Thm. 2]	interventional	nonparametric	finite-degree polynomial	V up to permutation, shift, and linear scaling from all single-node <i>perfect hard</i> interventions
Buchholz et al. [22]	interventional	linear Gaussian	nonparametric	G and V up to permutation from obs. dist. & all single-node <i>perfect</i> interventions
This Work (Thm. 3.2)	interventional	nonparametric	nonparametric	for $n = 2$: G and V up to \sim_{CRL} from all single-node <i>perfect</i> interventions, subject to genericity (3.2)
This Work (Thm. 3.4)	interventional	nonparametric	nonparametric	G and V up to \sim_{CRL} from two distinct, paired single-node <i>perfect</i> interventions per node
von Kügelgen et al. [123]	counterfactual	nonparametric	nonparametric	block of non-descendants $V_{\text{nd}(\mathcal{I})}$ up to invertible function from fat-hand <i>imperfect</i> interventions on $V_{\mathcal{I}}$
Brehmer et al. [18]	counterfactual	nonparametric	nonparametric	G and V up to \sim_{CRL} from all single-node <i>perfect</i> interventions

B Proof of Minimality of the CRL Equivalence Class \sim_{CRL}

First, let us recall the main statements from § 2.2.

Definition 2.6 (\sim_{CRL} -identifiability). Let \mathcal{H} be a space of unmixing functions $h : \mathcal{X} \rightarrow \mathbb{R}^n$ and let \mathcal{G} be the space of DAGs over n vertices. Let \sim_{CRL} be the equivalence relation on $\mathcal{H} \times \mathcal{G}$ defined as

$$(h_1, G_1) \sim_{\text{CRL}} (h_2, G_2) \iff (h_2, G_2) = (\mathbf{P}_{\pi^{-1}} \circ \phi \circ h_1, \pi(G_1)) \quad (\text{B.3})$$

for some element-wise diffeomorphism $\phi(\mathbf{v}) = (\phi_1(v_1), \dots, \phi_n(v_n))$ of \mathbb{R}^n and a permutation π of $[n]$ such that $\pi : G_1 \mapsto G_2$ is a graph isomorphism and \mathbf{P}_{π} the corresponding permutation matrix.

Proposition 3.1 (Minimality of \sim_{CRL} ; informal). *Let \mathbf{Z} be any representation that is \sim_{CRL} equivalent to \mathbf{V} , with $G' = \pi(G)$ the associated DAG. Then for any intervention on $\mathcal{V}_{\mathcal{I}} \subseteq \mathcal{V}$ in G , there exists an equally sparse intervention on $\mathcal{Z}_{\pi(\mathcal{I})} \subseteq \mathcal{Z}$ in G' inducing the same observed distribution on \mathcal{X} .*

We now restate this result more formally.

Proposition B.1 (Minimality of \sim_{CRL}). *Let $(h, G') \sim_{\text{CRL}} (f^{-1}, G)$ with π denoting the graph isomorphism mapping G to G' (i.e., a permutation that preserves the partial topological order of G). Let $\mathbf{Z} = h(\mathbf{X})$ be the inferred representation with distribution $Q_{\mathbf{Z}} = h_*(P_{\mathbf{X}})$ Markov w.r.t. G' and associated density q . Let $\mathcal{I}^e \subseteq [n]$ denote a set of intervention targets, and consider an intervention that changes $p_i(v_i | \mathbf{v}_{\text{pa}(i)})$ to some intervened mechanism $\tilde{p}_i(v_i | \mathbf{v}_{\text{pa}(i)})$ for all $i \in \mathcal{I}^e$, giving rise to the interventional distributions $P_{\mathbf{V}}^e$ and $P_{\mathbf{X}}^e = f_*(P_{\mathbf{V}}^e)$. Then there exist appropriately chosen $\tilde{q}_{\pi(i)}(z_{\pi(i)} | \mathbf{z}_{\text{pa}(\pi(i), G')})$ for $i \in \mathcal{I}^e$ such that the resulting interventional distribution $Q_{\mathbf{Z}}^e$ gives rise to the same observed distributions, that is, $P_{\mathbf{X}}^e = h_*^{-1}(Q_{\mathbf{Z}}^e)$.*

Proof. Since $(h, G') \sim_{\text{CRL}} (f^{-1}, G)$, by Defn. 2.6 we have

$$\mathbf{Z} = \mathbf{P}_{\pi^{-1}} \circ \phi(\mathbf{V}) \quad (\text{B.1})$$

for some element-wise diffeomorphism ϕ with inverse $\psi = \phi^{-1}$. Then (B.1) implies that for all $i \in [n]$

$$V_i = \psi_i(Z_{\pi(i)}) \quad (\text{B.2})$$

According to (B.2), each conditional in the Markov factorization of $Q_{\mathbf{Z}}$ is given in terms of p by

$$q_{\pi(i)}(z_{\pi(i)} | \mathbf{z}_{\text{pa}(\pi(i), G')}) = p_i \left(\psi_i(z_{\pi(i)}) | \psi_{\text{pa}(i)}(\mathbf{z}_{\text{pa}(\pi(i), G')}) \right) \left| \frac{d\psi_i}{dz_{\pi(i)}}(z_{\pi(i)}) \right| \quad (\text{B.3})$$

where we have used the change of variables in (B.2) and the fact that $\pi(\text{pa}(i)) = \text{pa}(\pi(i), G')$ since $\pi : G \mapsto G'$ is a graph isomorphism.

Consider an intervention that changes $p_i(v_i | \mathbf{v}_{\text{pa}(i)})$ to some intervened mechanism $\tilde{p}_i(v_i | \mathbf{v}_{\text{pa}(i)})$ for all $i \in \mathcal{I}^e$. Denote the corresponding interventional joint distribution by $P_{\mathbf{V}}^e$ with joint density p^e given by

$$p^e(\mathbf{v}) = \prod_{i \in \mathcal{I}^e} \tilde{p}_i(v_i | \mathbf{v}_{\text{pa}(i)}) \prod_{j \in [n] \setminus \mathcal{I}^e} p_j(v_j | \mathbf{v}_{\text{pa}(j)}) \quad (\text{B.4})$$

Denote by $Q_{\mathbf{Z}}^e = (\mathbf{P}_{\pi^{-1}} \circ \phi)_*(P_{\mathbf{V}}^e)$ the corresponding distribution over \mathbf{Z} with joint density given by q^e

$$q^e(\mathbf{z}) = p^e(\psi \circ \mathbf{P}_{\pi}(\mathbf{z})) |\det \mathbf{J}_{\psi \circ \mathbf{P}_{\pi}}(\mathbf{z})| \quad (\text{B.5})$$

$$= \prod_{i \in \mathcal{I}^e} \tilde{p}_i \left(\psi_i(z_{\pi(i)}) | \psi_{\text{pa}(i)}(\mathbf{z}_{\text{pa}(\pi(i), G')}) \right) \left| \frac{d\psi_i}{dz_{\pi(i)}}(z_{\pi(i)}) \right| \prod_{j \in [n] \setminus \mathcal{I}^e} q_{\pi(j)}(z_{\pi(j)} | \mathbf{z}_{\text{pa}(\pi(j), G')}) , \quad (\text{B.6})$$

where we have used (B.3), (B.4), and the fact that \mathbf{J}_{ψ} is diagonal.

By defining

$$\tilde{q}_{\pi(i)}(z_{\pi(i)} | \mathbf{z}_{\text{pa}(\pi(i), G')}) := \tilde{p}_i \left(\psi_i(z_{\pi(i)}) | \psi_{\text{pa}(i)}(\mathbf{z}_{\text{pa}(\pi(i), G')}) \right) \left| \frac{d\psi_i}{dz_{\pi(i)}}(z_{\pi(i)}) \right| \quad (\text{B.7})$$

we finally arrive at

$$q^e(\mathbf{z}) = \prod_{i \in \mathcal{I}^e} \tilde{q}_{\pi(i)} \left(z_{\pi(i)} \mid \mathbf{z}_{\text{pa}(\pi(i); G')} \right) \prod_{j \in [n] \setminus \mathcal{I}^e} q_{\pi(j)} \left(z_{\pi(j)} \mid \mathbf{z}_{\text{pa}(\pi(j); G')} \right). \quad (\text{B.8})$$

This shows that any intervention on $\{V_i\}_{i \in \mathcal{I}^e} \subseteq \mathbf{V}$ which replaces

$$\left\{ p_i(v_i \mid \mathbf{v}_{\text{pa}(i)}) \right\}_{i \in \mathcal{I}^e} \mapsto \left\{ \tilde{p}_i(v_i \mid \mathbf{v}_{\text{pa}(i)}) \right\}_{i \in \mathcal{I}^e}, \quad (\text{B.9})$$

can equivalently be captured by an intervention on $\{Z_{\pi(i)}\}_{i \in \mathcal{I}^e} \subseteq \mathbf{Z}$ which replaces

$$\left\{ q_{\pi(i)} \left(z_{\pi(i)} \mid \mathbf{z}_{\text{pa}(\pi(i); G')} \right) \right\}_{i \in \mathcal{I}^e} \mapsto \left\{ \tilde{q}_{\pi(i)} \left(z_{\pi(i)} \mid \mathbf{z}_{\text{pa}(\pi(i); G')} \right) \right\}_{i \in \mathcal{I}^e}. \quad (\text{B.10})$$

with \tilde{q}_i defined according to (B.7). \square

C Identifiability Proofs

C.1 Auxiliary Lemmata

Lemma C.1 (Lemma 2 of Brehmer et al. [18]). *Let $A = C = \mathbb{R}$ and $B = \mathbb{R}^n$. Let $f : A \times B \rightarrow C$ be differentiable. Define two differentiable measures p_A on A and p_C on C . Let $\forall b \in B$, $f(\cdot, b) : A \rightarrow C$ be measure-preserving, in the sense that the pushforward of p_A is always p_C . Then $f(\cdot, b)$ is constant in b on B .*

Proof. See Appendix A.2 of Brehmer et al. [18]. \square

Lemma C.2 (Preservation of conditional independence under invertible transformation.). *Let X and Y be continuous real-valued random variables, and let \mathbf{Z} be a continuous random vector taking values in \mathbb{R}^n . Suppose that (X, Y, \mathbf{Z}) have a joint density w.r.t. the Lebesgue measure. Let $f : \mathbb{R} \rightarrow \mathbb{R}$, $g : \mathbb{R} \rightarrow \mathbb{R}$, and $h : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be diffeomorphisms. Then $X \perp\!\!\!\perp Y \mid \mathbf{Z} \implies f(X) \perp\!\!\!\perp g(Y) \mid h(\mathbf{Z})$.*

Proof. Denote by $p(x, y, \mathbf{z})$ the density of (X, Y, \mathbf{Z}) . Then $X \perp\!\!\!\perp Y \mid \mathbf{Z}$ implies that for all (x, y, \mathbf{z}) , p can be factorized as follows:

$$p(x, y, \mathbf{z}) = p_{\mathbf{z}}(\mathbf{z}) p_x(x \mid \mathbf{z}) p_y(y \mid \mathbf{z}). \quad (\text{C.1})$$

Let $(A, B, C) = (f(X), g(Y), h(\mathbf{Z}))$, and write $\tilde{f} = f^{-1}$, $\tilde{g} = g^{-1}$, and $\tilde{h} = h^{-1}$.

Then (A, B, C) also has a density $q(a, b, \mathbf{c})$, which for all (a, b, \mathbf{c}) is given by the change of variable formula:

$$q(a, b, \mathbf{c}) = p \left(\tilde{f}(a), \tilde{g}(b), \tilde{h}(\mathbf{c}) \right) \left| \frac{d\tilde{f}}{da}(a) \frac{d\tilde{g}}{db}(b) \det \mathbf{J}_{\tilde{h}}(\mathbf{c}) \right| \quad (\text{C.2})$$

$$= p_{\mathbf{z}} \left(\tilde{h}(\mathbf{c}) \right) \left| \det \mathbf{J}_{\tilde{h}}(\mathbf{c}) \right| p_x \left(\tilde{f}(a) \mid \tilde{h}(\mathbf{c}) \right) \left| \frac{d\tilde{f}}{da}(a) \right| p_y \left(\tilde{g}(b) \mid \tilde{h}(\mathbf{c}) \right) \left| \frac{d\tilde{g}}{db}(b) \right| \quad (\text{C.3})$$

where in (C.2) we have used the fact that $(X, Y, \mathbf{Z}) \mapsto (f(X), g(Y), h(\mathbf{Z}))$ has block-diagonal Jacobian, and in (C.3) that p factorises as in (C.1). Next, define

$$q_{\mathbf{c}}(\mathbf{c}) := p_{\mathbf{z}} \left(\tilde{h}(\mathbf{c}) \right) \left| \det \mathbf{J}_{\tilde{h}}(\mathbf{c}) \right|, \quad (\text{C.4})$$

$$q_a(a \mid \mathbf{c}) := p_x \left(\tilde{f}(a) \mid \tilde{h}(\mathbf{c}) \right) \left| \frac{d\tilde{f}}{da}(a) \right|, \quad (\text{C.5})$$

$$q_b(b \mid \mathbf{c}) := p_y \left(\tilde{g}(b) \mid \tilde{h}(\mathbf{c}) \right) \left| \frac{d\tilde{g}}{db}(b) \right|. \quad (\text{C.6})$$

Since $p_{\mathbf{z}}$, p_x , and p_y are valid densities (non-negative and integrating to one), so are $q_{\mathbf{c}}$, q_a , and q_b . Substitution into (C.3) yields for all (a, b, \mathbf{c}) ,

$$q(a, b, \mathbf{c}) = q_{\mathbf{c}}(\mathbf{c}) q_a(a \mid \mathbf{c}) q_b(b \mid \mathbf{c}), \quad (\text{C.7})$$

which implies that $A \perp\!\!\!\perp B \mid C$, concluding the proof. \square

C.2 Proof of Thm. 3.2

Theorem 3.2 (Bivariate identifiability up to \sim_{CRL} from one perfect stochastic intervention per node). *Suppose that we have access to multiple environments $\{P_{\mathbf{X}}^e\}_{e \in \mathcal{E}}$ generated as described in § 2 under Asms. 2.2, 2.5, 2.8 and 2.9 with $n = 2$. Let (h, G') be any candidate solution such that the inferred latent distributions $Q_{\mathbf{Z}}^e = h_*(P_{\mathbf{X}}^e)$ of $\mathbf{Z} = h(\mathbf{X})$ and the inferred mixing function h^{-1} satisfy the above assumptions w.r.t. the candidate causal graph G' . Assume additionally that*

- (A1) all densities p^e and q^e are continuously differentiable and fully supported on \mathbb{R}^n ;
- (A2) we have access to a known observational environment e_0 and one single node perfect intervention for each node, with unknown targets: there exist $n+1$ environments $\mathcal{E} = \{e_i\}_{i=0}^n$ such that $\mathcal{I}^{e_0} = \emptyset$ and for each $i \in [n]$ we have $\mathcal{I}^{e_i} = \{\pi(i)\}$ for an unknown permutation π of $[n]$;
- (A3) for all $i \in [n]$, the intervened mechanisms $\tilde{p}_i(v_i)$ differ from the corresponding base mechanisms $p_i(v_i | \mathbf{v}_{\text{pa}(i)})$ everywhere, in the sense that

$$\forall \mathbf{v} : \quad \frac{\partial}{\partial v_i} \frac{\tilde{p}_i(v_i)}{p_i(v_i | \mathbf{v}_{\text{pa}(i)})} \neq 0; \quad (3.1)$$

- (A4) (“**genericity**”) the base and intervened mechanisms are not fine-tuned to each other, in the sense that there exists a continuous function $\varphi : \mathbb{R}^+ \rightarrow \mathbb{R}$ for which

$$\mathbb{E}_{\mathbf{v} \sim P_{\mathbf{V}}^{e_0}} \left[\varphi \left(\frac{\tilde{p}_2(v_2)}{p_2(v_2 | v_1)} \right) \right] \neq \mathbb{E}_{\mathbf{v} \sim P_{\mathbf{V}}^{e_1}} \left[\varphi \left(\frac{\tilde{p}_2(v_2)}{p_2(v_2 | v_1)} \right) \right] \quad (3.2)$$

Then the ground truth is identified in the sense of Defn. 2.6, that is, $(f^{-1}, G) \sim_{\text{CRL}} (h, G')$.

Proof. From the assumption of a shared mixing f and shared encoder h across all environments, we have that

$$\mathbf{Z} = h(\mathbf{X}) = h(f(\mathbf{V})) = h \circ f(\mathbf{V}). \quad (C.8)$$

Let $\psi = f^{-1} \circ h^{-1} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ such that

$$\mathbf{V} = \psi(\mathbf{Z}).$$

By Asm. 2.5, f , h , and thus also $h \circ f$ are diffeomorphisms. Hence, ψ is well-defined and also diffeomorphic.

By the change of variable formula, for all e and all \mathbf{z} the density $q^e(\mathbf{z})$ is given by

$$q^e(\mathbf{z}) = p^e(\psi(\mathbf{z})) |\det \mathbf{J}_{\psi}(\mathbf{z})| \quad (C.9)$$

where $(\mathbf{J}_{\psi}(\mathbf{z}))_{ij} = \frac{\partial \psi_i}{\partial z_j}(\mathbf{z})$ denotes the Jacobian of ψ .

We now consider two separate cases, depending on whether the intervention targets in q^{e_i} for $e_i \in \{e_1, e_2\}$ match those in p^{e_i} (Case 1) or not (Case 2).

Case 1: Aligned Intervention Targets. According to Asm. 2.8 and (A2), (C.9) applied to the known observational environment e_0 and the interventional environments e_1, e_2 leads to the system of equations:

$$q_1(z_1)q_2(z_2 | z_{\text{pa}(2; G')}) = p_1(\psi_1(\mathbf{z})) p_2(\psi_2(\mathbf{z}) | \psi_{\text{pa}(2)}(\mathbf{z})) |\det \mathbf{J}_{\psi}(\mathbf{z})| \quad (C.10)$$

$$\tilde{q}_1(z_1)q_2(z_2 | z_{\text{pa}(2; G')}) = \tilde{p}_1(\psi_1(\mathbf{z})) p_2(\psi_2(\mathbf{z}) | \psi_{\text{pa}(2)}(\mathbf{z})) |\det \mathbf{J}_{\psi}(\mathbf{z})| \quad (C.11)$$

$$q_1(z_1)\tilde{q}_2(z_2) = p_1(\psi_1(\mathbf{z})) \tilde{p}_2(\psi_2(\mathbf{z})) |\det \mathbf{J}_{\psi}(\mathbf{z})| \quad (C.12)$$

where $z_{\text{pa}(2; G')}$ denotes the parents of z_2 in G' , and $\text{pa}(2)$ denotes the parents of V_2 in G .

Note that neither side of the previous equations can be zero due to the full support assumption¹² (A1) and ψ being diffeomorphic implying the determinant is non-zero.

¹²This can also be relaxed to fully supported on a Cartesian product of intervals.

Taking quotients of (C.11) and (C.10), yields

$$\frac{\tilde{q}_1}{q_1}(z_1) = \frac{\tilde{p}_1}{p_1}(\psi_1(z)). \quad (\text{C.13})$$

Next, we take the partial derivative w.r.t. z_2 on both sides and use the chain rule to obtain:

$$0 = \left(\frac{\tilde{p}_1}{p_1} \right)' (\psi_1(z)) \frac{\partial \psi_1}{\partial z_2}(z). \quad (\text{C.14})$$

Now, by assumption (A3), the first term on the RHS of (C.14) is non-zero everywhere. Hence,

$$\forall \mathbf{z} : \frac{\partial \psi_1}{\partial z_2}(z) = 0. \quad (\text{C.15})$$

It follows that ψ_1 is, in fact, a scalar function, and

$$V_1 = \psi_1(Z_1). \quad (\text{C.16})$$

Since ψ is a diffeomorphism, ψ_1 must also be diffeomorphic. Hence, by the change of variable formula applied to (C.16), the marginal density $q_1(z_1)$ is given by

$$q_1(z_1) = p_1(\psi_1(z_1)) \left| \frac{\partial \psi_1}{\partial z_1}(z_1) \right|. \quad (\text{C.17})$$

Further, since \mathbf{J}_ψ is triangular due to (C.15), its determinant is given by

$$|\det \mathbf{J}_\psi(\mathbf{z})| = \left| \frac{\partial \psi_1}{\partial z_1}(z_1) \frac{\partial \psi_2}{\partial z_2}(z_1, z_2) \right|. \quad (\text{C.18})$$

Substituting (C.17) and (C.18) into (C.12) yields (after cancellation of equal terms):

$$\tilde{q}_2(z_2) = \tilde{p}_2(\psi_2(z_1, z_2)) \left| \frac{\partial \psi_2}{\partial z_2}(z_1, z_2) \right|. \quad (\text{C.19})$$

The expression in (C.19) implies that, for all z_1 , the mapping $\psi_2(z_1, \cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is measure preserving for the differentiable \tilde{q}_2 and \tilde{p}_2 . By Lemma C.1 (Lemma 2 of Brehmer et al. [18, § A.2]), it then follows that ψ_2 must, in fact, be constant in z_1 , that is

$$\forall \mathbf{z} : \frac{\partial \psi_2}{\partial z_1}(z) = 0. \quad (\text{C.20})$$

Note that this last step is where the assumption of perfect interventions (Asm. 2.9) is leveraged: the conclusion would not hold for arbitrary imperfect interventions for which (3.8) would involve $\tilde{q}_2(z_2 | z_1)$ and $p_2(\psi_2(z_1, z_2) | \psi_1(z_1))$.

Hence, we have shown that ψ is an element-wise function:

$$\mathbf{V} = (V_1, V_2) = \psi(\mathbf{Z}) = (\psi_1(Z_1), \psi_2(Z_2)). \quad (\text{C.21})$$

Finally, since ψ is a diffeomorphism, (C.21) implies that

$$V_1 \perp\!\!\!\perp V_2 \iff Z_1 \perp\!\!\!\perp Z_2. \quad (\text{C.22})$$

It then follows from the faithfulness assumption (Asm. 2.2) that we also must have $G = G'$.

This concludes the proof of Case 1, as we have shown that $(h, G') \sim_{\text{CRL}} (f^{-1}, G)$ with $G' = \pi(G) = G$ (π being the identity permutation) and $h \circ f = \psi^{-1} =: \phi$ an element-wise diffeomorphism.

Case 2: Misaligned Intervention Targets. Writing down the system of equations similar to (C.10)–(C.12), but for the case with misaligned intervention targets across p and q yields:

$$q_1(z_1)q_2(z_2 | z_{\text{pa}(2;G')}) = p_1(\psi_1(\mathbf{z})) p_2(\psi_2(\mathbf{z}) | \psi_{\text{pa}(2)}(\mathbf{z})) |\det \mathbf{J}_\psi(\mathbf{z})| \quad (\text{C.23})$$

$$\tilde{q}_1(z_1)q_2(z_2 | z_{\text{pa}(2;G')}) = p_1(\psi_1(\mathbf{z})) \tilde{p}_2(\psi_2(\mathbf{z})) |\det \mathbf{J}_\psi(\mathbf{z})| \quad (\text{C.24})$$

$$q_1(z_1)\tilde{q}_2(z_2) = \tilde{p}_1(\psi_1(\mathbf{z})) p_2(\psi_2(\mathbf{z}) | \psi_{\text{pa}(2)}(\mathbf{z})) |\det \mathbf{J}_\psi(\mathbf{z})|. \quad (\text{C.25})$$

Taking ratios of e_1 and e_2 with e_0 yields

$$\frac{\tilde{q}_1}{q_1}(z_1) = \frac{\tilde{p}_2(\psi_2(\mathbf{z}))}{p_2(\psi_2(\mathbf{z}) \mid \psi_{\text{pa}(2)}(\mathbf{z}))} \quad (\text{C.26})$$

$$\frac{\tilde{q}_2(z_2)}{q_2(z_2 \mid z_{\text{pa}(2;G')})} = \frac{\tilde{p}_1}{p_1}(\psi_1(\mathbf{z})) . \quad (\text{C.27})$$

We separate the remainder of the proof of Case 2 into different subcases depending on G and G' : as we will see, we can use a similar reasoning as in Case 1, except for the case where both G and G' are missing no edge.

Case 2a: $V_1 \not\rightarrow V_2$ in G , that is $\text{pa}(2) = \emptyset$. Then we can proceed similarly to Case 1. First, we take the partial derivative of (C.26) w.r.t. z_2 to arrive at:

$$0 = \left(\frac{\tilde{p}_2}{p_2} \right)'(\psi_2(\mathbf{z})) \frac{\partial \psi_2}{\partial z_2}(\mathbf{z}) . \quad (\text{C.28})$$

Using (A3), this implies that ψ_2 does not depend on Z_2 , that is, $V_2 = \psi_2(Z_1)$.

Next, we again write $q(z_1)$ in terms of $p_2(\psi_2(z_1))$ using the univariate change of variable formula, substitute into (C.25), cancel the corresponding terms, and arrive at:

$$\tilde{q}_2(z_2) = \tilde{p}_1(\psi_1(z_1, z_2)) \left| \frac{\partial \psi_1}{\partial z_2}(z_1, z_2) \right| \quad (\text{C.29})$$

Lemma C.1 applied to $\psi_1(z_1, \cdot)$ which preserves \tilde{q}_2 and \tilde{p}_1 for all z_1 shows that ψ_1 is constant in Z_1 , that is

$$\mathbf{V} = (V_1, V_2) = \psi(\mathbf{Z}) = (\psi_1(Z_2), \psi_2(Z_1)) . \quad (\text{C.30})$$

Since $V_1 \perp\!\!\!\perp V_2$ by the assumption of Case 2a, it follows from the invertible element-wise reparametrisation above that $Z_1 \perp\!\!\!\perp Z_2$ and hence, by faithfulness, $Z_1 \not\rightarrow Z_2$ in G' .

Finally, note that there is no partial order on the empty graph and so $G' = \pi(G) = G$ and $\mathbf{Z} = \mathbf{P}_{\pi^{-1}} \cdot \psi^{-1}(\mathbf{V})$ where π is the nontrivial permutation of $\{1, 2\}$.

Case 2b: $V_1 \rightarrow V_2$ in G , that is $\text{pa}(2) = \{1\}$. If $G' \neq G$, that is $Z_1 \not\rightarrow Z_2$ in G' , then the same argument as in Case 2a, this time starting by taking the partial derivative of (C.27) w.r.t. z_1 , can be used to reach the same conclusion in (C.30). However, this contradicts faithfulness since $V_1 \not\perp\!\!\!\perp V_2$ in G .

Hence, we must have $G' = G$, and the following two equations must hold for all \mathbf{z} :

$$\frac{\tilde{q}_1(z_1)}{q_1(z_1)} = \frac{\tilde{p}_2(\psi_2(\mathbf{z}))}{p_2(\psi_2(\mathbf{z}) \mid \psi_1(\mathbf{z}))} \quad (\text{C.31})$$

$$\frac{\tilde{q}_2(z_2)}{q_2(z_2 \mid z_1)} = \frac{\tilde{p}_1(\psi_1(\mathbf{z}))}{p_1(\psi_1(\mathbf{z}))} \quad (\text{C.32})$$

The remainder of the proof consists of exploring the implications of (C.32) and (C.31), ultimately resulting in a violation of the genericity condition (A4).

To ease notation, define the following auxiliary functions:

$$a(z_1) := \frac{\tilde{q}_1(z_1)}{q_1(z_1)} , \quad (\text{C.33})$$

$$b(\mathbf{v}) := \frac{\tilde{p}_2(v_2)}{p_2(v_2 \mid v_1)} , \quad (\text{C.34})$$

$$c(\mathbf{z}) := \frac{\tilde{q}_2(z_2)}{q_{2|1}(z_2 \mid z_1)} , \quad (\text{C.35})$$

$$d(v_1) := \frac{\tilde{p}_1(v_1)}{p_1(v_1)} . \quad (\text{C.36})$$

With this, (C.31) and (C.32) take the following form:

$$a(z_1) = b(\psi(\mathbf{z})). \quad (\text{C.37})$$

$$c(\mathbf{z}) = d(\psi_1(\mathbf{z})), \quad (\text{C.38})$$

Next, define the following maps:

$$\kappa : \mathbf{z} \mapsto \begin{bmatrix} a(z_1) \\ c(\mathbf{z}) \end{bmatrix} \quad (\text{C.39})$$

$$\rho : \mathbf{v} \mapsto \begin{bmatrix} b(\mathbf{v}) \\ d(v_1) \end{bmatrix} \quad (\text{C.40})$$

Then, (C.37) and (C.38) together imply that

$$\kappa = \rho \circ \psi. \quad (\text{C.41})$$

Recalling that by (A1) all densities are continuously differentiable, the Jacobians of κ and ρ are given by:

$$\mathbf{J}_\kappa(\mathbf{z}) = \begin{bmatrix} a'(z_1) & 0 \\ \frac{\partial c}{\partial z_1}(\mathbf{z}) & \frac{\partial c}{\partial z_2}(\mathbf{z}) \end{bmatrix}, \quad (\text{C.42})$$

$$\mathbf{J}_\rho(\mathbf{v}) = \begin{bmatrix} \frac{\partial b}{\partial v_1}(\mathbf{v}) & \frac{\partial b}{\partial v_2}(\mathbf{v}) \\ d'(v_1) & 0 \end{bmatrix}, \quad (\text{C.43})$$

and the corresponding determinants are given by

$$|\det \mathbf{J}_\kappa(\mathbf{z})| = \left| a'(z_1) \frac{\partial c}{\partial z_2}(\mathbf{z}) \right| \neq 0 \quad (\text{C.44})$$

$$|\det \mathbf{J}_\rho(\mathbf{v})| = \left| d'(v_1) \frac{\partial b}{\partial v_2}(\mathbf{v}) \right| \neq 0 \quad (\text{C.45})$$

where the inequalities for all \mathbf{z} follow since, by assumption (A3), the derivatives of ratios of intervened and original mechanisms are non-vanishing everywhere:

$$a'(z_1) \neq 0 \neq \frac{\partial c}{\partial z_2}(\mathbf{z}) \quad \text{and} \quad d'(v_1) \neq 0 \neq \frac{\partial b}{\partial v_2}(\mathbf{v}), \quad (\text{C.46})$$

This implies that the following families of maps are continuously differentiable, monotonic, and invertible,

$$a : z_1 \mapsto a(z_1), \quad (\text{C.47})$$

$$b_{v_1} : v_2 \mapsto b(v_1, v_2), \quad (\text{C.48})$$

$$c_{z_1} : z_2 \mapsto c(z_1, z_2), \quad (\text{C.49})$$

$$d : v_1 \mapsto d(v_1), \quad (\text{C.50})$$

with continuously differentiable inverses

$$a^{-1} : w_1 \mapsto a^{-1}(w_1), \quad (\text{C.51})$$

$$b_{v_1}^{-1} : w_1 \mapsto b_{v_1}^{-1}(w_1), \quad (\text{C.52})$$

$$c_{z_1}^{-1} : w_2 \mapsto c_{z_1}^{-1}(w_2), \quad (\text{C.53})$$

$$d^{-1} : w_2 \mapsto d^{-1}(w_2). \quad (\text{C.54})$$

This implies that ρ and κ are valid diffeomorphisms onto their image and their inverses are given by:

$$\kappa^{-1} : \mathbf{w} \mapsto \begin{bmatrix} a^{-1}(w_1) \\ c_{a^{-1}(w_1)}^{-1}(w_2) \end{bmatrix}, \quad (\text{C.55})$$

$$\rho^{-1} : \mathbf{w} \mapsto \begin{bmatrix} d^{-1}(w_2) \\ b_{d^{-1}(w_2)}^{-1}(w_1) \end{bmatrix}. \quad (\text{C.56})$$

Since $\mathbf{V} = \psi(\mathbf{Z})$, by (C.41) we have

$$\mathbf{W} := \rho(\mathbf{V}) = \rho \circ \psi(\mathbf{Z}) = \kappa(\mathbf{Z}). \quad (\text{C.57})$$

Denote the distributions of \mathbf{W} by $R_{\mathbf{W}}$ and its density by $r(\mathbf{w})$. Since for all e , we have

$$P_{\mathbf{V}}^e = \psi_*(Q_{\mathbf{Z}}^e) \quad (\text{C.58})$$

it follows from (C.57) that

$$R_{\mathbf{W}}^e = \rho_*(P_{\mathbf{V}}^e) = \kappa_*(Q_{\mathbf{Z}}^e). \quad (\text{C.59})$$

This provides two different ways of applying the change of variable formula to compute $r(\mathbf{w})$.

First, we consider the pushforward of $Q_{\mathbf{Z}}^{e_0}$ by κ :

$$r(\mathbf{w}) = q\left(a^{-1}(\mathbf{w})\right) \left| \det \mathbf{J}_{\kappa^{-1}}(\mathbf{w}) \right| \quad (\text{C.60})$$

$$= q_1\left(a^{-1}(w_1)\right) q_2\left(c_{a^{-1}(w_1)}^{-1}(w_2) \mid a^{-1}(w_1)\right) \left| \frac{d}{dw_1} a^{-1}(w_1) \frac{d}{dw_2} c_{a^{-1}(w_1)}^{-1}(w_2) \right| \quad (\text{C.61})$$

By integrating this joint density with respect to w_2 , we obtain the following expression for the marginal $r_1(w_1)$:

$$r_1(w_1) = \left| \frac{d}{dw_1} a^{-1}(w_1) \right| q_1\left(a^{-1}(w_1)\right) \int q_2\left(c_{a^{-1}(w_1)}^{-1}(w_2) \mid a^{-1}(w_1)\right) \left| \frac{d}{dw_2} c_{a^{-1}(w_1)}^{-1}(w_2) \right| dw_2. \quad (\text{C.62})$$

By the diffeomorphic change of variable $z_2 = c_{a^{-1}(w_1)}^{-1}(w_2)$,¹³ this can be written as

$$r_1(w_1) = \left| \frac{d}{dw_1} a^{-1}(w_1) \right| q_1\left(a^{-1}(w_1)\right) \int q_2\left(z_2 \mid a^{-1}(w_1)\right) dz_2 \quad (\text{C.63})$$

$$= \left| \frac{d}{dw_1} a^{-1}(w_1) \right| q_1\left(a^{-1}(w_1)\right) \quad (\text{C.64})$$

Next, we carry out the same calculation for the pushforward of $P_{\mathbf{V}}^{e_0}$ by ρ :

$$r(\mathbf{w}) = p\left(\rho^{-1}(\mathbf{w})\right) \left| \det \mathbf{J}_{\rho^{-1}}(\mathbf{w}) \right| \quad (\text{C.65})$$

$$= p_1\left(d^{-1}(w_2)\right) p_2\left(b_{d^{-1}(w_2)}^{-1}(w_1) \mid d^{-1}(w_2)\right) \left| \frac{d}{dw_2} d^{-1}(w_2) \frac{d}{dw_1} b_{d^{-1}(w_2)}^{-1}(w_1) \right|, \quad (\text{C.66})$$

leading to the marginal

$$r_1(w_1) = \int p_1\left(d^{-1}(w_2)\right) p_2\left(b_{d^{-1}(w_2)}^{-1}(w_1) \mid d^{-1}(w_2)\right) \left| \frac{d}{dw_1} b_{d^{-1}(w_2)}^{-1}(w_1) \right| \left| \frac{d}{dw_2} d^{-1}(w_2) \right| dw_2 \quad (\text{C.67})$$

$$= \int p_1(v_1) p_2\left(b_{v_1}^{-1}(w_1) \mid v_1\right) \left| \frac{d}{dw_1} b_{v_1}^{-1}(w_1) \right| dv_1, \quad (\text{C.68})$$

where the second line is obtained by the diffeomorphic change of variable $v_1 = d^{-1}(w_2)$.

Equating the two expressions for $r(w_1)$ in e_0 in (C.68) and (C.64), we obtain for all w_1 :

$$\left| \frac{d}{dw_1} a^{-1}(w_1) \right| q_1\left(a^{-1}(w_1)\right) = \int p_1(v_1) p_2\left(b_{v_1}^{-1}(w_1) \mid v_1\right) \left| \frac{d}{dw_1} b_{v_1}^{-1}(w_1) \right| dv_1. \quad (\text{C.69})$$

Applying the same approach to the environment in which V_1 is intervened upon changing p_1 to \tilde{p}_1 while Z_2 is intervened upon leaving q_1 invariant, yields for all w_1 :

$$\left| \frac{d}{dw_1} a^{-1}(w_1) \right| q_1\left(a^{-1}(w_1)\right) = \int \tilde{p}_1(v_1) p_2\left(b_{v_1}^{-1}(w_1) \mid v_1\right) \left| \frac{d}{dw_1} b_{v_1}^{-1}(w_1) \right| dv_1. \quad (\text{C.70})$$

¹³Note that: $\int q_2(z_2(w_2)) \left| \frac{dz_2}{dw_2} \right| dw_2 = \int q_2(z_2) dz_2$.

Finally, by equating (C.69) and (C.70), we arrive at the following expression which must hold for all w_1 :

$$\int p_1(v_1) p_2 \left(b_{v_1}^{-1}(w_1) \mid v_1 \right) \left| \frac{d}{dw_1} b_{v_1}^{-1}(w_1) \right| dw_1 = \int \tilde{p}_1(v_1) p_2 \left(b_{v_1}^{-1}(w_1) \mid v_1 \right) \left| \frac{d}{dw_1} b_{v_1}^{-1}(w_1) \right| dw_1 \quad (\text{C.71})$$

which we can rewrite as

$$\int (\tilde{p}_1(v_1) - p_1(v_1)) p_2 \left(b_{v_1}^{-1}(w_1) \mid v_1 \right) \left| \frac{d}{dw_1} b_{v_1}^{-1}(w_1) \right| dw_1 = 0. \quad (\text{C.72})$$

Multiplying by any continuous function $\varphi(w_1)$, integrating w.r.t. w_1 and applying the diffeomorphic change of variable $v_2 = b_{v_1}^{-1}(w_1)$, this can be expressed as:

$$0 = \int \varphi(w_1) \int (\tilde{p}_1(v_1) - p_1(v_1)) p_2 \left(b_{v_1}^{-1}(w_1) \mid v_1 \right) \left| \frac{d}{dw_1} b_{v_1}^{-1}(w_1) \right| dv_1 dw_1 \quad (\text{C.73})$$

$$= \int \int \varphi(b_{v_1}(v_2)) (\tilde{p}_1(v_1) - p_1(v_1)) p_2(v_2 \mid v_1) dv_2 dv_1 \quad (\text{C.74})$$

$$= \int \int \varphi \left(\frac{\tilde{p}_2(v_2)}{p_2(v_2 \mid v_1)} \right) (\tilde{p}_1(v_1) - p_1(v_1)) p_2(v_2 \mid v_1) dv_2 dv_1 \quad (\text{C.75})$$

where we have resubstituted the expression for $b_{v_1}(v_2)$ in the last line.

Equivalently, this can be written as: for any continuous function φ ,

$$\mathbb{E}_{v \sim P_V^{\varepsilon_0}} \left[\varphi \left(\frac{\tilde{p}_2(v_2)}{p_2(v_2 \mid v_1)} \right) \right] = \mathbb{E}_{v \sim P_V^{\varepsilon_1}} \left[\varphi \left(\frac{\tilde{p}_2(v_2)}{p_2(v_2 \mid v_1)} \right) \right]. \quad (\text{C.76})$$

However, the genericity condition (A4) precisely rules this out, since the above equality must be violated for at least one φ , concluding this last case.

To sum up, all cases either lead to a contradiction, or imply the conclusion that $(f^{-1}, G) \sim_{\text{CRL}} (h, G')$, concluding the proof. \square

C.3 Proof of Thm. 3.4

Theorem 3.4 (Identifiability up to \sim_{CRL} from two paired perfect stochastic interventions per node). *Suppose that we have access to multiple environments $\{P_{\mathbf{X}}^e\}_{e \in \mathcal{E}}$ generated as described in § 2 under Asms. 2.2, 2.3, 2.5, 2.8 and 2.9. Let (h, G') be any candidate solution such that the inferred latent distributions $Q_{\mathbf{Z}}^e = h_*(P_{\mathbf{X}}^e)$ of $\mathbf{Z} = h(\mathbf{X})$ and the inferred mixing function h^{-1} satisfy the above assumptions w.r.t. the candidate causal graph G' . Assume additionally that*

- (A1) all densities p^e and q^e are continuously differentiable and fully supported on \mathbb{R}^n ;
- (A2') we have access to at least one pair of single-node perfect interventions per node, with unknown targets: there exist $m \geq n$ known pairs of environments $\mathcal{E} = \{(e_j, e'_j)\}_{j=1}^m$ such that for each $i \in [n]$ there exists some unknown $j \in [m]$ for which $\mathcal{I}^{e_j} = \mathcal{I}^{e'_j} = \{i\}$;
- (A3') for all $i \in [n]$, the intervened mechanisms $\tilde{p}_i(v_i)$ and $\tilde{\tilde{p}}_i(v_i)$ differ everywhere, in the sense that

$$\forall v_i : \left(\frac{\tilde{\tilde{p}}_i}{\tilde{p}_i} \right)'(v_i) \neq 0; \quad (3.10)$$

Then the ground truth is identified in the sense of Defn. 2.6, that is, $(f^{-1}, G) \sim_{\text{CRL}} (h, G')$.

Proof. First, we show that we can extract from the $m \geq n$ available pairs of environments a suitable subset \mathcal{E}_n of exactly n pairs, containing one pair of interventional environments for each node.

Let $\mathcal{E}_n \subseteq \mathcal{E}$ be a subset of n pairs of environments which are assumed to correspond to distinct targets in the model q , and suppose for a contradiction that this is not actually the case for the ground truth p (i.e., there are duplicate and missing interventions w.r.t. p). Then there must be two pairs of environments $(e_a, e'_a), (e_b, e'_b) \in \mathcal{E}_n$, both corresponding to interventions on some V_i in p , but which

are modelled as interventions on distinct nodes Z_j and Z_k with $j \neq k$ in q . We show that this implies that V_i must simultaneously be a deterministic function of only Z_j and only Z_k . Similar to the proof of Thm. 3.2, we obtain the following equations,

$$\frac{\tilde{q}_j}{\tilde{q}_j}(z_j) = \frac{\tilde{p}_i}{\tilde{p}_i}(\psi_i(\mathbf{z})) , \quad (\text{C.77})$$

$$\frac{\tilde{q}_k}{\tilde{q}_k}(z_k) = \frac{\hat{p}_i}{\hat{p}_i}(\psi_i(\mathbf{z})) . \quad (\text{C.78})$$

By taking partial derivatives w.r.t. z_l and applying assumption (A3'), we find that

$$\frac{\partial \psi_i}{\partial z_l} = 0 \quad \forall l \neq j , \quad (\text{C.79})$$

$$\frac{\partial \psi_i}{\partial z_l} = 0 \quad \forall l \neq k . \quad (\text{C.80})$$

Since $j \neq k$, this implies that $\partial \psi_i / \partial z_l = 0$ for all l which contradicts invertibility of ψ . Thus, by contradiction, we find that \mathcal{E}_n must contain exactly one pair of intervention per node also w.r.t. p . For the remainder of the proof, we only consider \mathcal{E}_n .

W.l.o.g., for any $(e_i, e'_i) \in \mathcal{E}_n$ we now fix the intervention targets in p to $\mathcal{I}^{e_i} = \mathcal{I}^{e'_i} = \{i\}$ and let π be a permutation of $[n]$ such that $\pi(i)$ denotes the inferred intervention target in q that by (A2') is shared across (e_i, e'_i) . (We will show later that not all permutations are admissible, but only ones that preserve the partial order of G .)

The first part of the proof is similar to Case 1 in the proof of Thm. 3.2. Consider the densities in environments e_i and e'_i , which are related through the change of variable formula by:

$$\tilde{q}_{\pi(i)}(z_{\pi(i)}) \prod_{j \in [n] \setminus \{\pi(i)\}} q_j(z_j | \mathbf{z}_{\text{pa}(j; G')}) = \tilde{p}_i(\psi_i(\mathbf{z})) \prod_{j \in [n] \setminus \{i\}} p_j(\psi_j(\mathbf{z}) | \psi_{\text{pa}(j)}(\mathbf{z})) |\det \mathbf{J}_\psi(\mathbf{z})| , \quad (\text{C.81})$$

$$\tilde{q}_{\pi(i)}(z_{\pi(i)}) \prod_{j \in [n] \setminus \{\pi(i)\}} q_j(z_j | \mathbf{z}_{\text{pa}(j; G')}) = \tilde{p}_i(\psi_i(\mathbf{z})) \prod_{j \in [n] \setminus \{i\}} p_j(\psi_j(\mathbf{z}) | \psi_{\text{pa}(j)}(\mathbf{z})) |\det \mathbf{J}_\psi(\mathbf{z})| , \quad (\text{C.82})$$

where $\mathbf{Z}_{\text{pa}(j; G')} \subseteq \mathbf{Z} \setminus \{Z_j\}$ denotes the parents of Z_j in G' .

Taking the quotient of the two equations yields

$$\frac{\tilde{q}_{\pi(i)}}{\tilde{q}_{\pi(i)}}(z_{\pi(i)}) = \frac{\tilde{p}_i}{\tilde{p}_i}(\psi_i(\mathbf{z})) . \quad (\text{C.83})$$

Next, for any $j \neq \pi(i)$, taking partial derivatives w.r.t. z_j on both sides yields

$$0 = \left(\frac{\tilde{p}_i}{\tilde{p}_i} \right)'(\psi_i(\mathbf{z})) \frac{\partial \psi_i}{\partial z_j}(\mathbf{z}) . \quad (\text{C.84})$$

By assumption (A3'), the first term on the RHS is non-zero everywhere. Hence, (C.84) implies

$$\forall j \neq \pi(i), \forall \mathbf{z} : \frac{\partial \psi_i}{\partial z_j}(\mathbf{z}) = 0 \quad (\text{C.85})$$

from which we can conclude that

$$V_i = \psi_i(Z_{\pi(i)}) \quad (\text{C.86})$$

for all $i \in [n]$. That is, ψ is the composition of the permutation π with an element-wise reparametrisation.

It remains to show that π must, in fact, be a graph isomorphism, which is equivalent to the statement

$$V_i \rightarrow V_j \text{ in } G \iff Z_{\pi(i)} \rightarrow Z_{\pi(j)} \text{ in } G' . \quad (\text{C.87})$$

(\implies) Suppose for a contradiction that there exist (i, j) such that $V_i \rightarrow V_j$ in G , but $Z_{\pi(i)} \not\rightarrow Z_{\pi(j)}$ in G' .

The main idea is to demonstrate that the lack of such direct arrow implies a certain conditional independence which, by faithfulness, would contradict the unconditional dependence of V_i and V_j .

Consider environment e_i in which there are perfect interventions on $Z_{\pi(i)}$ and V_i , which has the effect of removing all incoming arrows to $Z_{\pi(i)}$ and V_i in the respective post-intervention graphs $G'_{\overline{Z_{\pi(i)}}}$ and $G'_{\overline{V_i}}$.

As a result of this and the lack of direct arrow by assumption, any d-connecting path between $Z_{\pi(i)}$ and $Z_{\pi(j)}$ must enter the latter via $\mathbf{Z}_{\text{pa}(\pi(j); G')}$ [95].

It then follows from Markovianity of q w.r.t. G' that the following holds in $Q_{\mathbf{Z}}^{e_i}$:

$$Z_{\pi(i)} \perp\!\!\!\perp Z_{\pi(j)} \mid \mathbf{Z}_{\text{pa}(\pi(j); G')}. \quad (\text{C.88})$$

We now consider the corresponding implication for $P_{\mathbf{V}}^{e_i}$. Define

$$\tilde{\mathbf{V}} = \left\{ V_k = \psi_k \left(Z_{\pi(k)} \right) : Z_{\pi(k)} \in \mathbf{Z}_{\text{pa}(\pi(j); G')} \right\} \subseteq \mathbf{V} \setminus \{V_i, V_j\}, \quad (\text{C.89})$$

and note that by assumption, $Z_{\pi(i)} \notin \mathbf{Z}_{\text{pa}(\pi(j); G')}$ and hence $V_i \notin \tilde{\mathbf{V}}$.

By applying the corresponding diffeomorphic functions ψ_i from (C.86) to (C.88), it follows from Lemma C.2 that

$$V_i \perp\!\!\!\perp V_j \mid \tilde{\mathbf{V}} \quad (\text{C.90})$$

in $P_{\mathbf{V}}^{e_i}$. However, this violates faithfulness (Asm. 2.2) of $P_{\mathbf{V}}$ to G since V_i and V_j are d-connected in $G_{\overline{V_i}}$.

Thus, by contradiction, we must have $Z_{\pi(i)} \rightarrow Z_{\pi(j)}$ in G' .

(\impliedby) Now, suppose for a contradiction that there exist (i, j) such that $Z_{\pi(i)} \rightarrow Z_{\pi(j)}$ in G' , but $V_i \not\rightarrow V_j$ in G .

By the same argument as before, we find that

$$V_i \perp\!\!\!\perp V_j \mid \mathbf{V}_{\text{pa}(j)} \quad (\text{C.91})$$

in $P_{\mathbf{V}}^{e_i}$, and thus by Lemma C.2

$$Z_{\pi(i)} \perp\!\!\!\perp Z_{\pi(j)} \mid \tilde{\mathbf{Z}} \quad (\text{C.92})$$

in $Q_{\mathbf{Z}}^{e_i}$ where

$$\tilde{\mathbf{Z}} = \left\{ Z_{\pi(k)} : V_k \in \mathbf{V}_{\text{pa}(j)} \right\} \subseteq \mathbf{Z} \setminus \{Z_{\pi(i)}, Z_{\pi(j)}\}.$$

However, this contradicts faithfulness of $Q_{\mathbf{Z}}$ to G' . Hence, we must have that $V_i \rightarrow V_j$ in G .

This shows that π must be a graph isomorphism, thus concluding the proof. \square

C.4 Proof of Thm. 4.2

Theorem 4.2 (Preservation of causal influences under \sim_{CRL}). *Let $P_{\mathbf{V}}$ be Markovian w.r.t. G , let π be a graph isomorphism of G , and let ϕ be an element-wise diffeomorphism. Let $\mathbf{Z} = \mathbf{P}_{\pi^{-1}} \circ \phi(\mathbf{V})$ and denote its induced distribution by $Q_{\mathbf{Z}}$. Then for any $V_i \rightarrow V_j$ in G we have $\mathfrak{C}_{i \rightarrow j}^{P_{\mathbf{V}}} = \mathfrak{C}_{\pi(i) \rightarrow \pi(j)}^{Q_{\mathbf{Z}}}$.*

Proof. First, recall that according to Defn. 4.1,

$$\mathfrak{C}_{i \rightarrow j}^{P_{\mathbf{V}}} := D_{\text{KL}}(P_{\mathbf{V}} \parallel P_{\mathbf{V}}^{i \rightarrow j}), \quad (\text{C.93})$$

where $P_{\mathbf{V}}^{i \rightarrow j}$ denotes the interventional distribution obtained by replacing $p_j(v_j \mid \mathbf{v}_{\text{pa}(j)})$ with

$$p_j^{i \rightarrow j}(v_j \mid \mathbf{v}_{\text{pa}(j) \setminus \{i\}}) = \int_{\mathbf{V}_i} p_j(v_j \mid \mathbf{v}_{\text{pa}(j)}) p_i(v_i) dv_i. \quad (\text{C.94})$$

Writing out the KL divergence and noting that all terms except the interved mechanism j cancel inside the log, we obtain

$$\mathfrak{E}_{i \rightarrow j}^{P_{\mathbf{V}}} = \int_{\mathcal{V}} \log \left(\frac{p_j(v_j | \mathbf{v}_{\text{pa}(j)})}{\int_{\mathcal{V}_i} p_j(v_j | \mathbf{v}_{\text{pa}(j)}) p_i(v_i) dv_i} \right) p(\mathbf{v}) d\mathbf{v}. \quad (\text{C.95})$$

and similarly

$$\mathfrak{E}_{\pi(i) \rightarrow \pi(j)}^{Q_{\mathbf{Z}}} = \int_{\mathcal{Z}} \log \left(\frac{q_{\pi(j)}(z_{\pi(j)} | \mathbf{z}_{\text{pa}(\pi(j); G')})}{\int_{\mathcal{Z}_{\pi(i)}} q_{\pi(j)}(z_{\pi(j)} | \mathbf{z}_{\text{pa}(\pi(j); G')}) q_{\pi(i)}(z_{\pi(i)}) dz_{\pi(i)}} \right) q(\mathbf{z}) d\mathbf{z}. \quad (\text{C.96})$$

Since $\mathbf{Z} = \mathbf{P}_{\pi^{-1}} \circ \phi(\mathbf{V})$, we have $V_i = \psi_i(Z_{\pi(i)})$ for all $i \in [n]$ where $\psi = \phi^{-1}$.

Thus, by the change of variable formula, and using the fact that $\pi(\text{pa}(i)) = \text{pa}(\pi(i); G')$ since $\pi : G \mapsto G'$ is a graph isomorphism, we have for all $i \in [n]$:

$$q_{\pi(i)}(z_{\pi(i)} | \mathbf{z}_{\text{pa}(\pi(i); G')}) = p_i(\psi_i(z_{\pi(i)}) | \psi_{\text{pa}(i)}(\mathbf{z}_{\text{pa}(\pi(i); G')})) \left| \frac{d\psi_i}{dz_{\pi(i)}}(z_{\pi(i)}) \right|, \quad (\text{C.97})$$

as well as for the marginal density

$$q_{\pi(i)}(z_{\pi(i)}) = p_i(\psi_i(z_{\pi(i)})) \left| \frac{d\psi_i}{dz_{\pi(i)}}(z_{\pi(i)}) \right|, \quad (\text{C.98})$$

and

$$q(\mathbf{z}) = p(\psi \circ \mathbf{P}_{\pi}(\mathbf{z})) |\det \mathbf{J}_{\psi}(\mathbf{z})|. \quad (\text{C.99})$$

Substitution into the expression for $\mathfrak{E}_{\pi(i) \rightarrow \pi(j)}^{Q_{\mathbf{Z}}}$ yields:

$$\mathfrak{E}_{\pi(i) \rightarrow \pi(j)}^{Q_{\mathbf{Z}}} = \int_{\mathcal{Z}} \log \left(\frac{p_j(\psi_j(z_{\pi(j)}) | \psi_{\text{pa}(j)}(\mathbf{z}_{\text{pa}(\pi(j); G')}))}{\int_{\mathcal{Z}_{\pi(i)}} p_j(\psi_j(z_{\pi(j)}) | \psi_{\text{pa}(j)}(\mathbf{z}_{\text{pa}(\pi(j); G')})) p_i(\psi_i(z_{\pi(i)})) \left| \frac{d\psi_i}{dz_{\pi(i)}}(z_{\pi(i)}) \right| dz_{\pi(i)}} \right) \quad (\text{C.100})$$

$$p(\psi \circ \mathbf{P}_{\pi}(\mathbf{z})) |\det \mathbf{J}_{\psi}(\mathbf{z})| d\mathbf{z}. \quad (\text{C.101})$$

$$= \int_{\mathcal{V}} \log \left(\frac{p_j(v_j | \mathbf{v}_{\text{pa}(j)})}{\int_{\mathcal{V}_i} p_j(v_j | \mathbf{v}_{\text{pa}(j)}) p_i(v_i) dv_i} \right) p(\mathbf{v}) d\mathbf{v} \quad (\text{C.102})$$

$$= \mathfrak{E}_{i \rightarrow j}^{P_{\mathbf{V}}}. \quad (\text{C.103})$$

where the second to last line follows by integration by substitution, applied to both integrals. \square

D Experimental Details and Additional Results

In this appendix, we describe the experiments presented in § 6 in more details (Appx. D.1), and present additional results (Appx. D.2).

D.1 Experimental Details for § 6

Synthetic Data Generating Process. We consider linear Gaussian latent SCMs of the form

$$V_1 := U_1, \quad V_2 := \alpha V_1 + U_2, \quad (\text{D.1})$$

with standard normal U_1 and U_2 . As a mixing function, we use a three-layer multilayer perceptron (MLP),

$$f = \sigma \circ \mathbf{W}_3 \circ \sigma \circ \mathbf{W}_2 \circ \sigma \circ \mathbf{W}_1 \quad (\text{D.2})$$

where $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3 \in \mathbb{R}^{2 \times 2}$ are invertible weight matrices, and σ is an element-wise invertible nonlinear leaky-tanh activation function used in [41]:

$$\sigma(x) = \tanh(x) + 0.1x. \quad (\text{D.3})$$

To compute averages of our results over multiple runs, we construct different ground truth data generating processes as follows. We generate different latent SCMs by drawing α uniformly from $[-10, -2] \cup [2, 10]$. (We exclude $(-2, 2)$ to avoid sampling near unfaithful models.) We generate the corresponding mixing functions by uniformly sampling each element of the weight matrices, $(\mathbf{W}_k)_{ij} \sim U(0, 1)$. (To avoid the sampled weight matrices being too close to singular, we reject and resample if $|\det \mathbf{W}_k| < 0.1$.)

Interventional Environments. In line with Thm. 3.2, for each choice of latent SCM and mixing function, we generate three environments: one observational environment and one interventional environment for each perfect single-node intervention. For $i = 1, 2$, we model a perfect intervention on V_i by removing the influence of the parent variables and changing the exogenous noise by shifting its mean up or down. Specifically, we replace the corresponding assignment in (D.1) by

$$V_i := \tilde{U}_i, \quad \text{where } \tilde{U}_i \sim \mathcal{N}(m_i, 1) \quad (\text{D.4})$$

where the mean m_i of the shifted Gaussian noise is fixed per environment and sampled uniformly from $\{\pm 2\}$.

We label the observational environment as $e = 0$ and the environment arising from intervention on V_i by $e = i$ for $i = 1, 2$. Samples from p^e are then generated by sampling latents \mathbf{v} from the respective (un)intervened SCM and then applying the mixing function.

Model Architecture. We use normalizing flows [93] to model observations \mathbf{x} as the result of an invertible, differentiable transformation g of some latent (noise) variable \mathbf{z} ,

$$\mathbf{x} = g(\mathbf{z}). \quad (\text{D.5})$$

We apply a series of L such transformations $g^l : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ such that $g = g^L \circ \dots \circ g^1$ which we refer to as *flow layers*. We use Neural Spline Flows [30] for the invertible transformation, with a 3-layer feedforward neural network with hidden dimension 128 and permutation in each flow layer and $L = 12$ layers. The transformations g, g^1, \dots, g^L have learnable parameters (the weights and biases of the neural networks), which we omit to simplify notation.

Typically, simple distributions such as a uniform or isotropic Gaussian are used as base distribution $q(\mathbf{z})$ in normalizing flows. Here, we instead choose a base distribution that encodes information about the latent SCM. Specifically, we model the base mechanism as

$$q_1(z_1) = \mathcal{N}(\mu_1, \sigma_1^2), \quad q_2(z_2 | z_1) = \mathcal{N}(\hat{\alpha}z_1, \sigma_2^2), \quad q_2(z_2) = \mathcal{N}(\mu_2, \hat{\sigma}_2^2) \quad (\text{D.6})$$

and the intervened mechanism as

$$\tilde{q}_1(z_1) = \mathcal{N}(\tilde{\mu}_1, \tilde{\sigma}_1^2), \quad \tilde{q}_2(z_2) = \mathcal{N}(\tilde{\mu}_2, \tilde{\sigma}_2^2). \quad (\text{D.7})$$

Candidate Graphs and Intervention Targets. We train a separate normalizing-flow based model for each choice of candidate graph G' and inferred intervention targets. For the bivariate case with $n = 2$, this gives rise to four models, depending on whether G' matches G or not, and whether the intervention targets are aligned or misaligned w.r.t. the ground truth intervention targets. To model the setting $G' \neq G$ in which Z_1 and Z_2 are assumed independent, we use $q_2(z_2)$ in place of $q_2(z_2 | z_1)$ in (D.6). If the intervention targets are aligned, we use \tilde{q}_i instead of q_i in $e = i$ for $i = 1, 2$. Else, if they are misaligned, we use \tilde{q}_2 instead of q_2 in $e = 1$ and \tilde{q}_1 instead of q_1 in $e = 2$. By multiplying the respective mechanisms, we thus obtain three environment-specific joint base distributions $q^e(\mathbf{z})$ for $e = 0, 1, 2$.

Learning Objective. Given multi-environment data, the parameters $\mu_1, \sigma_1, \hat{\alpha}, \sigma_2, \mu_2, \hat{\sigma}_2, \tilde{\mu}_2, \tilde{\sigma}_2, \tilde{\mu}_1$ and $\tilde{\sigma}_1$ are jointly learned with the parameters of the invertible transformations g^l by maximising the log-likelihood of the data under our model, which is given by:

$$\sum_{e \in \mathcal{E}} \mathbb{E}_{\mathbf{x} \sim p^e(\mathbf{x})} [\log p_{\text{model}}^e(\mathbf{x})] = \sum_{e \in \mathcal{E}} \mathbb{E}_{\mathbf{x} \sim p^e(\mathbf{x})} [\log q^e(h(\mathbf{x})) + \log |\det \mathbf{J}_h(\mathbf{x})|] \quad (\text{D.8})$$

where the encoder $h := g^{-1}$ is the inverse of the normalizing flow which is readily available by construction; and where the expectations are empirical averages over the respective datasets in practice.

Training and Model Selection Details. Each environment comprises a total of 200k data points. We use the ADAM optimizer [67] with cosine annealing learning rate scheduling, starting with a learning rate of 5×10^{-3} and ending with 1×10^{-7} . We train the model for 200 epochs with a batch size of 4096. We split the dataset into 70% for training, and 15% for validation and held-out test data, each sampled randomly across all environments. For each drawn data generating process, we train three versions of each model with different random initializations and select the one with the highest validation log likelihood at the end of training for evaluation.

Evaluation Metrics. We evaluate the trained models w.r.t. *mean correlation coefficient* (MCC) on held-out data and *log-likelihood* on validation data (for model selection).

- The MCC measures the extent to which there is a one-to-one correspondence between the ground truth latents V_i and (a permuted version of) the inferred latents $Z_i = h_i(\mathbf{X})$. Its maximum value of one indicates a perfect correlation between the two. MCC is thus a proxy measure for the level of identifiability up to permutation and invertible reparametrisation. We report MCC based on Pearson (linear) correlation, though we found the results based on Spearman (nonlinear monotonic) correlation to be almost identical.
- The log-likelihood, on the other hand, measures how well a model explains or fits the data. Since the ground truth is typically unknown, a reasonable procedure when training multiple models is to select the one that attains the highest likelihood. For this reason, we report the difference in log-likelihood between misspecified models (ones assuming a wrong graph or intervention targets) to the correctly specified model. Whenever this difference is larger than zero, the correct model fits the data better and would thus be selected.

D.2 Additional Results: Learning Nonlinear Latent SCMs from Partial Causal Order

In this subsection, we present an additional experiment, in which we extend the setting investigated in § 6 and Appx. D.1 along the following axes.

- We fit generative models over three instead of two variables, corresponding to the setting of Thm. 3.4.
- The ground-truth SCM is now given by nonlinear mechanisms with non-additive, non-Gaussian noise.
- The generative model, including the learnt mechanisms, is now fully nonlinear.
- Despite Thm. 3.4 formally requiring two environments per single-node intervention, we only provide one interventional environment per node.
- Rather than searching over candidate graphs, we only fix the causal order and fit the reduced form of the SCM (see § 2.1) with a second normalizing flow.

Below, we describe these differences in more detail.

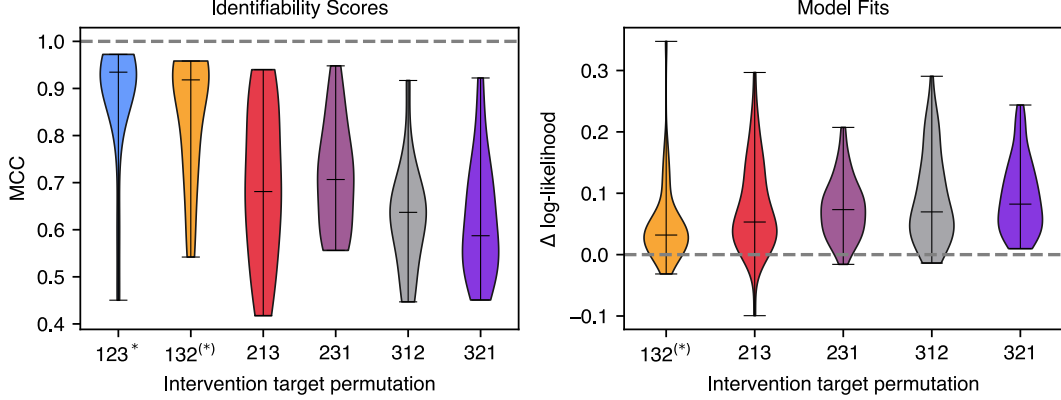


Figure 4: **Comparison of Correctly and Incorrectly Specified Models for $V_2 \leftarrow V_1 \rightarrow V_3$ with Fixed Causal Order and Nonlinear SCM.** Each violinplot corresponds to one setting where the intervention target labels are permuted. The blue plot (123^{*}) is the setting with correct intervention target labels. The yellow plot (132^{*}) has the targets for the two children V_2 and V_3 permuted, which also corresponds to a correct causal ordering and should thus be considered equivalent. We show mean correlation coefficients (MCCs) between the learned and ground truth latents (*Left*) and the difference in validation model log-likelihood between the well-specified (blue) and misspecified models (*Right*). Each violin plot is based on 20 different ground truth data generating processes; the horizontal lines indicate the minimum, median and maximum values.

Three-Variable Graph. The *unknown* ground truth graph is given by

$$V_2 \leftarrow V_1 \rightarrow V_3. \quad (\text{D.9})$$

This is consistent with the partial ordering $V_1 \preceq V_2 \preceq V_3$, which is assumed for all models a priori w.l.o.g., see § 2.2. Note that, due to the encoding of causal structure in the nonparametric model explained below, we only iterate over different permutations of the intervention targets and not over latent graph configurations. Due to the causal order implied by the graph (D.9), the permutations (1, 2, 3) (no permutation) and (1, 3, 2) (permutation of the two effects) are equivalent since the latter also implies the correct causal ordering.

Nonlinear, Non-Gaussian SCM. The mechanisms in the ground-truth SCM are now given by

$$V_i := \beta f_i^{\text{loc}}(\mathbf{V}_{\text{pa}(i)}) + f_i^{\text{scale}}(\mathbf{V}_{\text{pa}(i)})U_i \quad (\text{D.10})$$

for all i , where the location and scale functions $f_i^{\text{loc}}, f_i^{\text{scale}} : \mathbb{R}^{|\text{pa}(i)|} \rightarrow \mathbb{R}$ are parameterized by random 3-layer neural networks (sampled as the random mixing function in (D.2)) and the noise variables are Gaussian, $U_i \sim \mathcal{N}(0, 1)$. The factor β controls the influence of the parent variables relative to the exogenous noise. As β increases, variables tend to become more dependent, as also the mean shifts as a function of the parent variables. We set $\beta = 10$ for the experiments shown in Fig. 4.

Nonparametric Latent SCM. We use a second normalizing flow to learn a *reduced form of the latent SCM* via the transformation $g^{\text{SCM}} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ mapping an exogenous noise variable ϵ to the latent variable \mathbf{z} ,

$$\mathbf{z} = g^{\text{SCM}}(\epsilon). \quad (\text{D.11})$$

The distribution of the exogenous noise variable ϵ as well as the distribution of the intervened mechanisms $\tilde{q}_i(z_i)$ for $i = 1, 2, 3$ is fixed and standard (isotropic) Gaussian. The flow layers in g^{SCM} have an *upper triangular Jacobian* and thus allow us to encode assumptions about the causal graph: by passing the variables in topological order, which we can assume w.l.o.g., we ensure that an exogenous noise variable ϵ_i can only influence endogenous variables in \mathbf{z} that are descendants of z_i . The learned weights of the flow layers then implicitly encode which endogenous variables are connected. Therefore, only different choices of the permutations of the intervention targets need to be considered as candidate models. We use a similar architecture based on Neural Spline Flows. However, we omit permutation layers, which would violate the topological order of the variables.

Results. In Fig. 4, we present identifiability scores and model fits for both well-specified and misspecified models (corresponding to different intervention target choices). Notably, we observe that the well-specified model (in blue) or its equivalent (in yellow) yield the highest log-likelihood in the majority of cases, as depicted in Fig. 4 (Right). This demonstrates that, even in this nonparametric setting without fully specified graph, the log-likelihood remains a reliable criterion for selecting the correct intervention targets. Fig. 4 (Left) shows that the selected models (blue or yellow) approximately identify the ground-truth latent variables up to element-wise rescaling, whereas other choices lead to much lower MCCs.

It is worth noting that, compared to the parametric setting investigated in § 6 and Fig. 3, the nonparametric setting appears to be more challenging (as expected), as there is a less pronounced distinction between well-specified and misspecified models, both in terms of identifiability scores and model fits. Moreover, future work is needed to parse the implicitly learned causal relationships in the transformation g^{SCM} in (D.11): since only the (pre-imposed) causal order is specified, in practice, g^{SCM} may learn to use additional or fewer edges than in the true graph G .

E Discussion of the Role of Our Assumptions

Below, we summarize the rationale and intuition behind each assumption:

- Asm. 2.2 helps rule out degenerate cases (cancellation along different paths) in which variables are (conditionally) independent despite being causally related. It is a standard assumption in classical causal discovery from observational data, and therefore also helps in CRL to recover the true causal graph.
- Asm. 2.3 is required to know how many latent variables we are looking for. It is a standard assumption in identifiable representation learning (that is often made implicitly). However, it may be dropped when suitable techniques for estimating the intrinsic dimensionality of \mathcal{X} can be employed.
- Asm. 2.5 is needed for the mapping between latents and observations to be invertible in the first place. Without it, full recovery of the causal variables (up to CRL equivalence) is infeasible. This assumption is also standard for the simpler problem of nonlinear ICA.
- Asm. 2.8 is a characterisation of our generative setup. Sharing of some mechanisms and the mixing function is needed for the multi-environment setting to provide useful additional information: if everything may change across environments, the datasets can only be analysed in isolation, running into the non-identifiability of CRL from iid data.
- Asm. 2.9 and (A2) / (A2') are needed since with imperfect interventions or interventions not on all nodes, identifiability is not achievable even in the linear setting as shown by Squires et al. [117].
- Asm. (A1) is a technical assumption needed for our analysis. It is not strictly necessary (it can also be relaxed to fully supported on a Cartesian product of intervals) but substantially eases the readability and accessibility of the proof, without a major impact on the main causal aspects of the problem setup.
- Asm. (A3) / (A3') is needed to avoid spurious solutions based on applying a measure preserving transformation on a part of the domain unaffected by the intervention.
- Asm. (A4) is needed to rule out a fine-tuning of the ground-truth generating process that are possible due to fully non-parametric nature of the setup, see also Remark 4.2 and the following paragraph.